

Personalized medicine in cancer

Genome Variation | Data Formats | Resources | Sharing | Privacy

Michael Baudis **UZH SIB**
Computational Oncogenomics



University of
Zurich^{UZH}

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

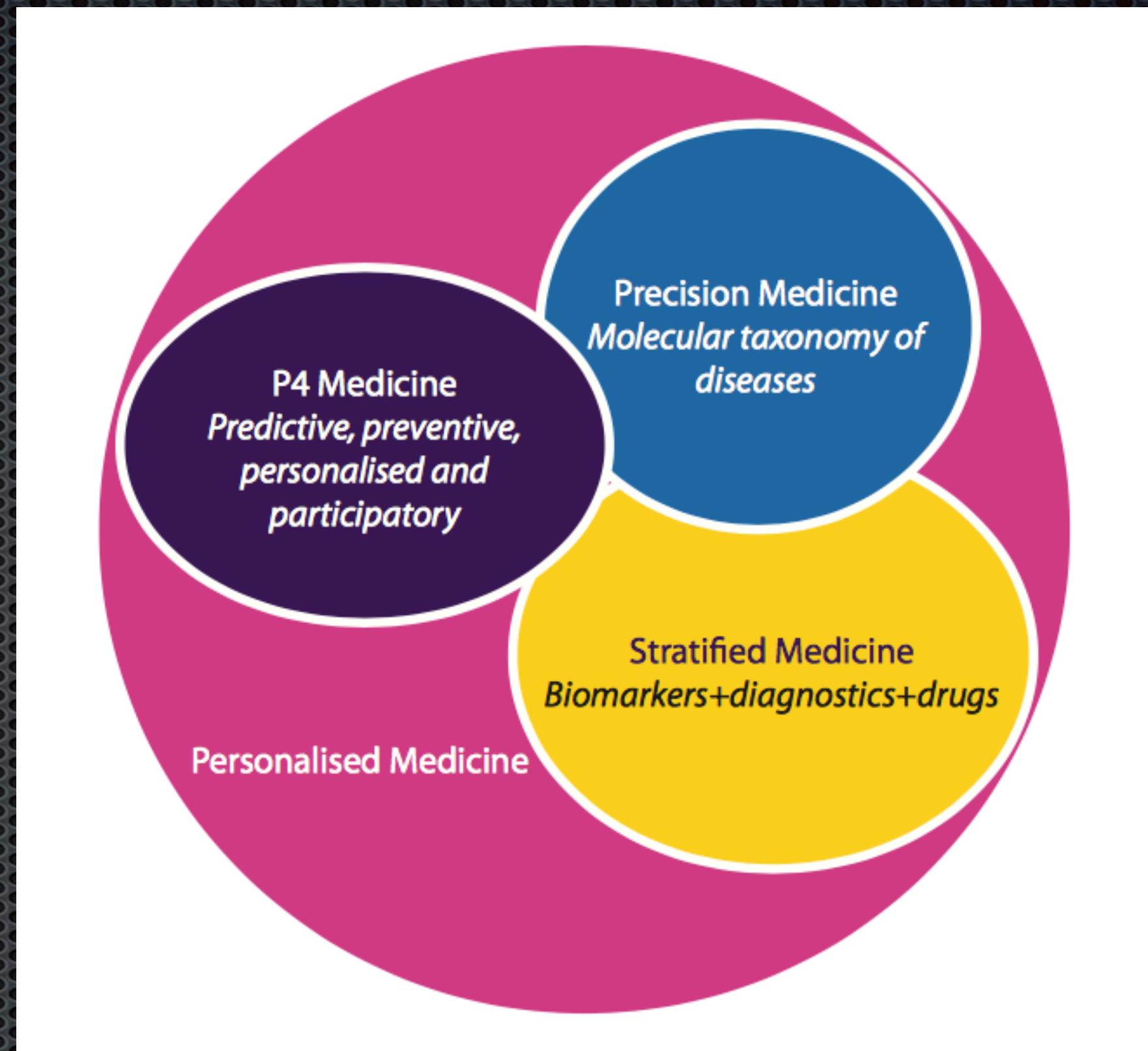
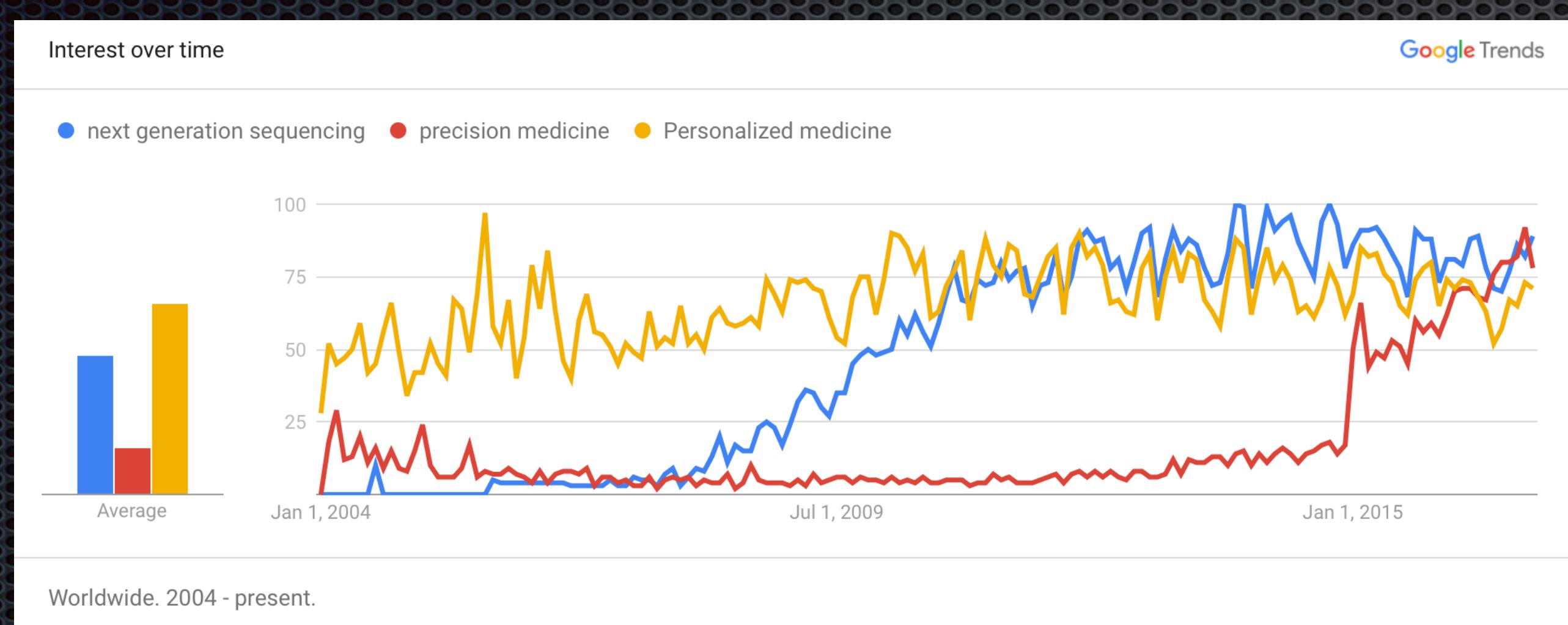
Assistant professor in paediatric haematology | molecular mechanisms of leukemogenesis | focus on bioinformatics for cancer genome data analysis

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

Many names for one concept or many concepts in one name?

Stratified, personalised, precision, individualised, P4 medicine or personalised healthcare – all are terms in use to describe notions often referred to as the future of medicine and healthcare. But what exactly is it all about, and are we all talking about the same thing?

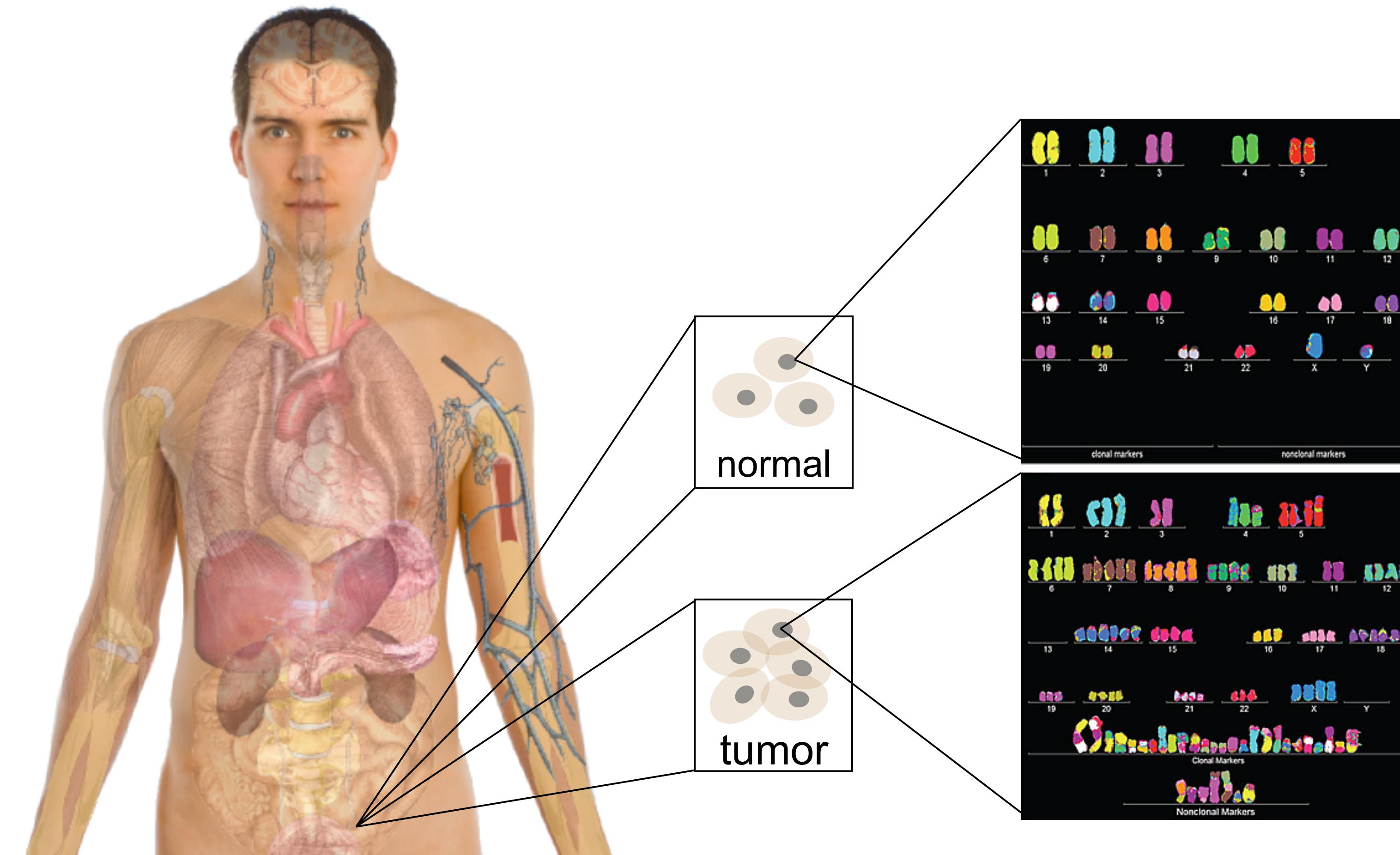


Source: PHG Foundation

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, concept based metadata and individually targeted therapies.

Personal Genomics as a Gateway into Biology

Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.



Genome analyses at the core of Personalized Health™

- Genome analyses (including transcriptome, metagenomics) are the **core technologies** for Personalized Health™ applications
- In the context of **academic medicine**, this requires
 - standard sample acquisition procedures & central **biobanking**
 - **core sequencing facility** (large throughput, cost efficiency, uniform sample and data handling procedures)
- secure **computing/analysis** platform
- Standardized **data formats** and **sample identification** procedures
- Metadata rich, reference **variant resource(s)** & expertise
- participation in reciprocal, international **data sharing** and **biocuration** efforts

Genome analyses at the core of Personalized Health™

Susceptibility, Pharmacogenomics, Classification, Infectious Diseases, Outcome Prediction, Lifestyle ...

doi:10.1038/nature19057

Analysis of protein-coding genetic variation in 60,706 humans

Monkol Lek^{1,2,3,4}, Konrad J. Karczewski^{1,2,*}, Eric V. Minikel^{1,2,5,*}, Kaitlin E. Samocha^{1,2,5,6*}, Eric Banks², Timothy Fennell², Anne H. O'Donnell-Luria^{1,2,7}, James S. Ware^{2,8,9,10,11}, Andrew J. Hill^{1,2,12}, Beryl B. Cummings^{1,2,5}, Taru Tukiainen^{1,2}, Daniel P. Birnbaum², Jack A. Kosmicki^{1,2,6,13}, Laramie E. Duncan^{1,2,6}, Karol Estrada^{1,2}, Fengmei Zhao^{1,2}, James Zou²,

Rapid whole genome sequencing and precision neonatology

Joshua E. Petrikis, MD^{a,*}, Laurel K. Willig, MD, FAAP^b, Laurie D. Smith, MD, PhD^c, and Stephen F. Kingsmore, MB, BAO, ChB, Dsc, FRCPPath^{d,e}

Barkur S. Shastry

SNP alleles in human disease and evolution



insight progress

Cancer genetics

Bruce A. J. Ponder

DISEASE MECHANISMS

Mechanisms underlying structural variant formation in genomic disorders

Claudia M. B. Carvalho^{1,2} and James R. Lupski^{1,3,4,5}

Abstract | With the recent burst of technological developments in genomics, and the clinical implementation of genome-wide assays, our understanding of the molecular basis of genomic disorders, specifically the contribution of structural variation to disease burden, is evolving

Genomic Classification of Cutaneous Melanoma

The Cancer Genome Atlas Network^{1,*,**}

¹Cancer Genome Atlas Program Office, National Cancer Institute at NIH, 31 Center Drive, Bldg. 31, Suite 3A20, Bethesda, MD 20892, USA

*Correspondence: irwatson@mdanderson.org (I.R.W.), jgershen@mdanderson.org (J.E.G.), lchin@mdanderson.org (L.C.)

<http://dx.doi.org/10.1016/j.cell.2015.05.044>



PCN Frontier Review

doi:10.1111/pcn.12128

Copy-number variation in the pathogenesis of autism spectrum disorder

Emiko Shishido, PhD^{1,2,3}, Branko Aleksić, MD, PhD³ and Norio Ozaki, MD, PhD^{3,*}

Open Access

Published online by the Promotion of Science, Japan

RESEARCH ARTICLE

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

Haoyang Cai^{1,2}, Nitin Kumar^{1,2}, Homayoun C Bagheri³, Christian von Mering^{1,2}, Mark D Robinson^{1,2*} and Michael Baudis^{1,2*}

Common gene variants, mortality and extreme longevity in humans

B.T. Heijmans^{a,b}, R.G.J. Westendorp^b, P.E. Slagboom^{a,*}

RESEARCH ARTICLE

Integrative genome-wide expression profiling identifies three distinct molecular subgroups of renal cell carcinoma with different patient outcome

Fred Beleut^{1,5*}, Philip Zimmermann², Michael Baudis³, Nicole Bruni⁴, Peter Bühlmann⁴, Oliver Laule², H-Duc Luu¹, Wilhelm Gruissem², Peter Schraml^{1,*} and Holger Moch¹

NEURODEVELOPMENT

Genes, circuits, and precision therapies for autism and related neurodevelopmental disorders

Mustafa Sahin* and Mrieganka Sur*

Activating Mutations in the Epidermal Growth Factor Receptor Underlying Responsiveness of Non-Small-Cell Lung Cancer to Gefitinib

Thomas J. Lynch, M.D., Daphne W. Bell, Ph.D., Raffaella Sordella, Ph.D., Sarada Gurubhagavatula, M.D., Ross A. Okimoto, B.S., Brian W. Brannigan, B.A., Patricia L. Harris, M.S., Sara M. Haserlat, B.A., Jeffrey G. Supko, Ph.D., Frank G. Haluska, M.D., Ph.D., David N. Louis, M.D., David C. Christiani, M.D., Jeff Settleman, Ph.D., and Daniel A. Haber, M.D., Ph.D.

N Engl J Med 2004; 350:2129-2139 | May 20, 2004 | DOI: 10.1056/NEJMoa040938

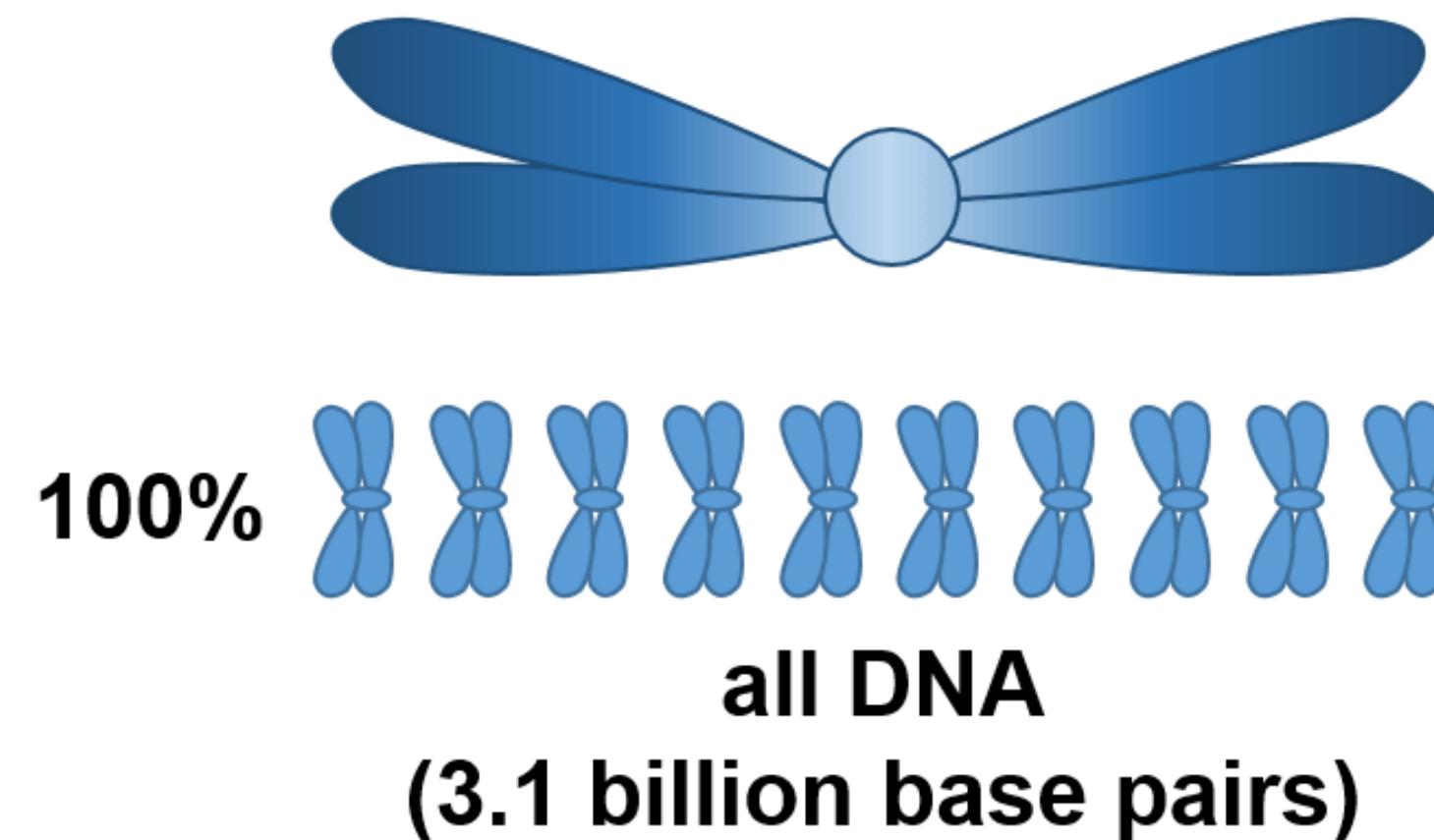
Open Access

The landscape of somatic copy-number alteration across human cancers

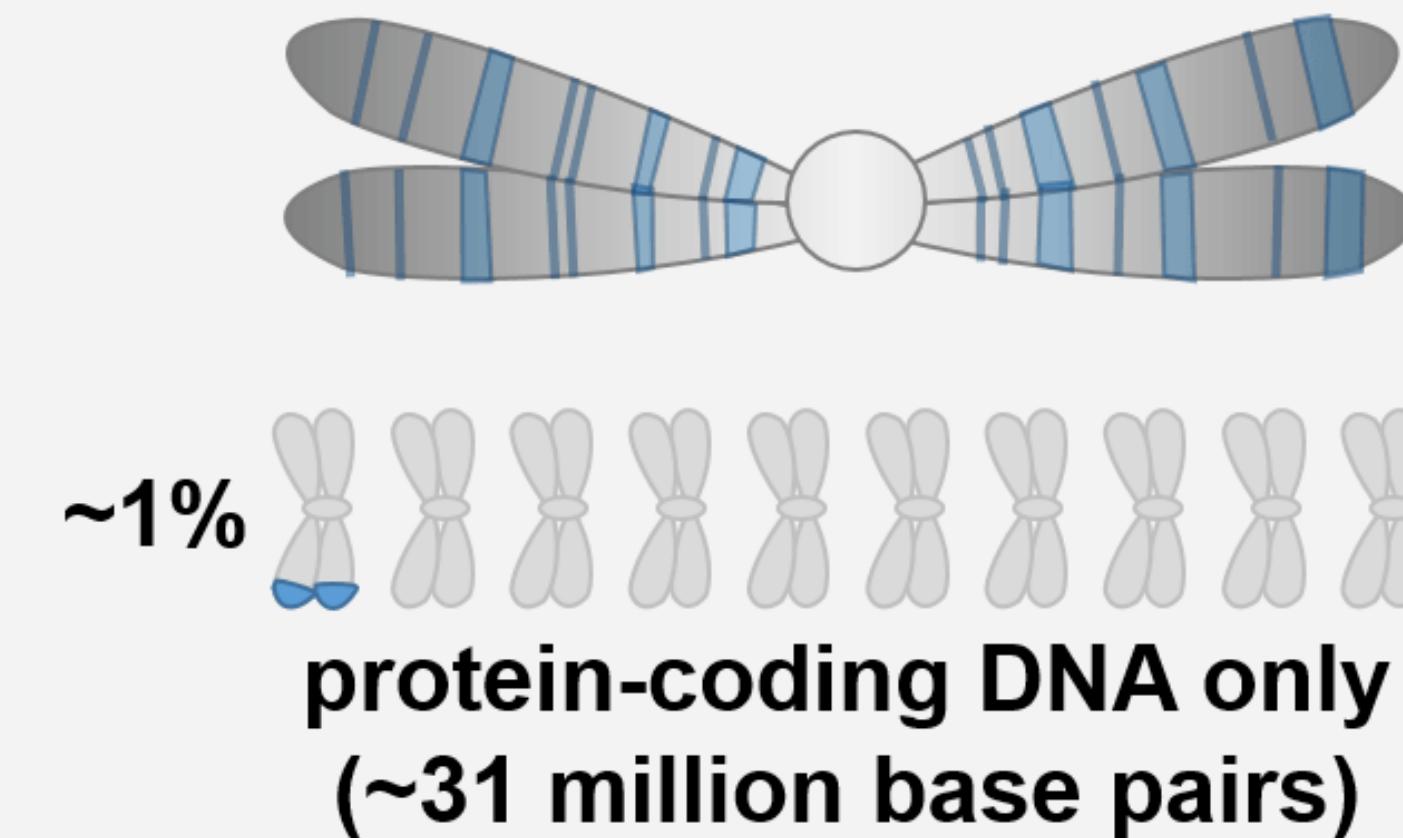
Rameen Beroukhim^{1,3,4,5,*}, Craig H. Mermel^{1,3,*}, Dale Porter⁸, Guo Wei¹, Soumya Raychaudhuri^{1,4}, Jerry Donovan⁸, Jordi Barretina^{1,3}, Jesse S. Boehm¹, Jennifer Dobson^{1,3}, Mitsuyoshi Urashima⁹, Kevin T. Mc Henry⁸, Reid M. Pinchback¹, Azra H. Ligon⁴, Yoon-Jae Cho⁶, Leila Haery^{1,3}, Heidi Greulich^{1,3,4,5}, Michael Reich¹, Wendy Winckler¹, Michael S. Lawrence¹, Barbara A. Weir^{1,3}, Kumiko E. Tanaka^{1,3}, Derek Y. Chiang^{1,3,13}, Adam J. Bass^{1,3,4}, Alice Loo⁸, Carter Hoffman^{1,3}, John Prentser^{1,3}, Ted Liefeld¹, Qing Gao¹, Derek Yecies³, Sabina Signoretti^{3,4}, Elizabeth Maher¹⁰, Frederic J. Kaye¹¹, Hidefumi Sasaki¹², Joel E. Tepper¹³, Jonathan A. Fletcher⁴, Josep Tabernero¹⁴, José Baselga¹⁴, Ming-Sound Tsao¹⁵, Francesca Demichelis¹⁶, Mark A. Rubin¹⁶, Pasi A. Janne^{3,4}, Mark J. Daly^{1,17}, Carmelo Nucera⁷, Ross L. Levine¹⁸, Benjamin L. Ebert^{1,4,5}, Stacey Gabriel¹, Anil K. Rustgi¹⁹, Cristina R. Antonescu¹⁸, Marc Ladanyi¹⁸, Anthony Letai³, Levi A. Garraway^{1,3}, Massimo Loda^{3,4}, David G. Beer²⁰, Lawrence D. True²¹, Aikou Okamoto²², Scott L. Pomeroy⁶, Samuel Singer¹⁸, Todd R. Golub^{1,3,23}, Eric S. Lander^{1,2,5}, Gad Getz¹, William R. Sellers⁸ & Matthew Meyerson^{1,3,5}

Genome Sequencing

whole genome sequencing (WGS)



exome sequencing



What does it cost to sequence a genome?

Human Genome

Project (HGP):

1991-2003

today:

2017

cost: \$2.7 billion

time: 12+ years

~\$1,500

< 2 days

today:

2017

~\$530

~3 days

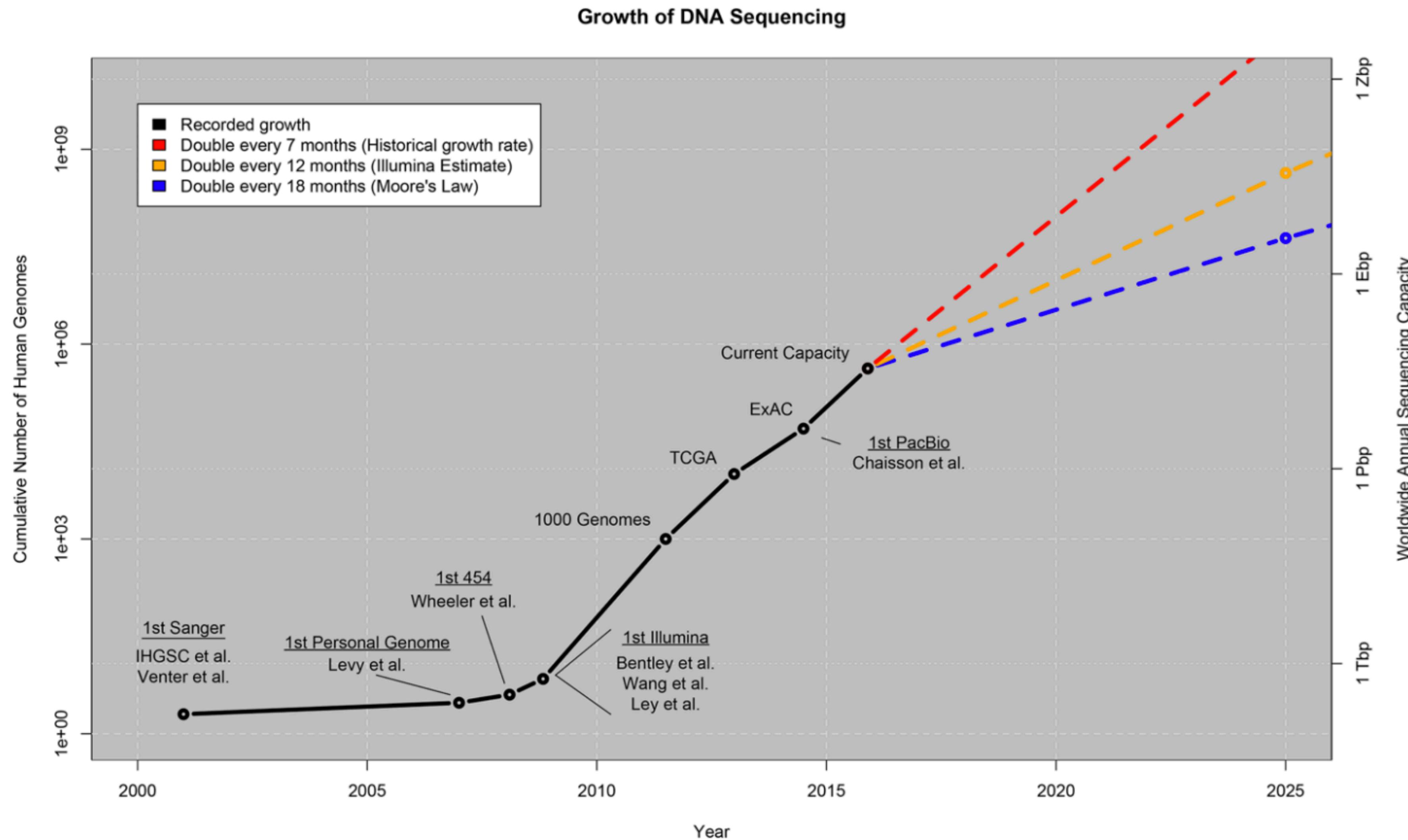
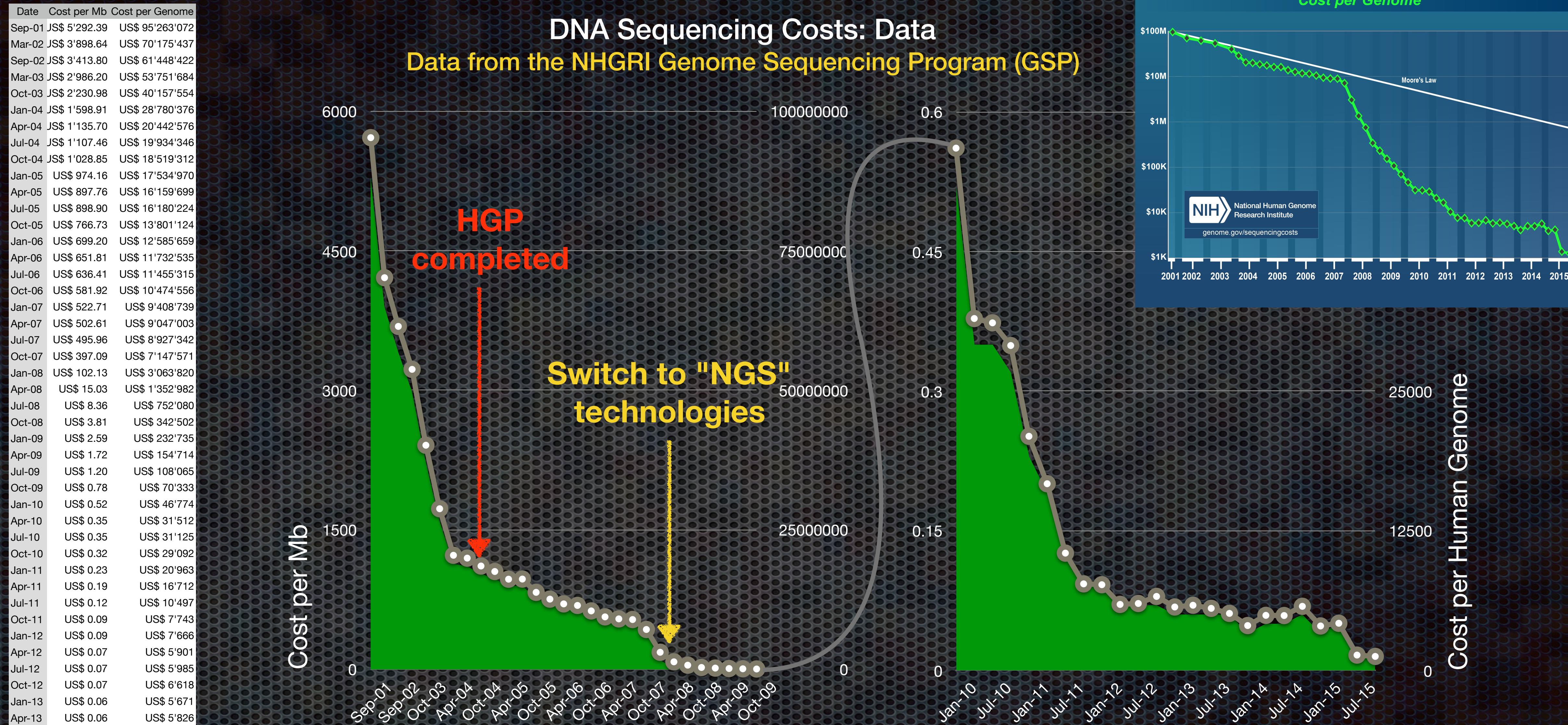
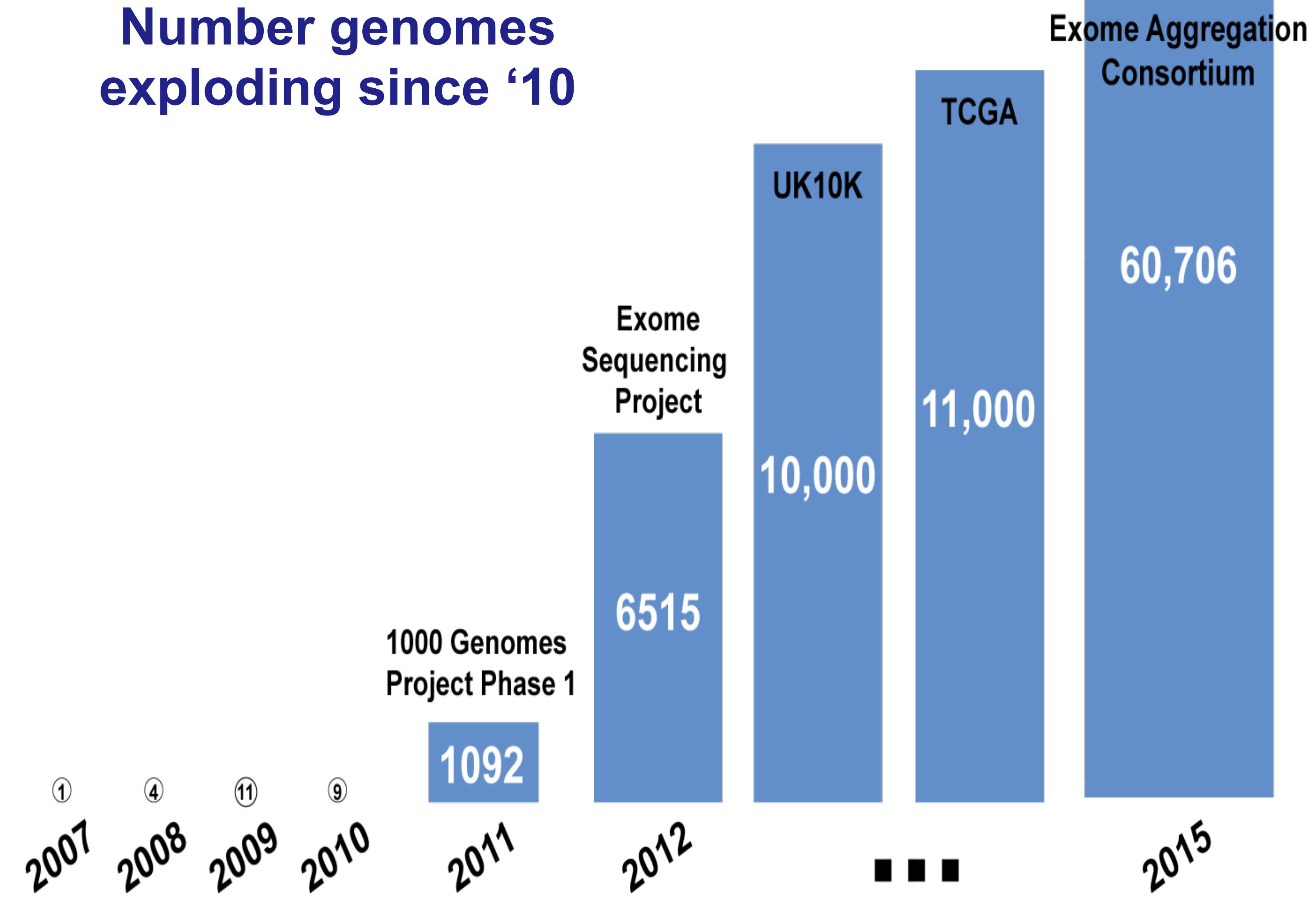


Fig 1. Growth of DNA sequencing. The plot shows the growth of DNA sequencing both in the total number of human genomes sequenced (left axis) as well as the worldwide annual sequencing capacity (right axis: Tera-basepairs (Tbp), Peta-basepairs (Pbp), Exa-basepairs (Ebp), Zetta-basepairs (Zbps)). The values through 2015 are based on the historical publication record, with selected milestones in sequencing (first Sanger through first PacBio human genome published) as well as three exemplar projects using large-scale sequencing: the 1000 Genomes Project, aggregating hundreds of human genomes by 2012 [3]; The Cancer Genome Atlas (TCGA), aggregating over several thousand tumor/normal genome pairs [4]; and the Exome Aggregation Consortium (ExAC), aggregating over 60,000 human exomes [5]. Many of the genomes sequenced to date have been whole exome rather than whole genome, but we expect the ratio to be increasingly favored towards whole genome in the future. The values beyond 2015 represent our projection under three possible growth curves as described in the main text.

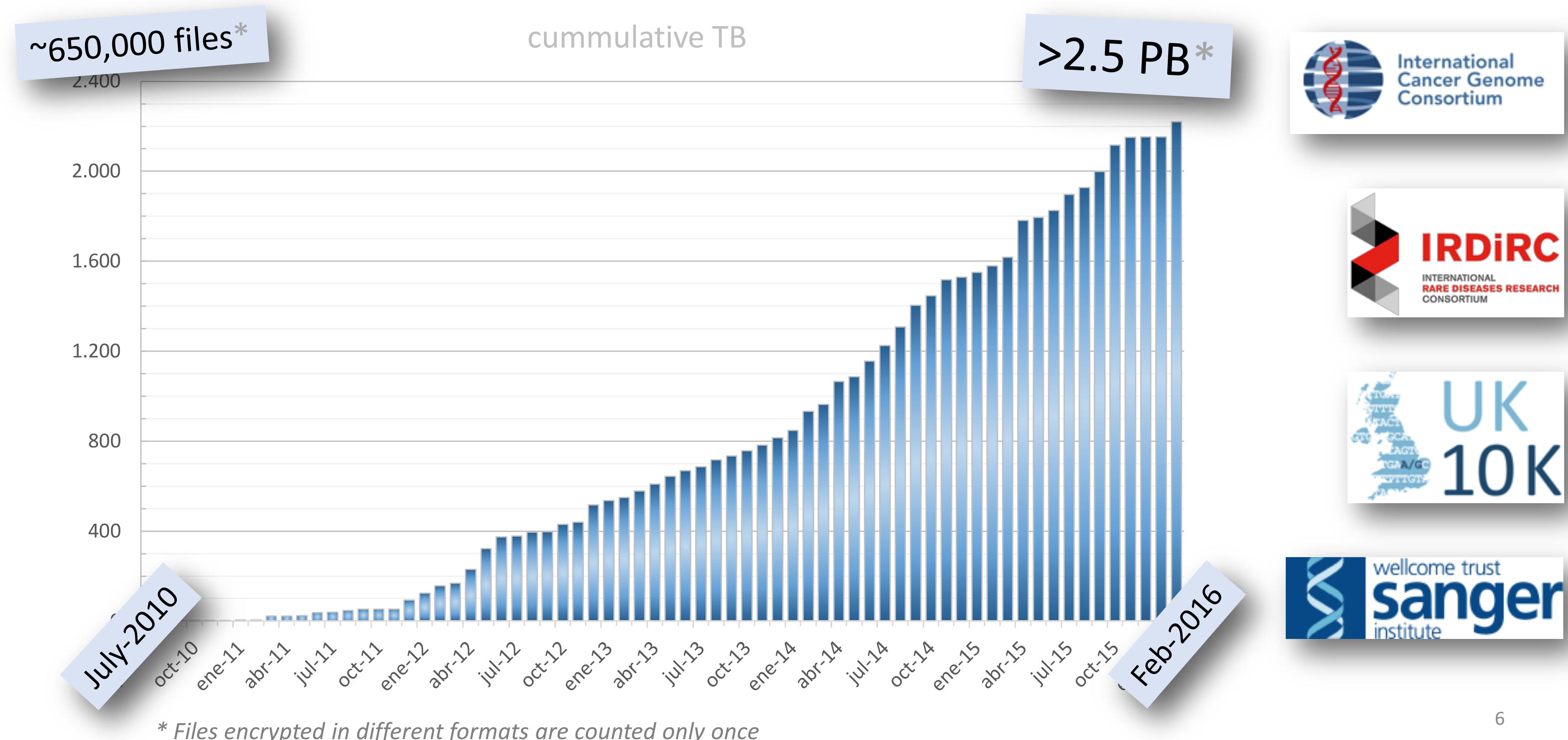


- Labor, administration, management, utilities, reagents, and consumables
- Sequencing instruments and other large equipment (amortized over three years)
- Informatics activities directly related to sequence production (e.g., laboratory information management systems and initial data processing)
- Submission of data to a public database
- Indirect Costs (<http://oamp.od.nih.gov/dfas/faq/indirect-costs#difference>) as they relate to the above items



Growth of Genome Data Repositories: Example EGA

The EGA contains a growing amount of data



What is a PB, for human genomes? It depends.

- 2 bits per base are sufficient to encode TCGA
 - using 00, 01, 10, 11
 - [TCGA]{3'000'000'000}
 - $2 * 3 * 10^9 \text{ b} = 6,000,000,000 \text{ b}$
 - perfect genome (no overhead): ~715 MB
 - 1PB => ~1400000 genomes
- according to Swiss online store (Sep 2017) ~45'000CHF (100x10TB disks)
- this is less than a PhD position per year in Switzerland ...
- (real costs are 2x that, + duplication, facilities, service ... => ~500'000CHF)
- **However: A single 30x BAM file => 100GB**
- Still: 500'000CHF => 1PB => 10'000 genomes => 50CHF/genome (BAM format)



Bioinformatics: File Formats

- **text** or **binary** file formats, optimised for specific types of biological data
- examples from genomics:
 - **BAM** - compressed binary version of Sequence Alignment/ Map (SAM)
 - **BED** (Browser Extensible Data) - flexible way to define the data lines in an genome browser annotation tracks
 - **VCF** (Variant Call Format)

GSM1904006.CEL 69.1 MB
Modified: 3 February 2016 at 17:46

Add Tags...

General:

Kind: FLC animation
Size: 69'078'052 bytes (69.1 MB on disk)
Where: arrayRAID • arraymapln • affyRaw
→ GSE73822 • GPL6801
Created: 3 February 2016 at 17:46
Modified: 3 February 2016 at 17:46
 Stationery pad
 Locked

More Info: Name & Extension: Comments: Open with:

QuickTime Player (default)
Use this application to open all documents like this one.

Change All

Preview:

not a movie...

itemRgb="On"

browser position chr7:127471196-127495720
browser hide all
track name="ItemRGBDemo" description="Item RGB Demo"
chr7 127471196 127472363 Pos1 0 + 127472363
chr7 127472363 127473530 Pos2 0 + 127472363 127473530 255,0,0
chr7 127473530 127474697 Pos3 0 + 127473530 127474697 255,0,0
chr7 127474697 127475864 Pos4 0 + 127474697 127475864 255,0,0
chr7 127475864 127477031 Neg1 0 - 127475864 127477031 0,0,255
chr7 127477031 127478198 Neg2 0 - 127477031 127478198 0,0,255
chr7 127478198 127479365 Neg3 0 - 127478198 127479365 0,0,255
chr7 127479365 127480532 Pos5 0 + 127479365 127480532 255,0,0
chr7 127480532 127481699 Neg4 0 - 127480532 127481699 0,0,255

BED file example

The VCF file format

Standard for variant representation

Example

VCF header

```

##fileformat=VCFv4.0
##fileDate=20100707
##source=VCFtools
##reference=NCBI36
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT SAMPLE1 SAMPLE2
1 1 . ACG A,AT . PASS .
1 2 rs1 C T,CT . PASS H2 ; AA=T 0|1:100 2/2:70
1 5 . G <DEL> . PASS .
1 100 . T . PASS SVTYPE=DEL ; END=300 GT:DP 1/2:13 0/0:29

```

Body

Mandatory header lines

Optional header lines (meta-data about the annotations in the VCF body)

Reference alleles (GT=0)

Alternate alleles (GT>0 is an index to the ALT column)

Phased data (G and C above are on the same chromosome)

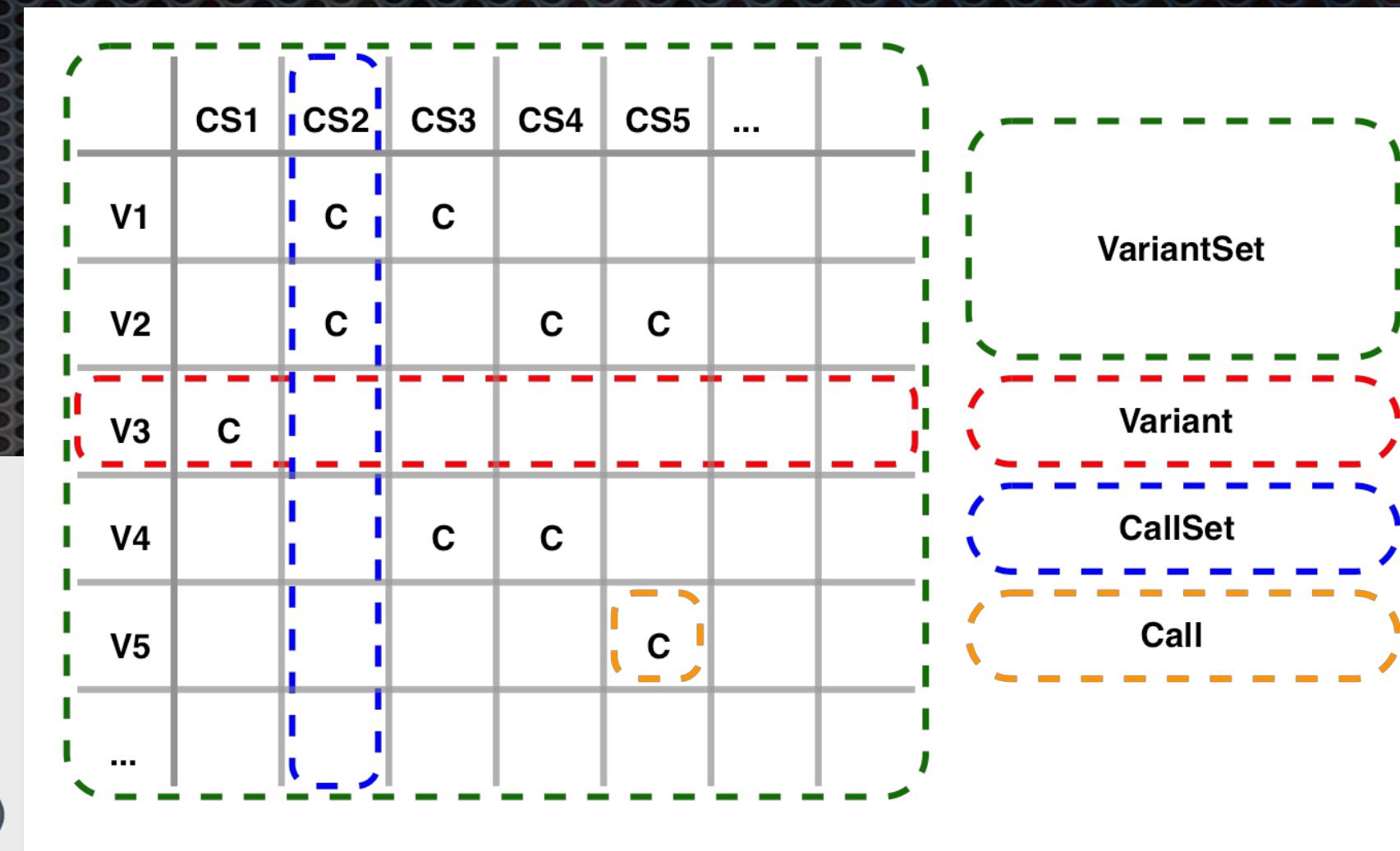
Deletion

SNP

Large SV

Insertion

Other event



Variant
Call
Format

- stores the results of a single or multiple interpretations of genome sequencing datasets, in comparison to a reference genome
- standard format for file-based storage of human genome variants



The trouble with human genome variation

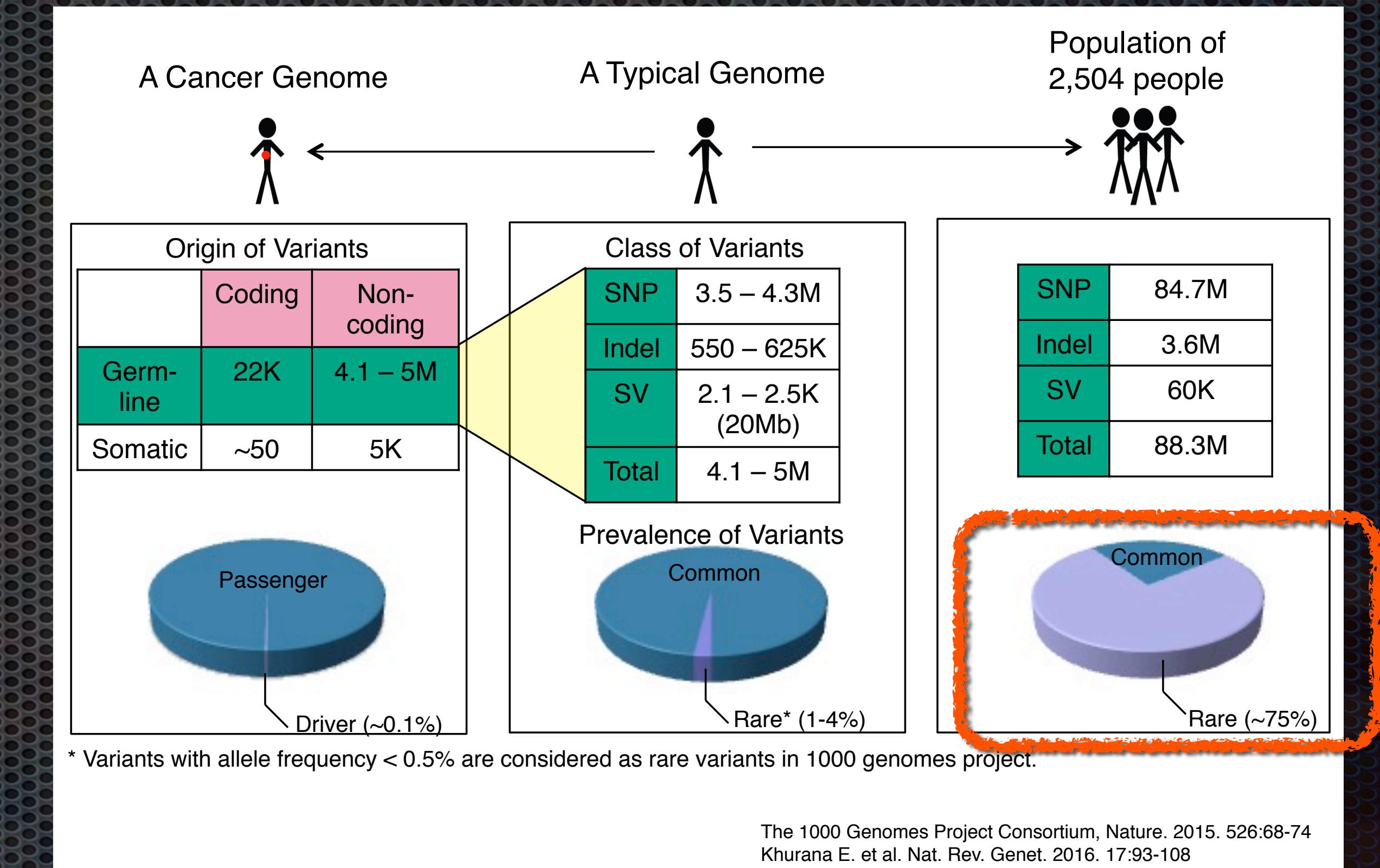


Conclusions from the analysis of variation in the human genome

- 1. Humans are all very similar to each other
 - Two humans will show about 99.9% sequence identity with each other. In other words, only about 1 in 1'000 bp is different between two individuals.
 - Humans show about 98% sequence identity to chimps. So two humans are still much more similar to each other than either is to the monkey.
- 2. Humans are very different from each other
 - Two typical humans will likely have over 1'000'000 independent sequence differences in their genomes.

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

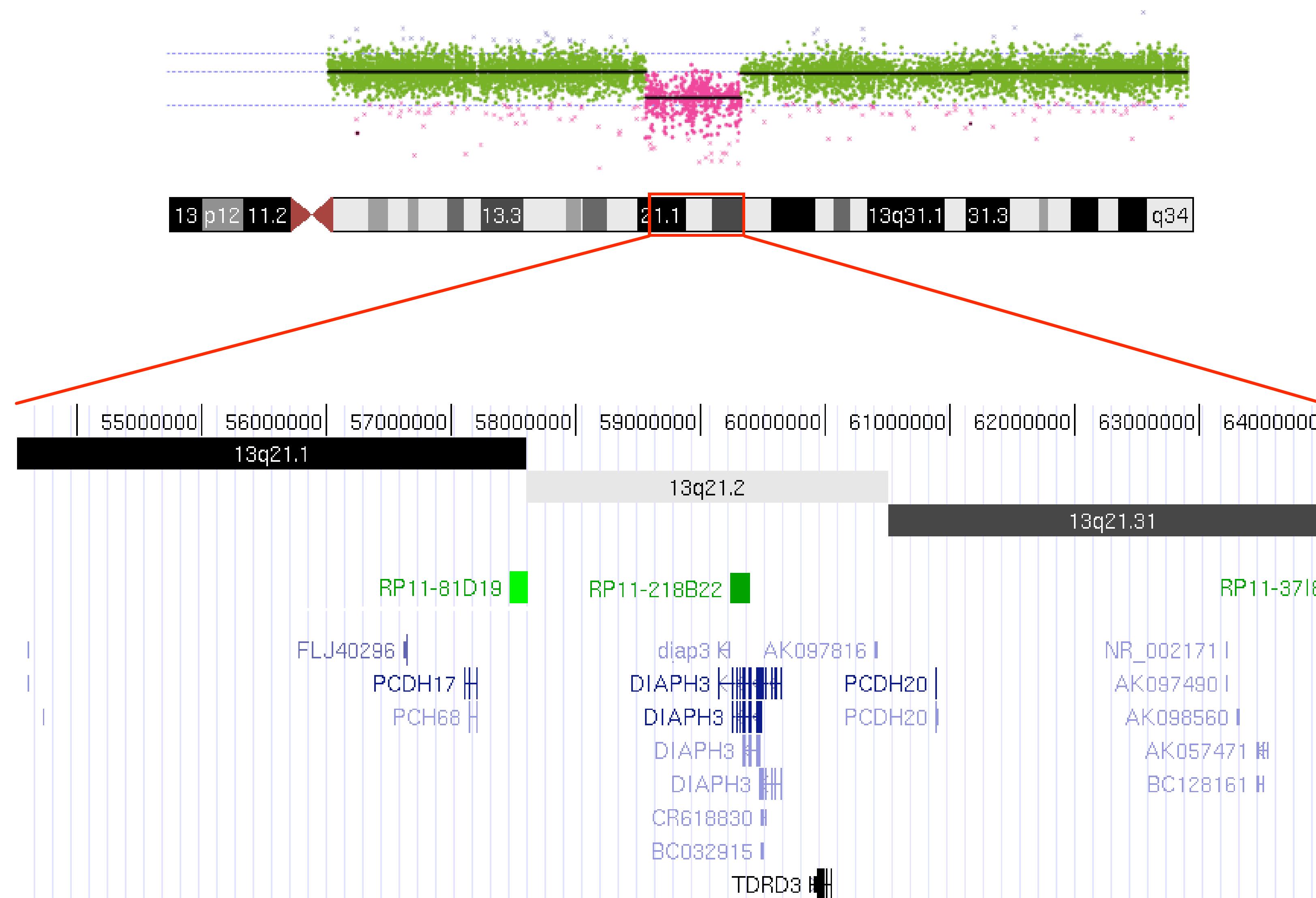
- a typical human genome (~3 billion base pairs) has ~5 million variants
- most of them are "**rare**"; i.e. can only be identified as recurring when sequencing thousands of people
- cancer cells accumulate additional variants, only **few** of which ("**drivers**") are relevant for the disease



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

Nobody is perfect (?)

A 10.7 Mb Interstitial Deletion of 13q21 Without Phenotypic Effect Defines a Further Non-Pathogenic Euchromatic Variant
Andreas Roos, Miriam Elbracht, Michael Baudis, Jan Senderek, Nadine Schönherr, Thomas Eggemann, and Herdit M. Schüler
American Journal of Medical Genetics Part A 146A:2417 – 2420 (2008)



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

Reference Resources for Human Genome Variants

- NCBI:dbSNP

- single nucleotide polymorphisms (SNPs) and multiple small-scale variations
 - including insertions/deletions, microsatellites, non-polymorphic variants

- NCBI:dbVAR

- genomic structural variation
 - insertions, deletions, duplications, inversions, multinucleotide substitutions, mobile element insertions, translocations, complex chromosomal rearrangements

- NCBI:ClinVar

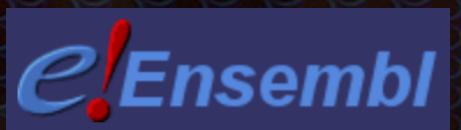
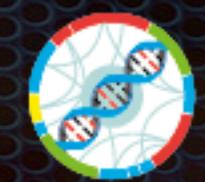
- aggregates information about genomic variation and its relationship to human health

- EMBL-EBI:EVA

- open-access database of all types of genetic variation data from all species

- Ensembl

- portal for many things genomic...



Personalized medicine in cancer

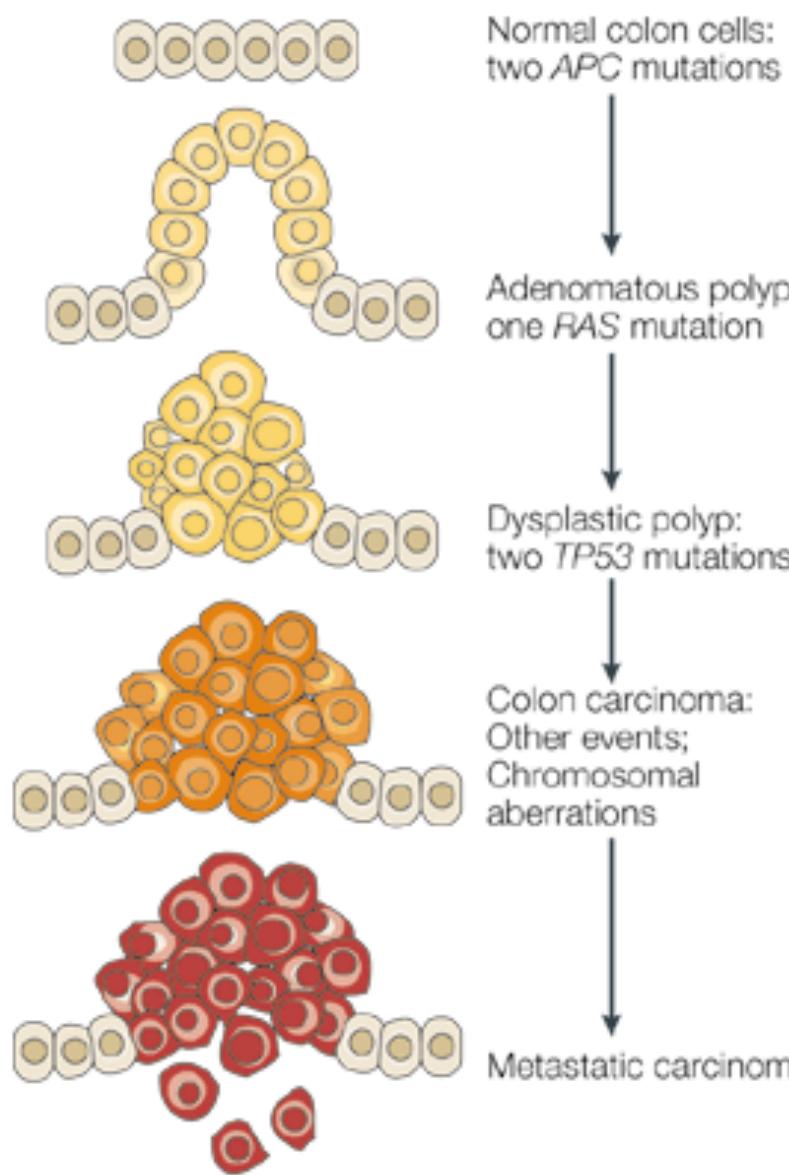
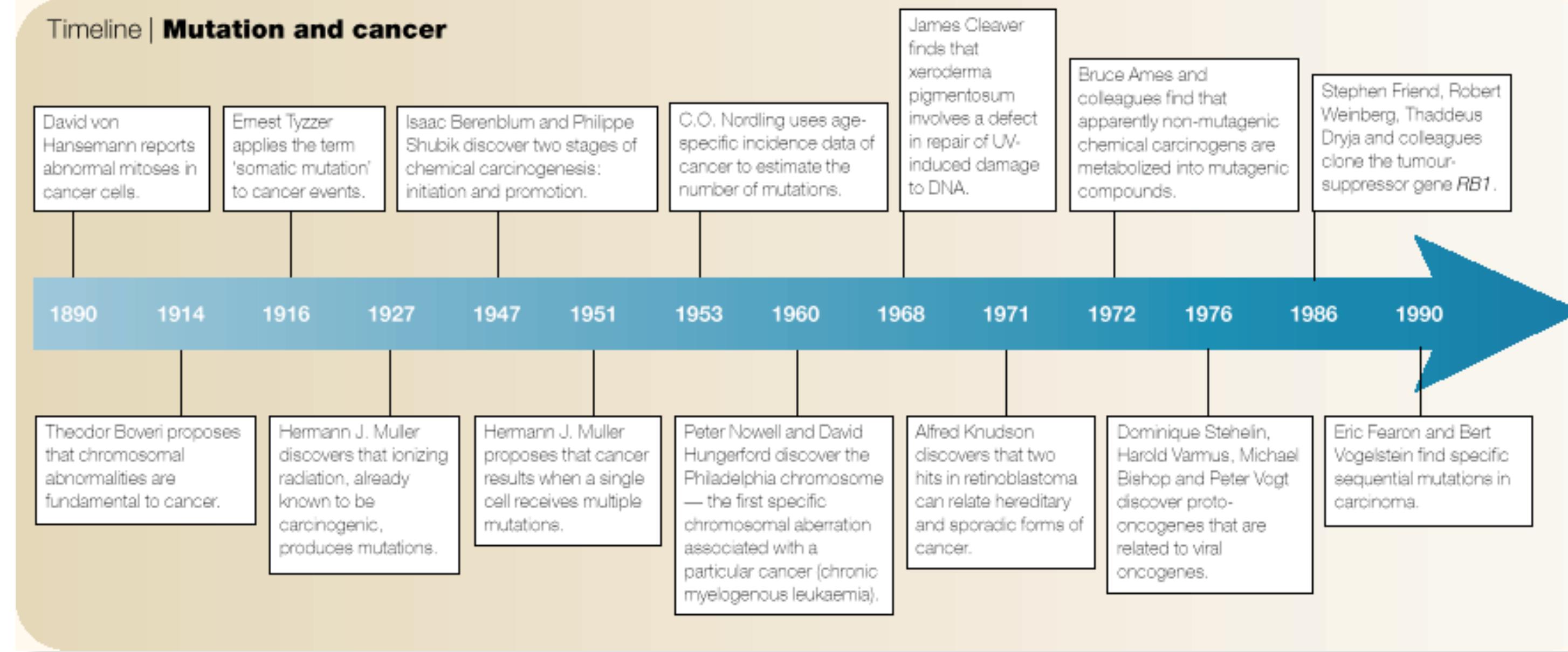
Data Formats | Genome Variation | Techniques | Resources | Sharing

Michael Baudis **UZH SIB**
Computational Oncogenomics



University of
Zurich^{UZH}

Timeline | Mutation and cancer

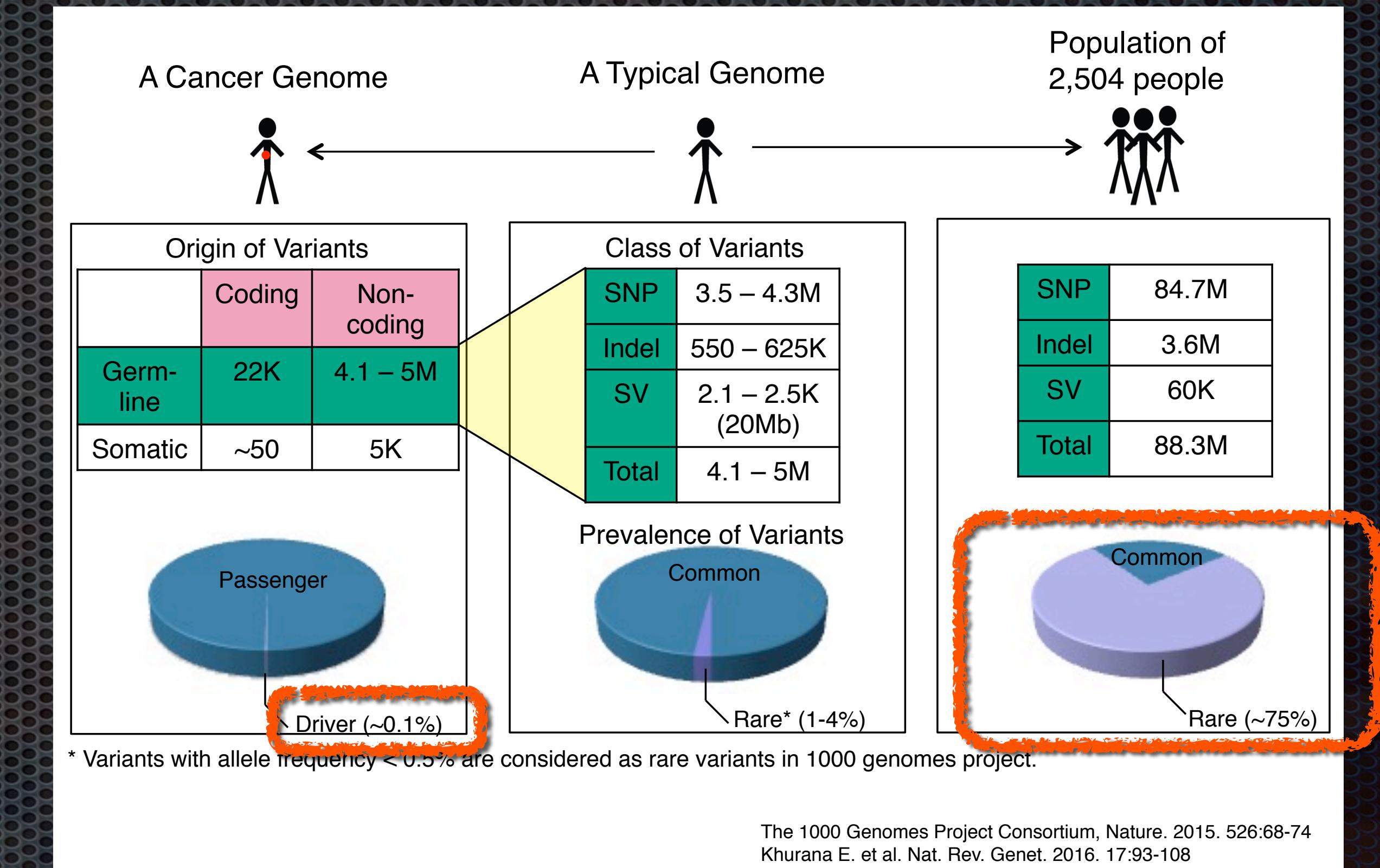


Cancers are based on acquired and inherited genomic mutations

Knudson, A. G. (2001). Two genetic hits (more or less) to cancer. Nature Reviews Cancer, 1(2), 157–162.

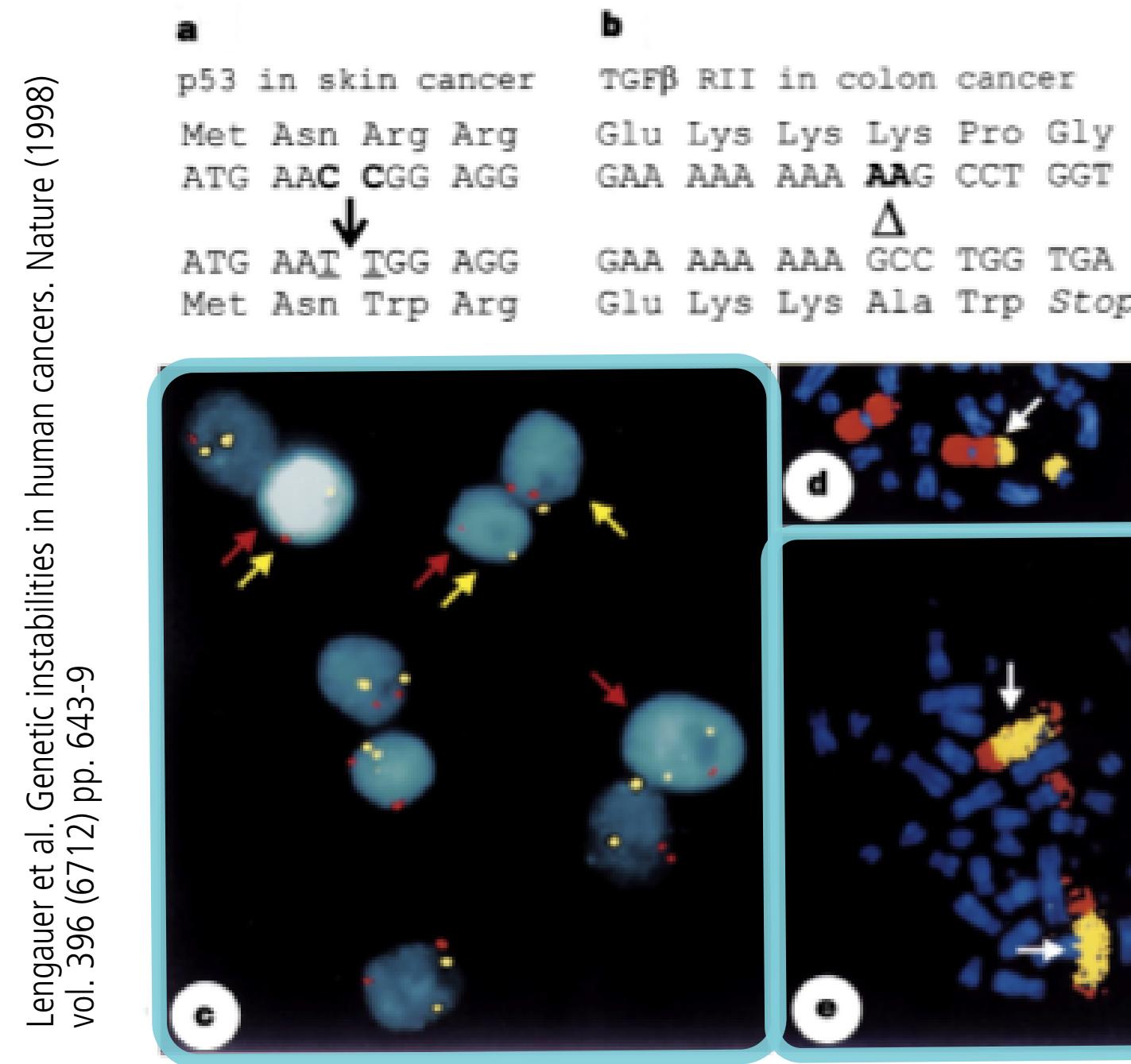
Finding Somatic Mutations In Cancer: Many Needles in a Large Haystack

- 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

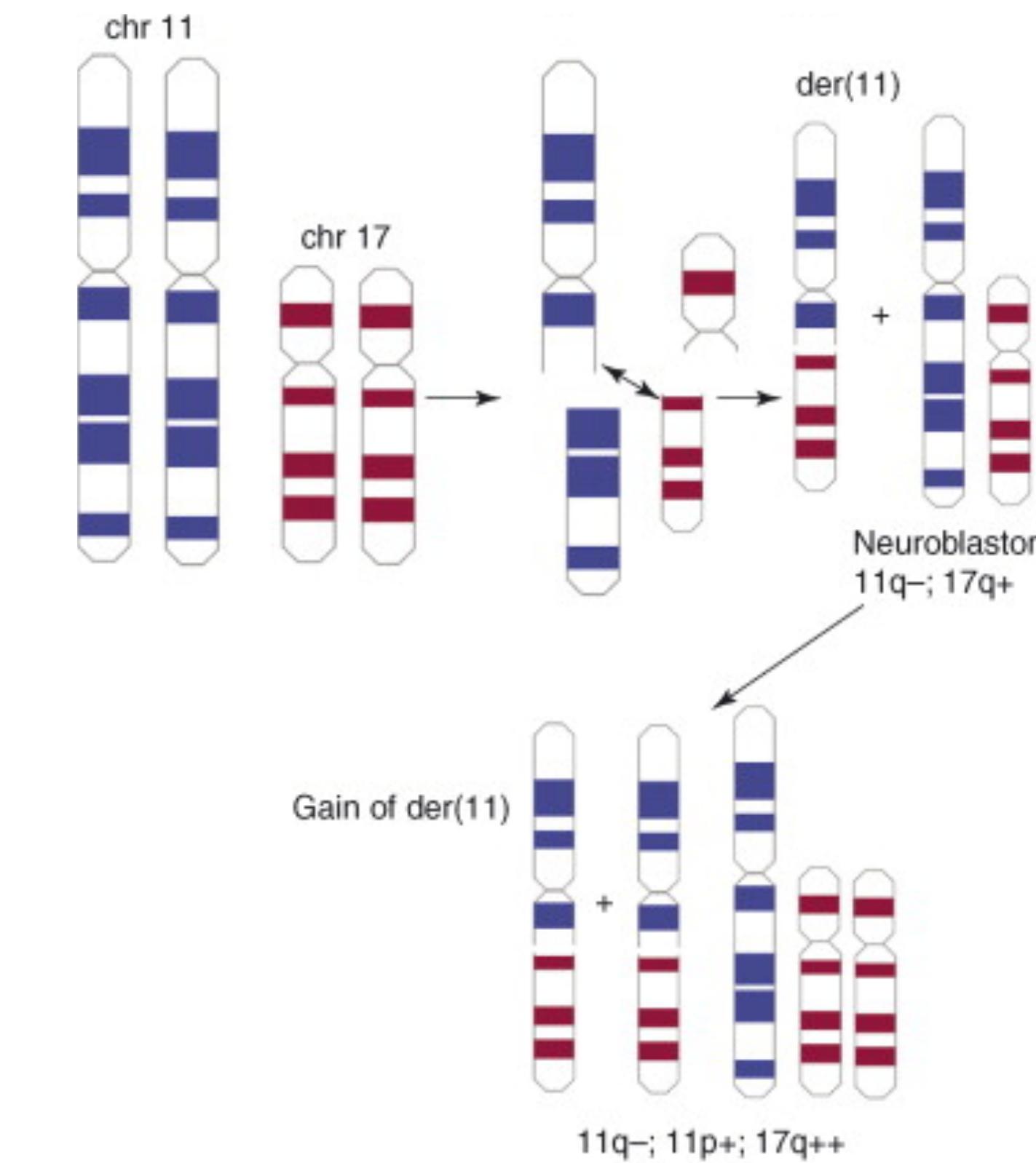


Graphic adapted from Mark Gerstein (GersteinLab.org; @markgerstein)

Mutations & genomic rearrangements in cancer

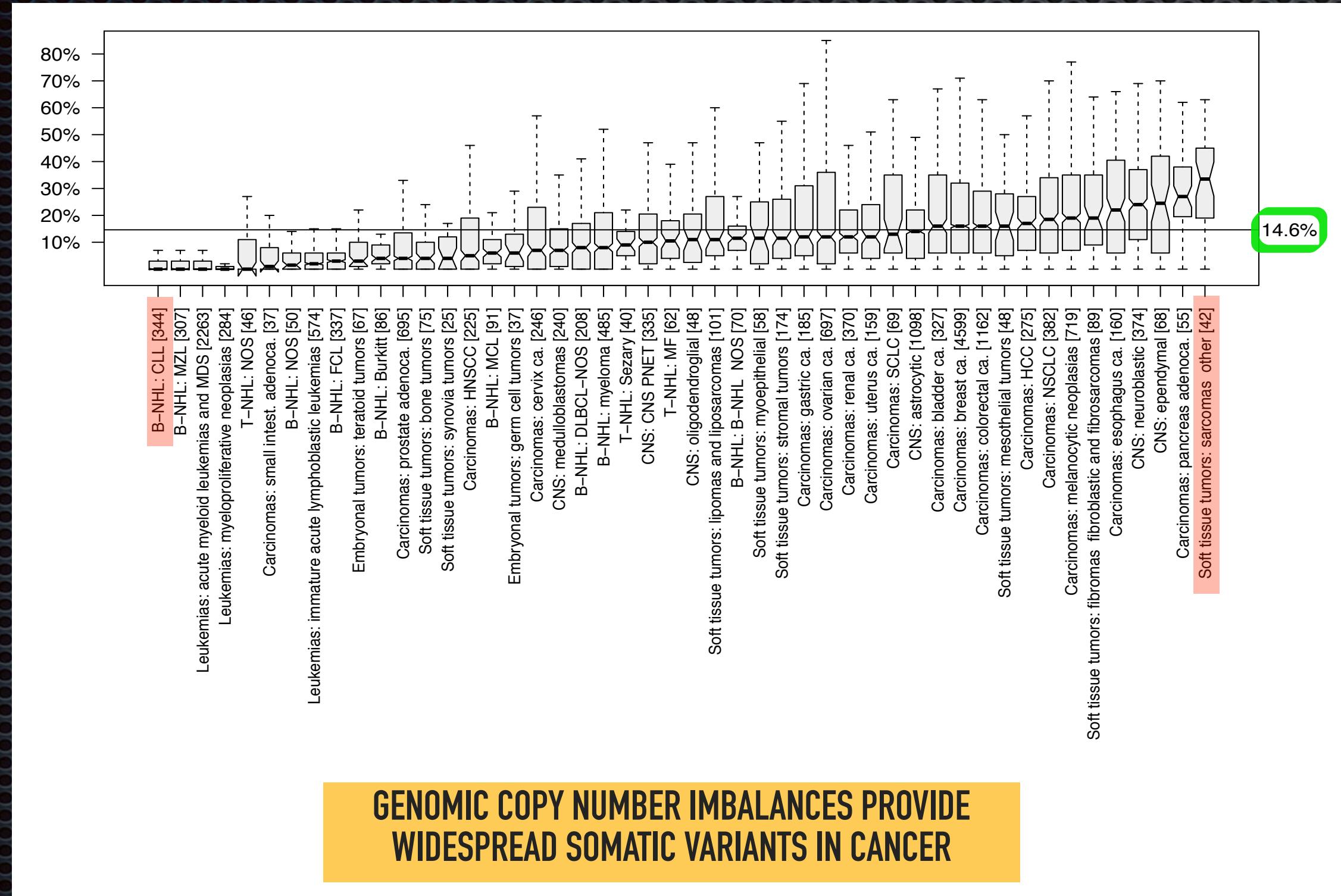


- a. small mutation (di-pyrimidine exchange at p53 in Xeroderma pigmentosum patient)
- b. two-base deletion in *TGFB* in a colorectal cancer patient with mismatch repair deficiency
- c. chromosomal losses (FISH; red=3, yellow=12) in CRC
- d. t(1;17) in neuroblastoma, whole-chromosomal painting
- e. *MYCN* gene amplification (multiple copies inserted into chromosome 1 derived marker)

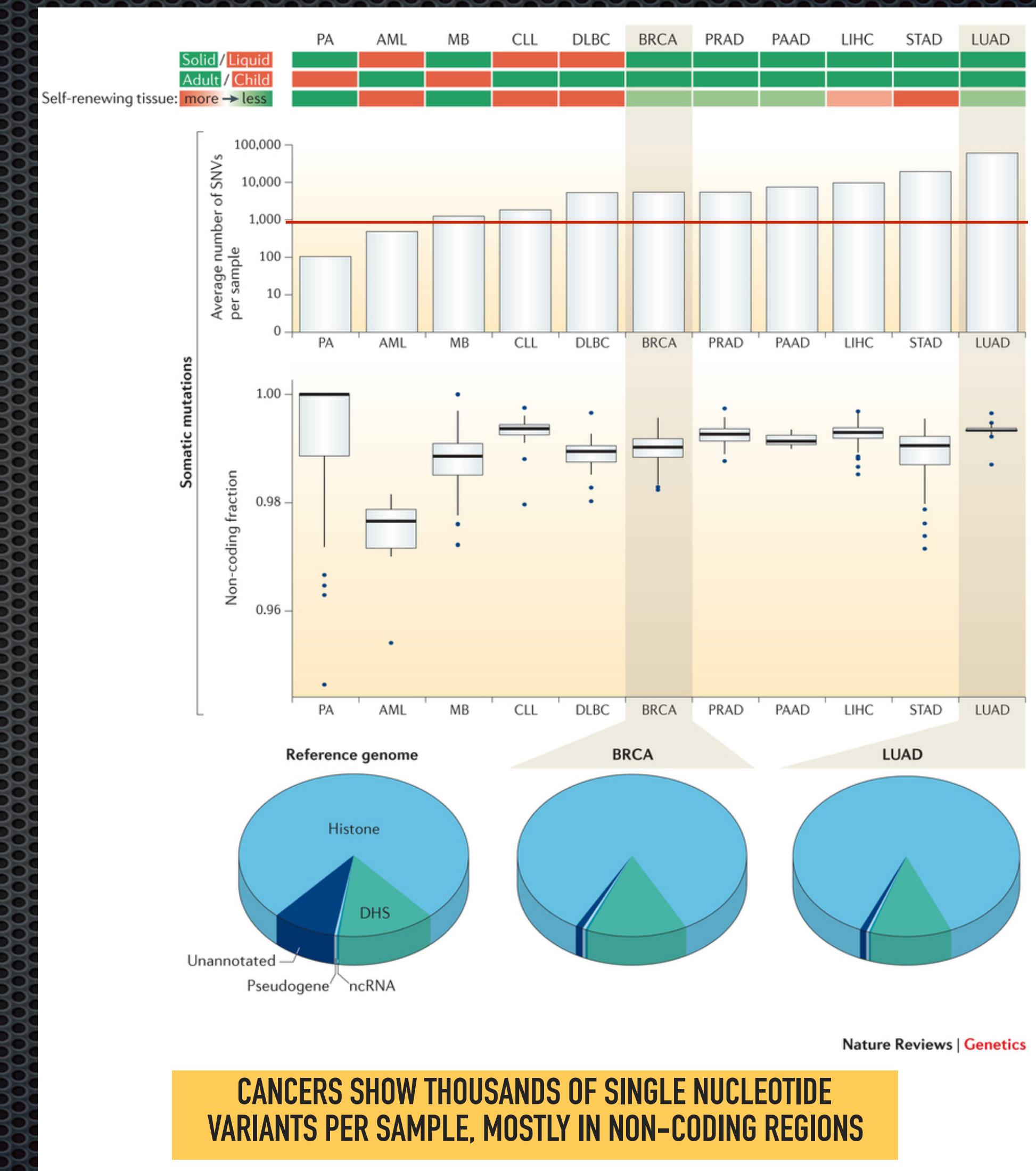


Generation of copy number imbalances in cancer through imbalanced cytogenetic rearrangements - partial deletion of 11q, gain of 11pterq21 and 2 addl. copies of 17q

Quantifying Somatic Mutations In Cancer

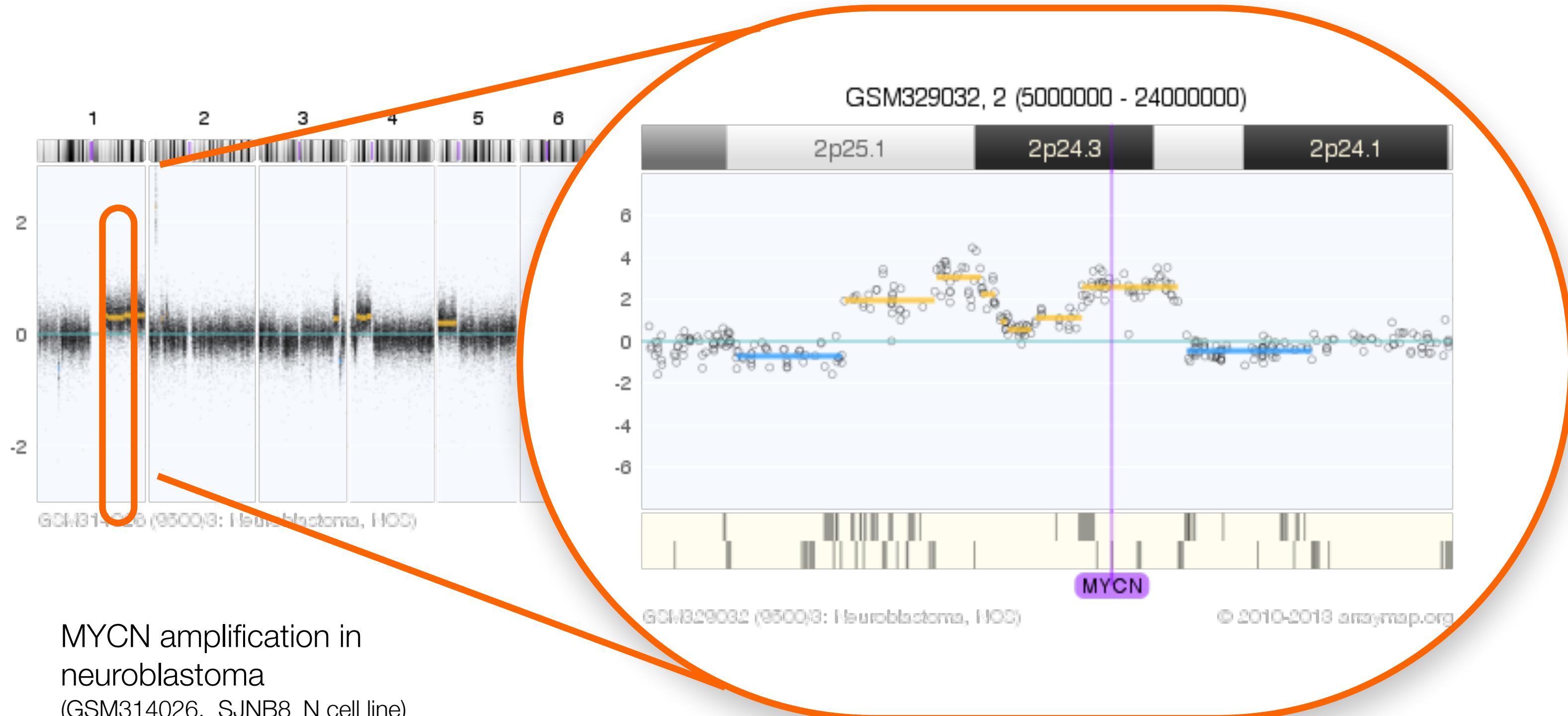
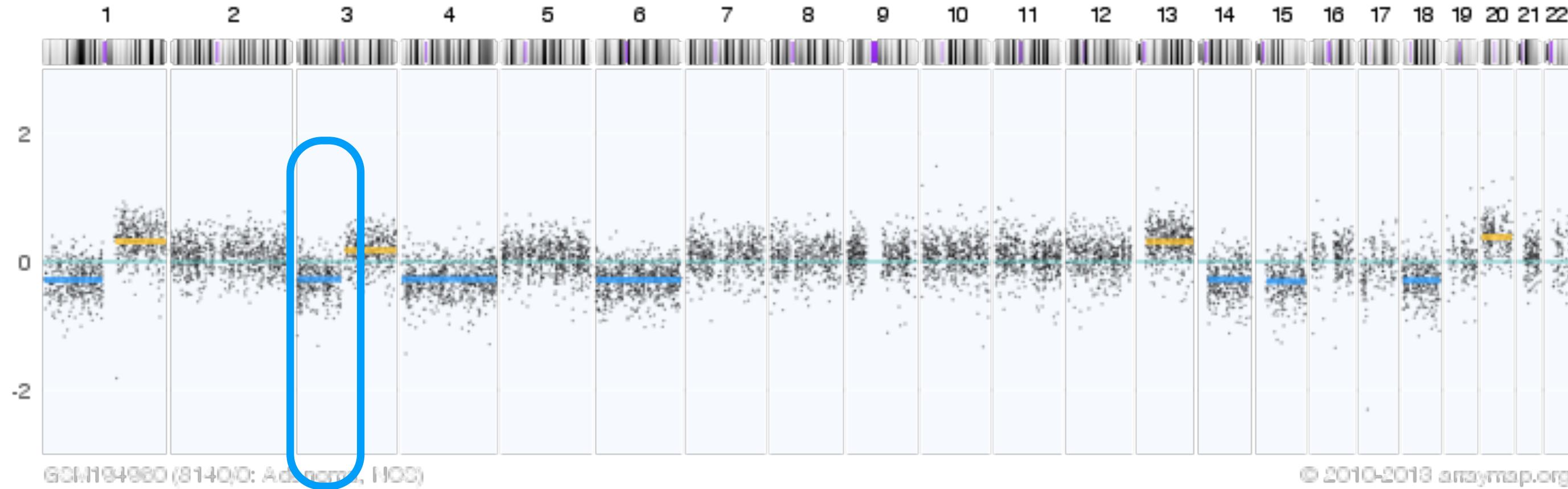


On average ~15% of a cancer genome are in an imbalanced state (more/less than 2 alleles);
Original data based on >30'000 cancer genomes from arraymap.org



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))

Genomic arrays: Many probes + bioinformatics determine copy number aberrations

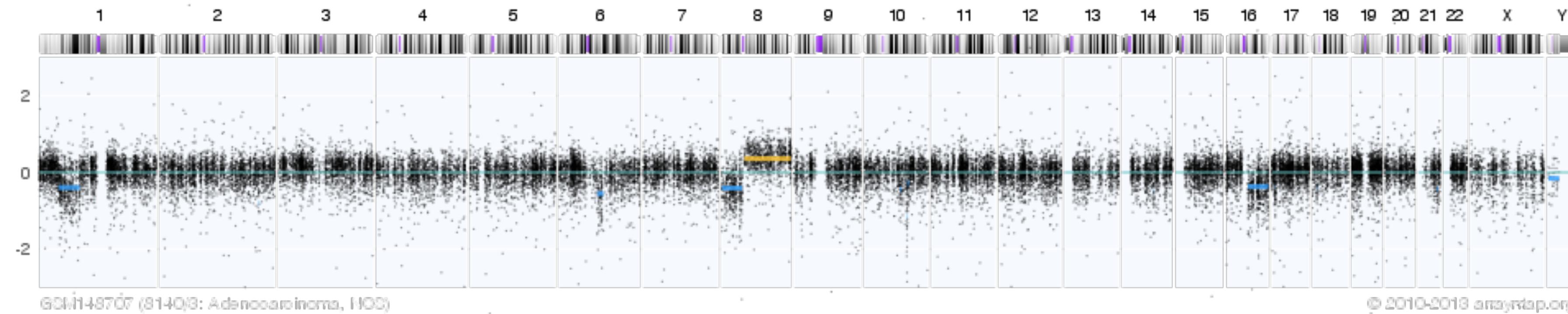
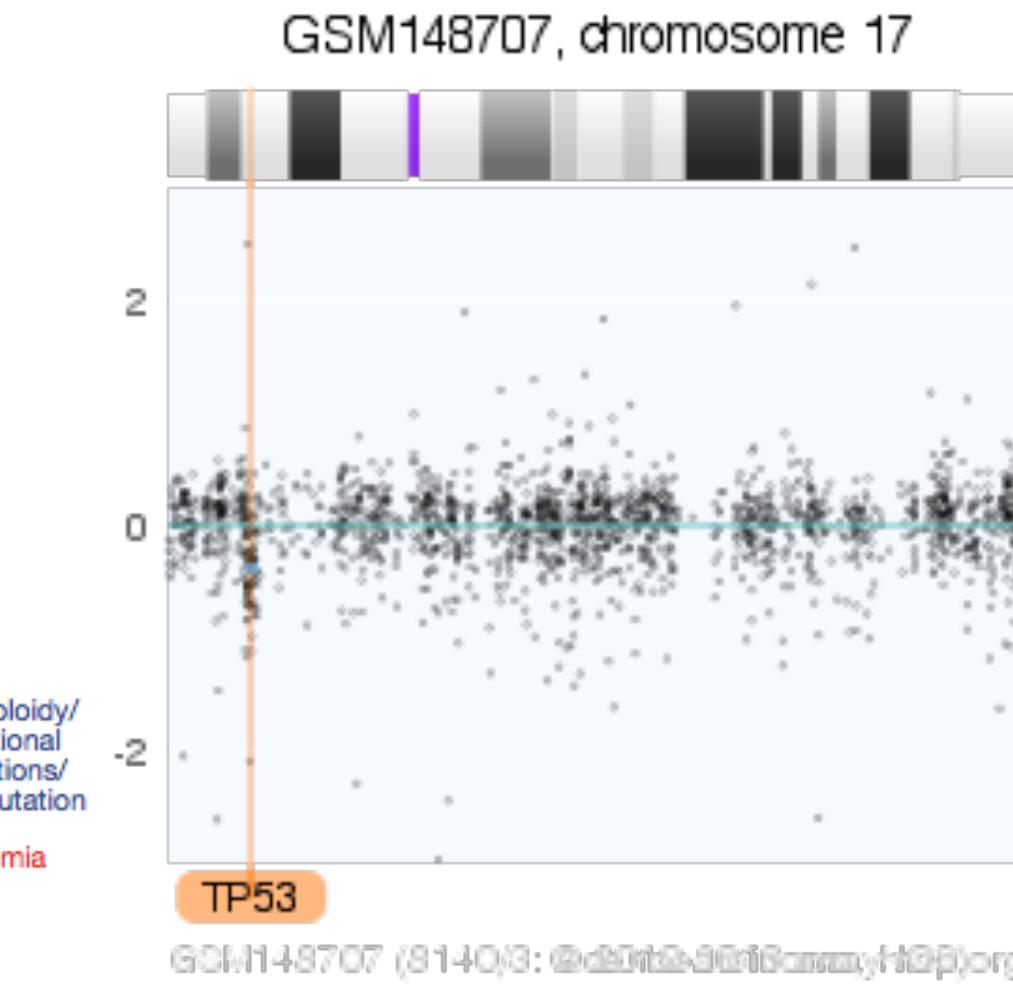
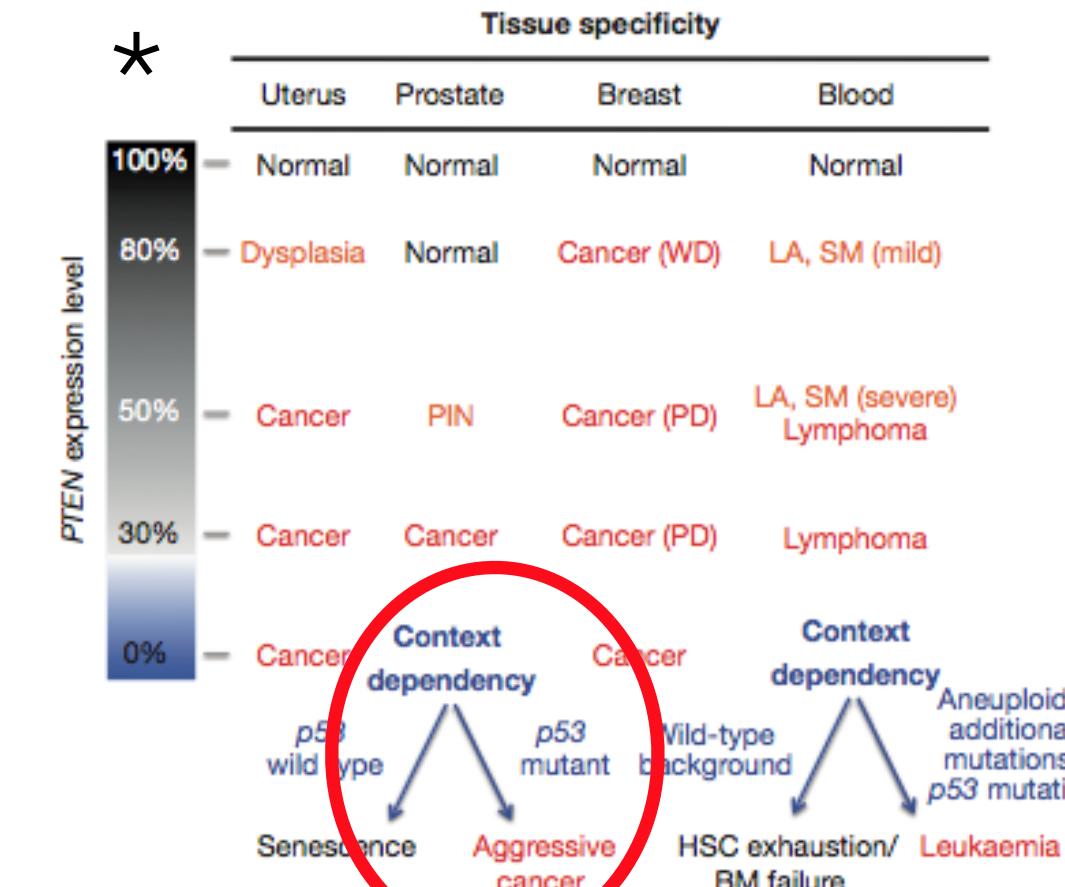
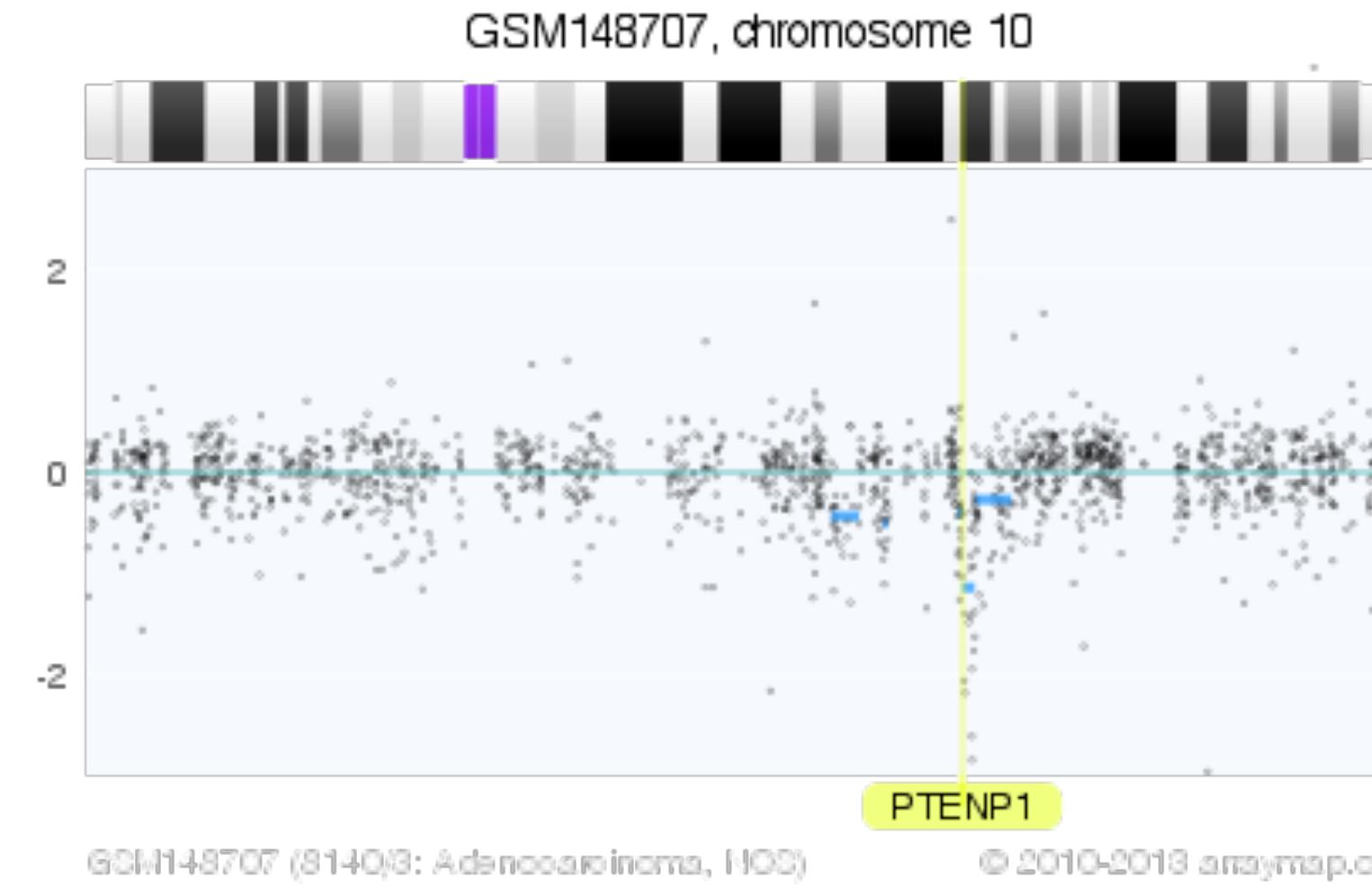


low level/high level copy number alterations (CNAs)

arrayMap



Gene dosage phenomena beyond simple on/off effects



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

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



Challenges in aCGH data collection: Constitutional CNVs vs. imbalances

Segmental copy number variations of unique sequences have been described as normal feature of human DNA.

These CNVs may be up to 2-3 megabases in size, and can involve coding regions.

Everybody has them ...

Lockwood et al. Recent advances in array comparative genomic hybridization technologies and their applications in human genetics.
Eur J Hum Genet (2006) vol. 14 (2) pp. 139-48

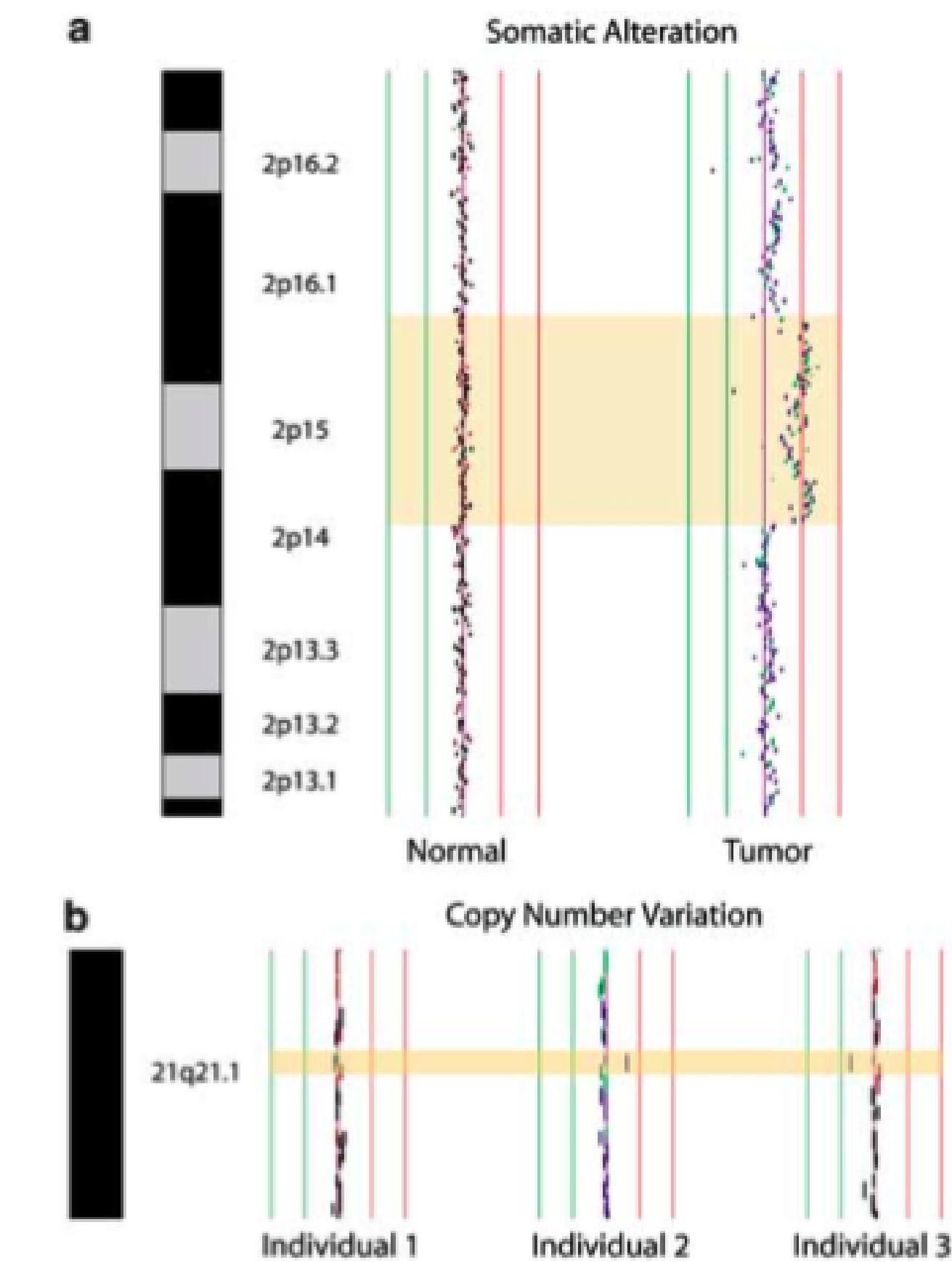


Figure 3 Somatic alterations and copy number variations. (a) Example of a segmental duplication observed at chromosome arm 2p present in the cancer cells but absent in the normal cells from the same individual. Each black dot represents a single BAC clone spotted on the array. The purple line represents equal fluorescent intensity ratio between sample and reference. Copy number gain (and loss) shifts the ratio to the right (and left). (b) Illustrates a copy number variation observed at chromosomal region 21q21.1. Three normal individuals exhibit equal, more and fewer copies relative to the reference DNA, indicating variation in the population.

CNV vs. CNA: Size matters

- no unambiguous criterium for CNV vs. CNA
- statistic argument: CNVs are recurring copy number variations found in the germline DNA of “healthy” individuals
- size argument: CNVs are rarely larger than 1Mb

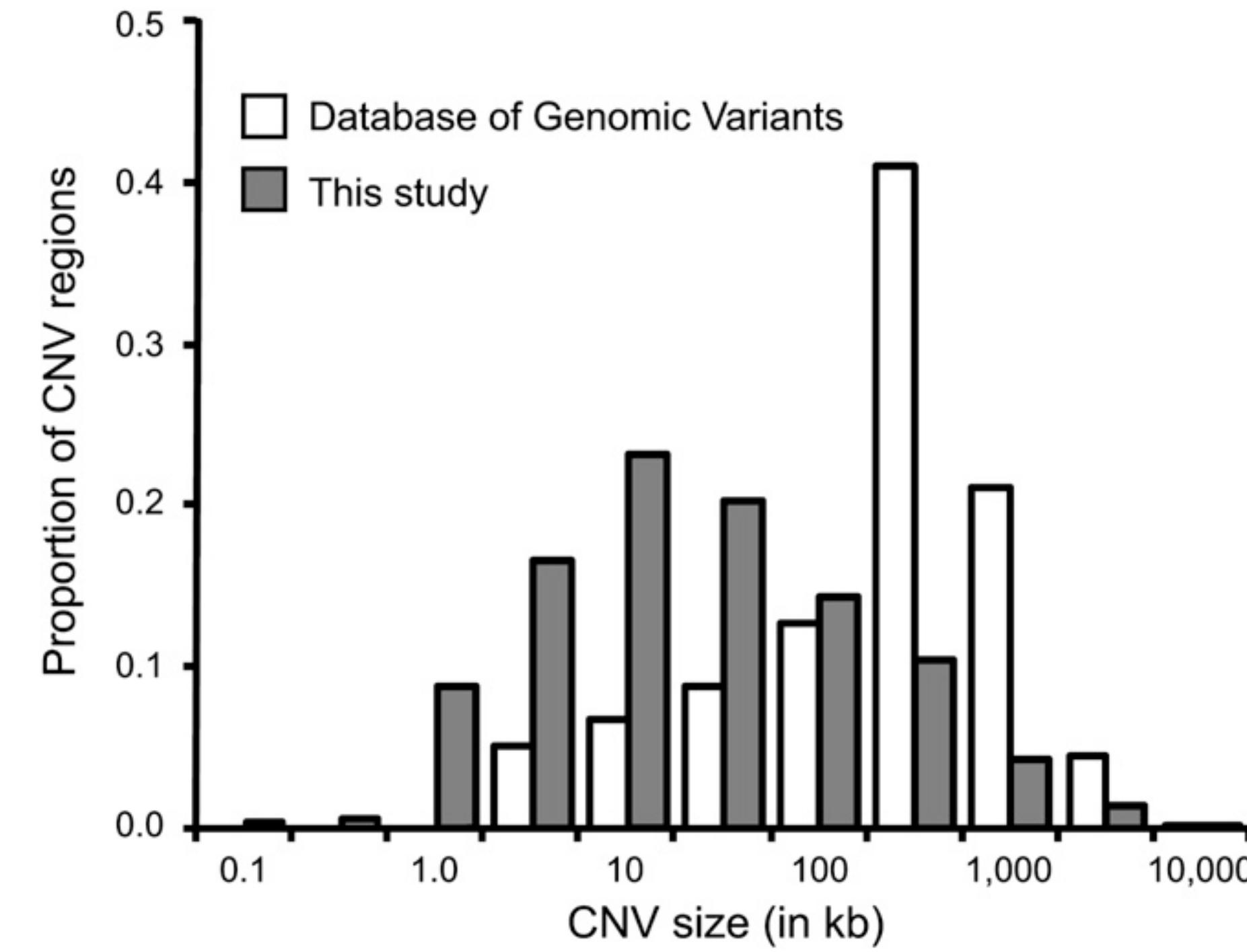


Figure 1. Size Distribution of CNVs from the Database of Genomic Variants, with Corresponding CNVs from This Study

We identified CNVs in at least one individual for 1153 of 2191 putative CNV regions annotated in the Database of Genomic Variants (DGV) as of 30 November 2006. Size distributions for these regions are shown in log scale, with 10-fold multiples of 1 and $\sqrt{10}$, based on the size of each region from DGV and the estimates from our study of the total amount of copy-number-variable sequence within and overlapping the DGV-defined region. Our estimates were smaller than the corresponding DGV region for 1020 of the 1153 loci (88%) and smaller by more than 50% for 876 regions (76%).

Where to find cancer genome data ...

RESOURCES FOR CANCER GENOMICS

COSMIC
Catalogue of somatic mutations in cancer

Home ▾ Resources ▾ Curation ▾ Tools ▾ Data ▾ News ▾ Help ▾ About ▾ Search COSMIC... Login ▾

COSMIC v79, released 14-NOV-16

COSMIC, the Catalogue Of Somatic Mutations In Cancer, is the world's largest and most comprehensive resource for exploring the impact of somatic mutations in human cancer.

Start using COSMIC by searching for a gene, cancer type, mutation, etc. below, or by browsing a region of the human genome using the map to the right.

eg: *Braf*, COLO-829, Carcinoma, V600E, BRCA-UK, Campbell **SEARCH**

R Resources

Key COSMIC resources

- Cell Lines Project
- COSMIC
- Whole Genomes
- Cancer Gene Census
- Drug Sensitivity
- Mutational Signatures
- GRCh37 Cancer Archive

T Tools

Additional tools to explore COSMIC

- Cancer Browser
- Genome Browser
- GA4GH Beacon
- CONAN

C Expert Curation

High quality curation by expert postdoctoral scientists

- Drug Resistance
- Cancer Gene Census
- Curated Genes
- Gene Fusions
- Genome-Wide Screens

D Data

Further details on using COSMIC's content

- Downloads
- License
- Submission
- Genome Annotation
- Datasheets
- Help
- FAQ

Browse the [genomic landscape of cancer](#)

Cancer Gene Census Update

7 genes have been added to the [Cancer Gene Census](#) -

- EPAS1 - Endothelial PAS domain protein 1.
- PTPRT - Protein tyrosine phosphatase, receptor type T.
- PPM1D - Protein phosphatase, Mg²⁺/Mn²⁺ dependent 1D.
- BTK - Bruton tyrosine kinase.
- PREX2 - Phosphatidylinositol-3,4,5-trisphosphate dependent Rac exchange factor 2.
- TP63 - Tumour protein p63.
- QKI - QKI, KH domain containing RNA binding.

For full details, see the [Datasheet](#).

RESOURCES FOR CANCER GENOMICS

National Cancer Institute U.S. National Institutes of Health | www.cancer.gov

CANCER GENOME ANATOMY PROJECT

CGAP How To

Tools

CGAP Info

- Educational Resources
- Slide Tour
- Team Members
- References

CGAP Data

Quick Links:

- ICG
- NCI Home
- NCICB Home
- NCBI Home
- OCG

Genes Genes | **Chromosomes** Chromosomes | **Tissues** Tissues | **SAGE Genie** SAGE Genie | **RNAi** RNAi | **Pathways** Pathways

Cancer Genome Anatomy Project (CGAP)

The NCI's Cancer Genome Anatomy Project sought to determine the gene expression profiles of normal, precancer, and cancer cells, leading eventually to improved detection, diagnosis, and treatment for the patient. Resources generated by the CGAP initiative are available to the broad cancer community. Interconnected modules provide access to all CGAP data, bioinformatic analysis tools, and biological resources allowing the user to find "in silico" answers to biological questions in a fraction of the time it once took in the laboratory.

The CGAP Website

Interconnected modules provide access to all CGAP data, bioinformatic analysis tools, and biological resources allowing the user to find "in silico" answers to biological questions in a fraction of the time it once took in the laboratory.

Genes Gene information, clone resources, SNP500Cancer, GAI, and transcriptome analysis.

Chromosomes FISH-mapped BAC clones, SNP500Cancer, and the Mitelman database of chromosome aberrations.

Tissues cDNA library information, methods, and EST-based gene expression analysis.

SAGE Genie Analysis of gene expression using long and short SAGE tag data for both human and mouse.

Pathways Diagrams of biological pathways and protein complexes, with links to genetic resources for each known protein.

Tools Direct access to all analytic and data mining tools developed for the project.

RNAi RNA-interference constructs, targeted specifically against cancer relevant genes. New addition: Validated set of shRNAs.

International Cancer Genome Consortium

Home Cancer Genome Projects Committees and Working Groups Policies and Guidelines Media

ICGC Cancer Genome Projects

Committed projects to date: 89

Sort by: Project

Biliary Tract Cancer Japan	Biliary Tract Cancer Singapore	Bladder Cancer China
Bladder Cancer United States	Blood Cancer China	Blood Cancer Singapore
Blood Cancer South Korea	Blood Cancer United States	Blood Cancer United States
Blood Cancer United States	Blood Cancer United States	Bone Cancer France
Bone Cancer United Kingdom	Bone Cancer United States	Brain Cancer Canada
Brain Cancer China	Brain Cancer United States	Brain Cancer United States
Breast Cancer China	Breast Cancer European Union / United Kingdom	Breast Cancer France
Breast Cancer Mexico	Breast Cancer South Korea	Breast Cancer South Korea

ICGC Goal: To obtain a comprehensive description of genomic, transcriptomic and epigenomic changes in 50 different tumor types and/or subtypes which are of clinical and societal importance across the globe.

[Read more »](#)

Launch Data Portal »

Apply for Access to Controlled Data »

Announcements

23/August/2016 - The ICGC Data Coordination Center (DCC) is pleased to announce ICGC data portal data release 22 (<http://dcc.icgc.org>).

ICGC data release 22 in total comprises data from more than 16,000 cancer donors spanning 70 projects and 21 tumour sites.

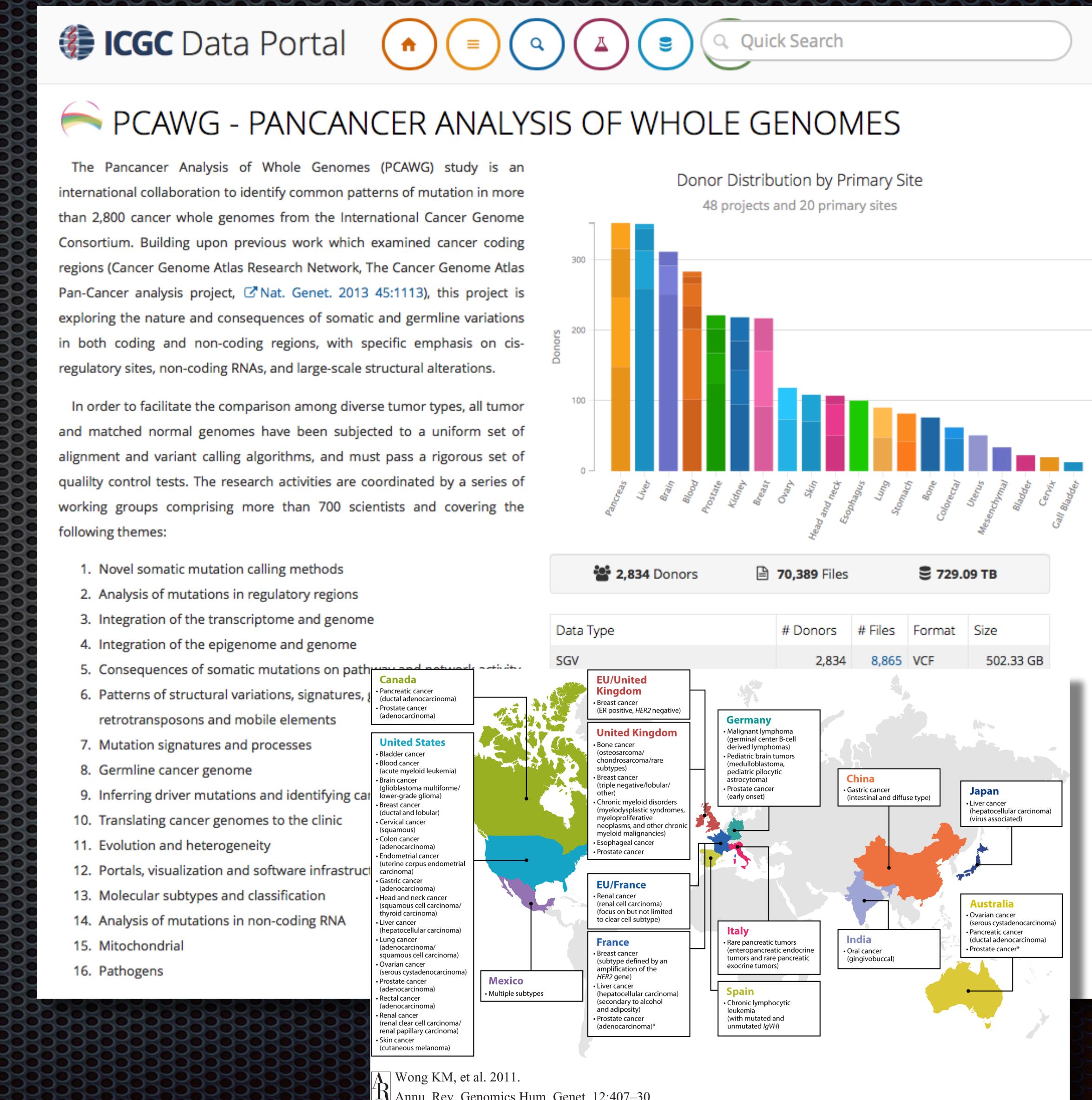
17/April/2016 - ICGCmed is pleased to announce the release of its white paper (<http://icgcmed.org>).

The International Cancer Genome Consortium for Medicine (ICGCmed) will link genomics data to clinical information, health and responses to therapies.

18/November/2015 - The International Cancer Genome Consortium (ICGC) PanCancer dataset generated by the PanCancer Analysis of Whole Genomes (PCAWG) study is now available on Amazon Web Services (AWS), giving cancer researchers access to over 2,400 consistently analyzed genomes corresponding to over 1,100 unique ICGC donors (<https://icgc.org/icgc-in-the-cloud>).

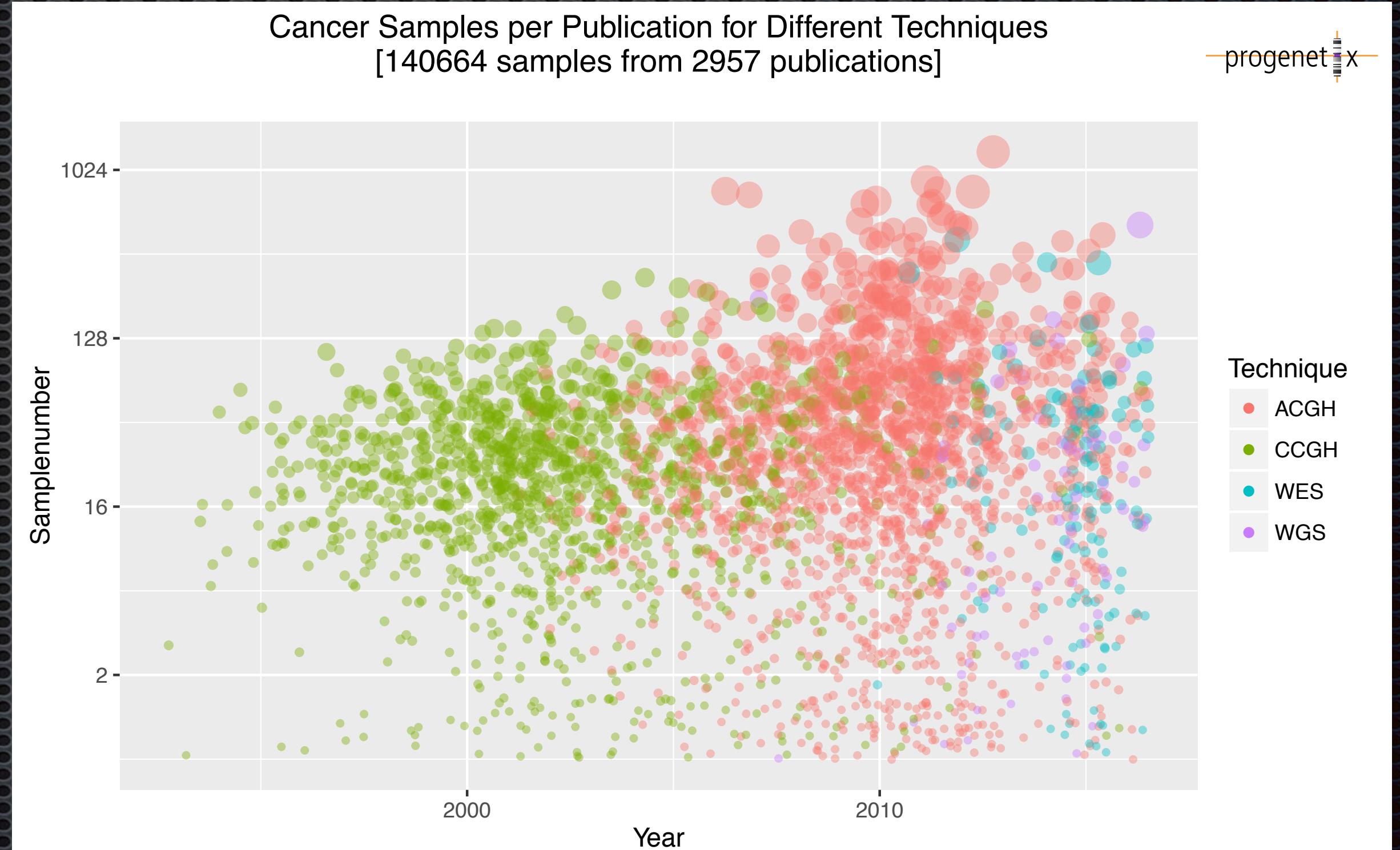
Genome-wide multi"omics" data generation for understanding tumor biology

- 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)

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 ...)

The arrayMap Cancer Genome Resource

arrayMap 

- [Search Samples](#)
- [Search Publications](#)
- [Gene CNA Frequencies](#)
- [User Data](#)
- [Array Visualization](#)
- [Progenetix](#)

 **University of Zurich**

- [Citation](#)
- [User Guide](#)
- [Registration & Licensing](#)
- [People](#)
- [External Links ↗](#)

FOLLOW US ON [twitter](#)

 130.60.23.21

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
-  **ICD-O** 141 ICD-O cancer entities
-  554 publications (Pubmed entries)

Genomic copy number imbalances on chromosome 9 in a case of Glioblastoma ([GSM491153](#)), indicating, among others, a homozygous deletion involving CDKN2A/B.

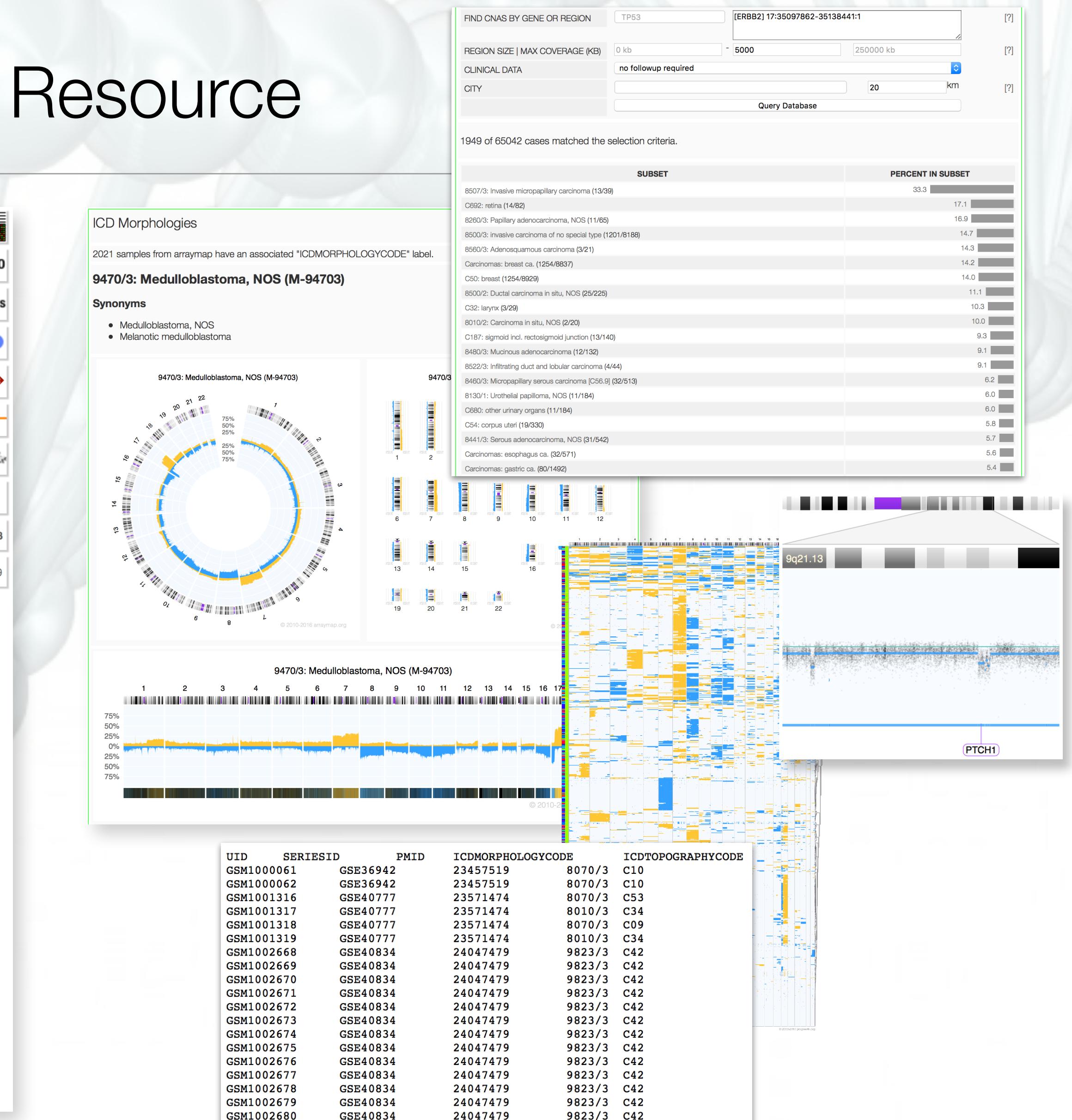
For the majority of the samples, probe level visualization as well as customized data representation facilitate gene level and genome wide data review. Results from multi-case selections can be connected to downstream data analysis and visualization tools, as we provide through our Progenetix project.

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

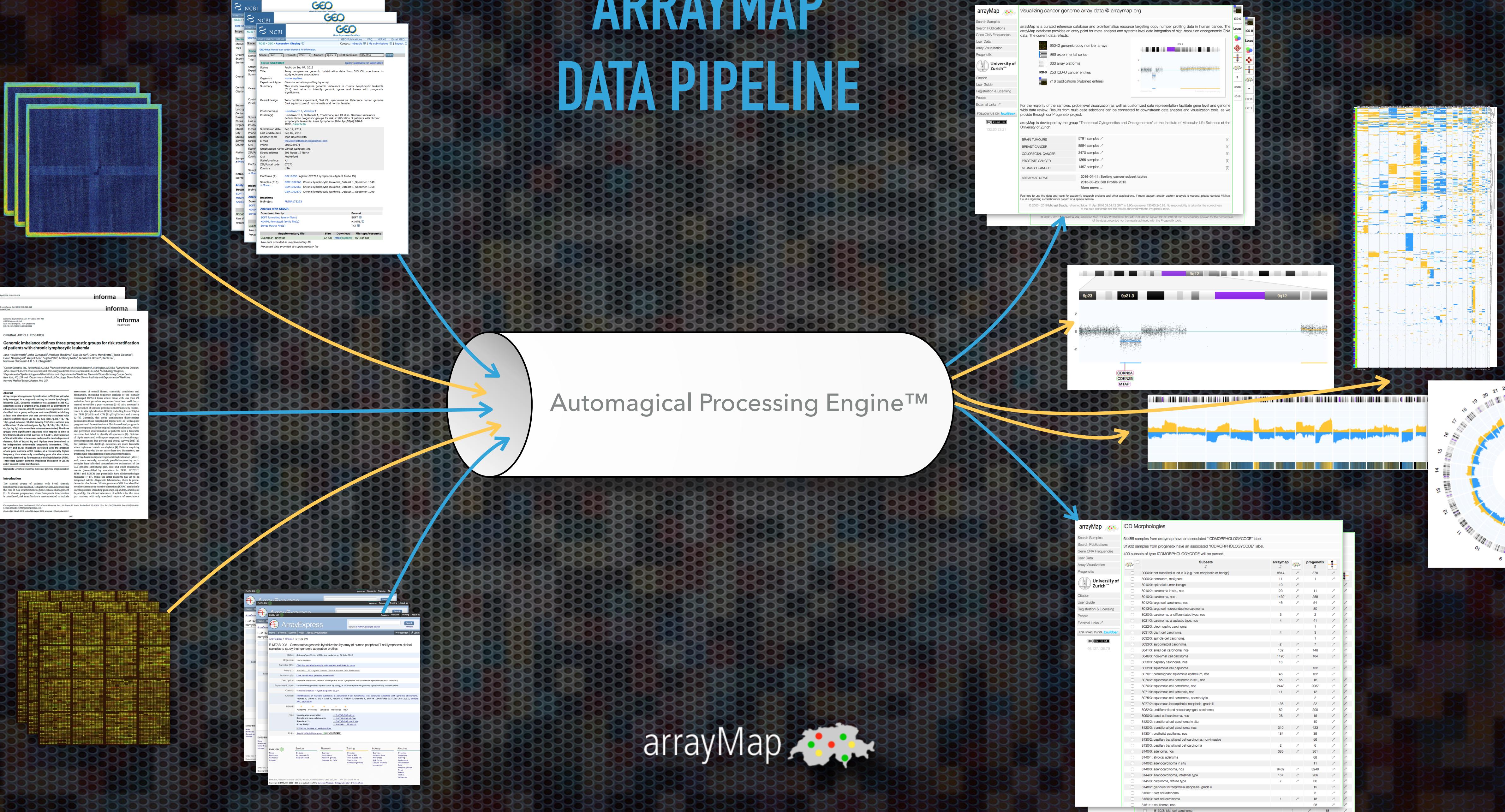
BRAIN TUMOURS	5653 samples ↗	[?]
BREAST CANCER	8329 samples ↗	[?]
COLORECTAL CANCER	3238 samples ↗	[?]
PROSTATE CANCER	991 samples ↗	[?]
STOMACH CANCER	1062 samples ↗	[?]
ARRAYMAP NEWS		
2016-08-03: SVG graphics		
2016-05-17: Transitioning to Europe PMC		
More news ...		

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.



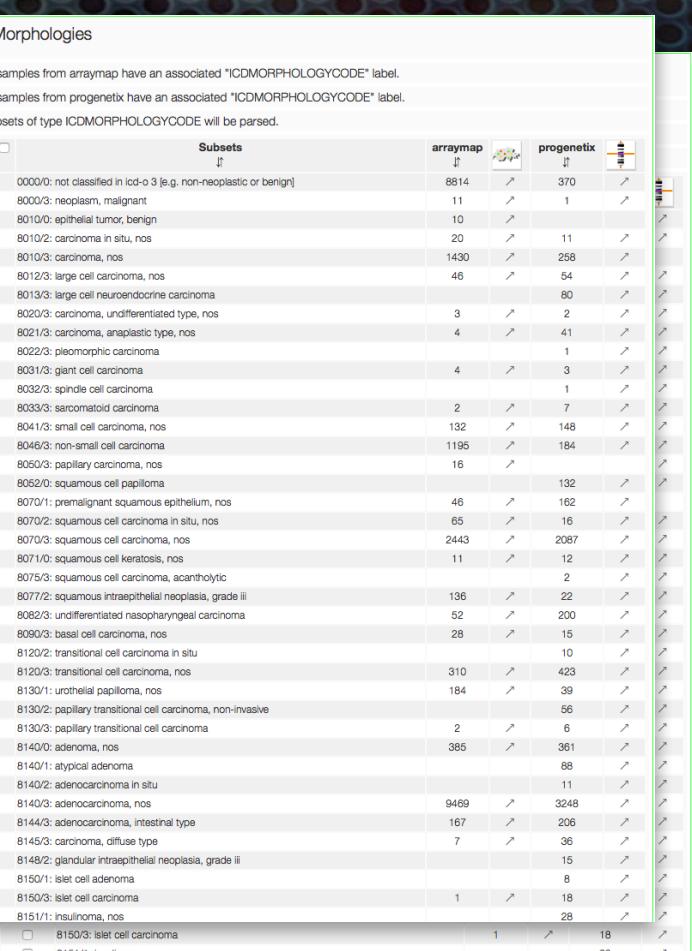
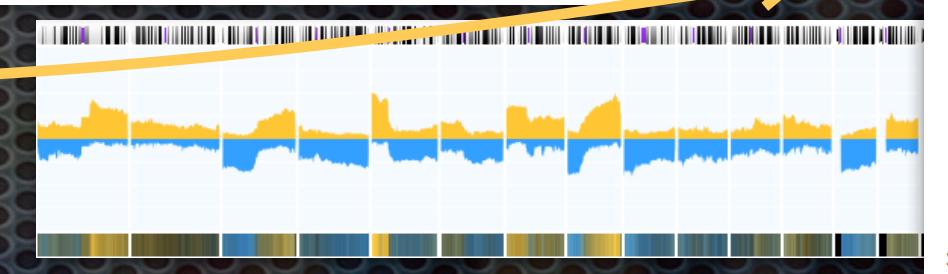
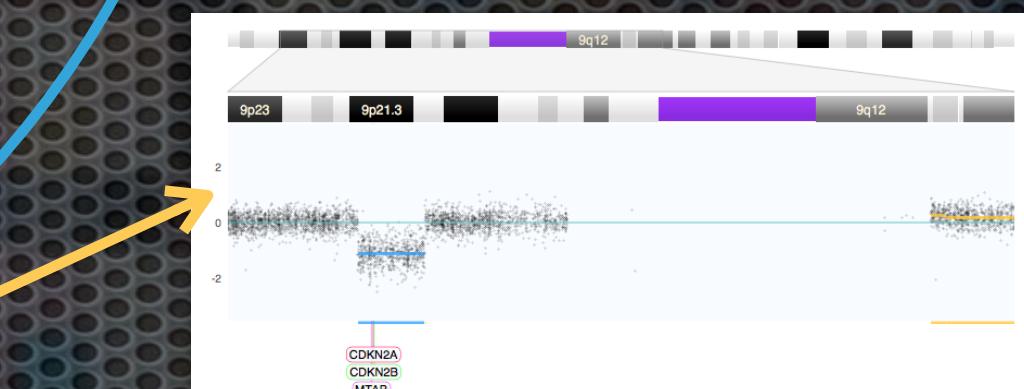
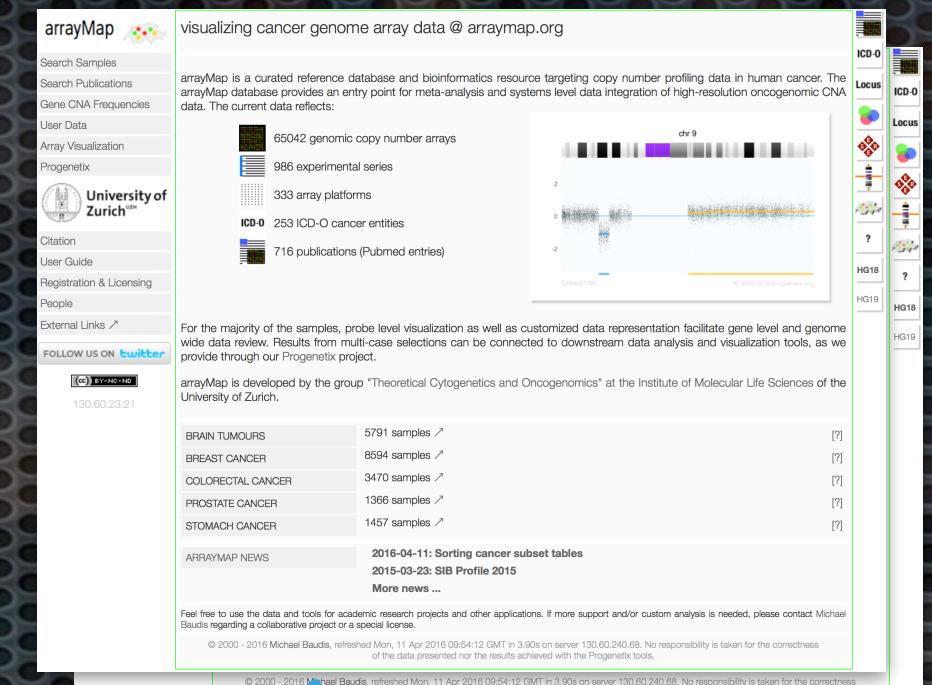
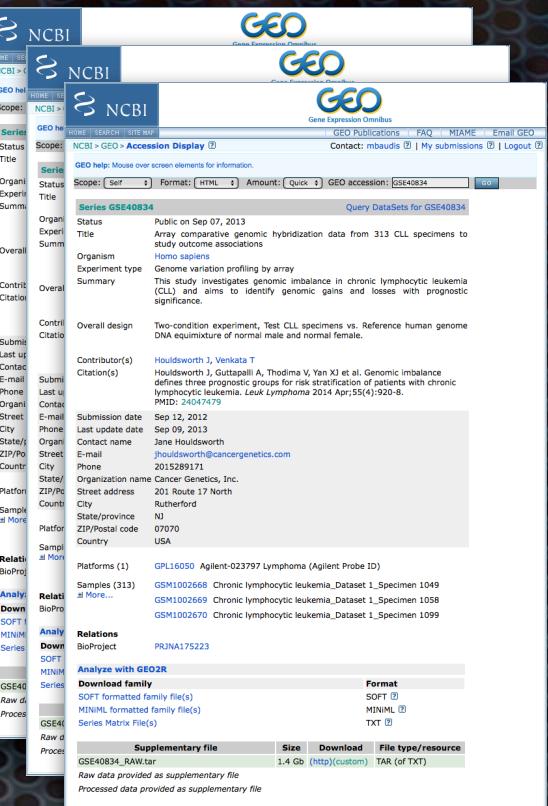
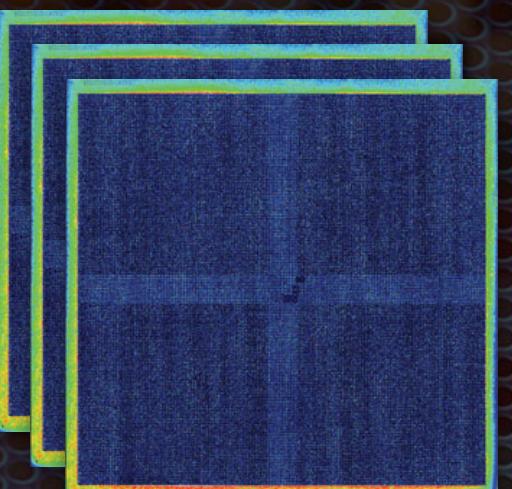
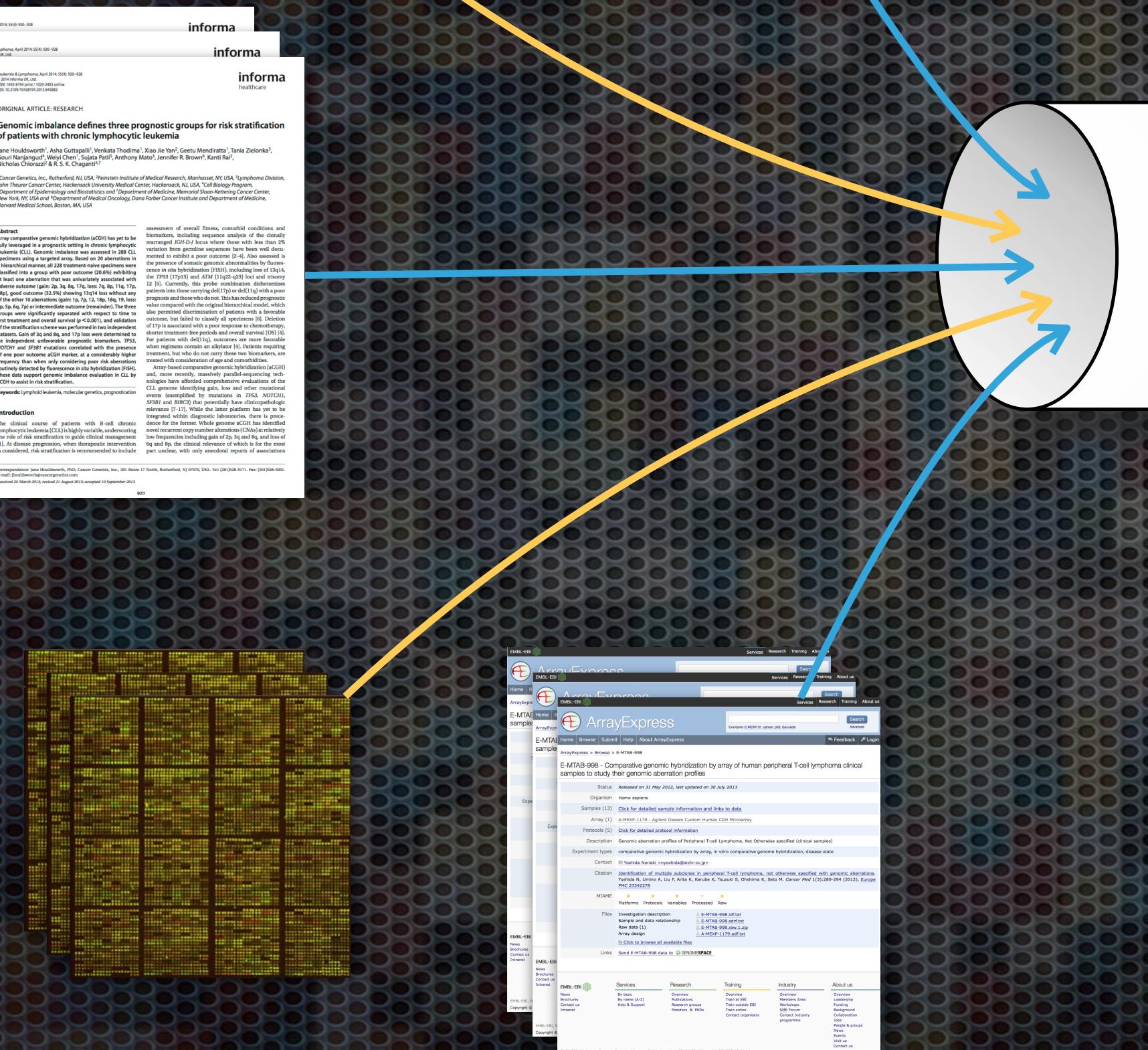
ARRAYMAP DATA PIPELINE



ARRAYMAP DATA PIPELINE

BIOCURATION

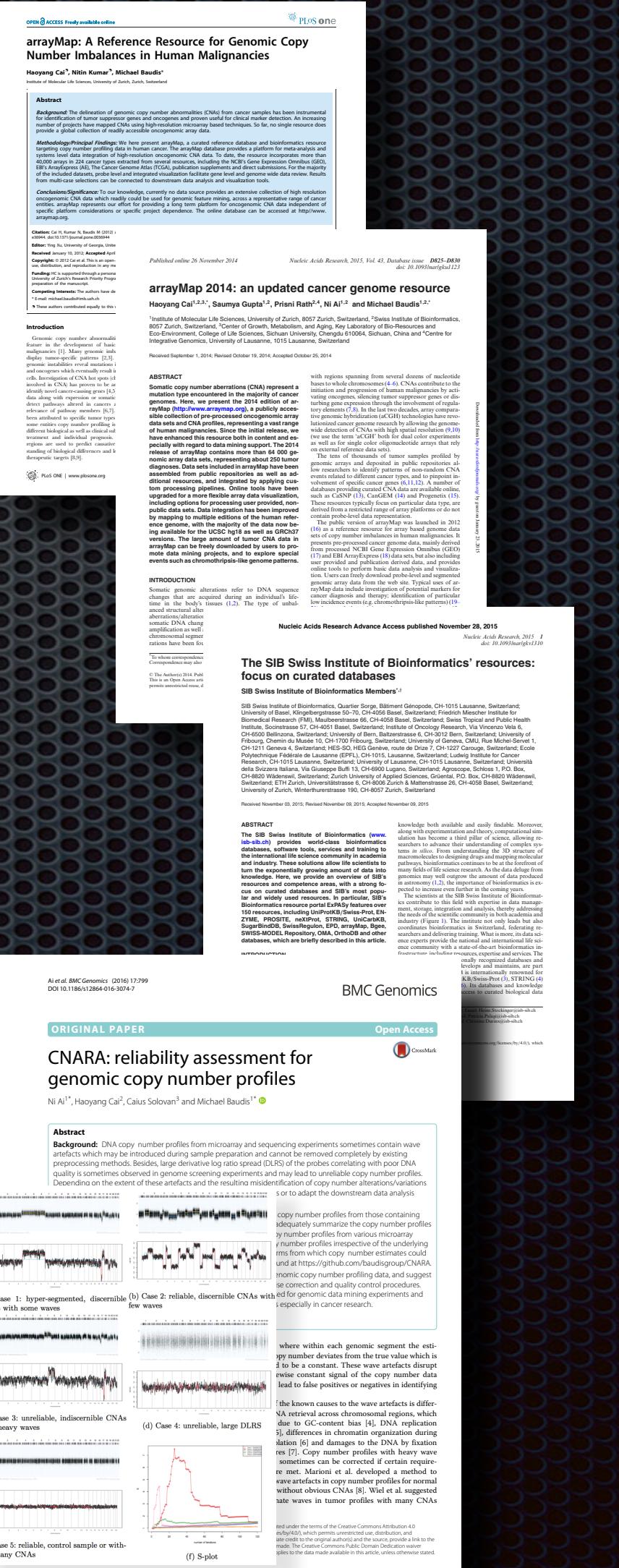
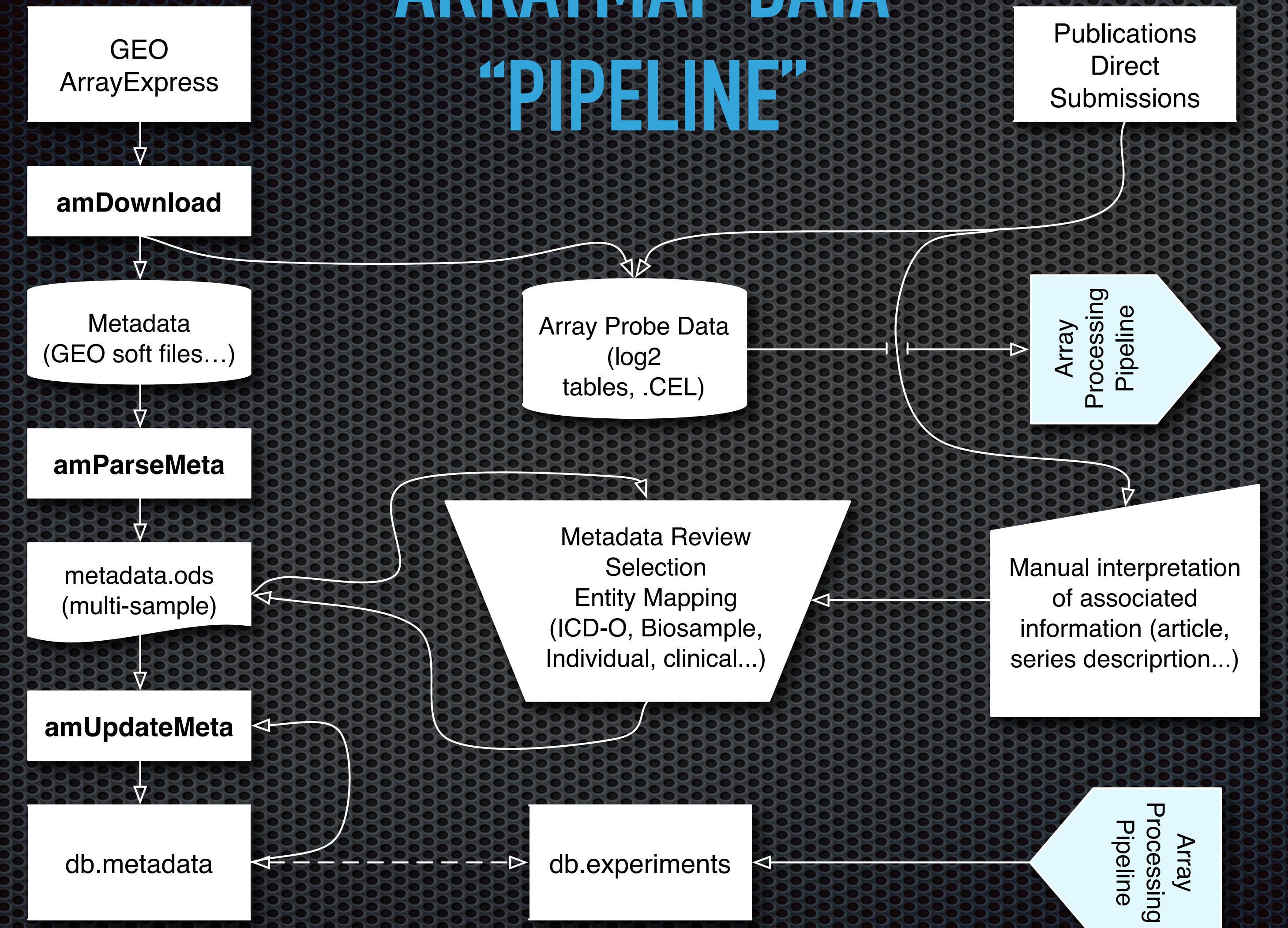
BIOINFORMATICS



arrayMap

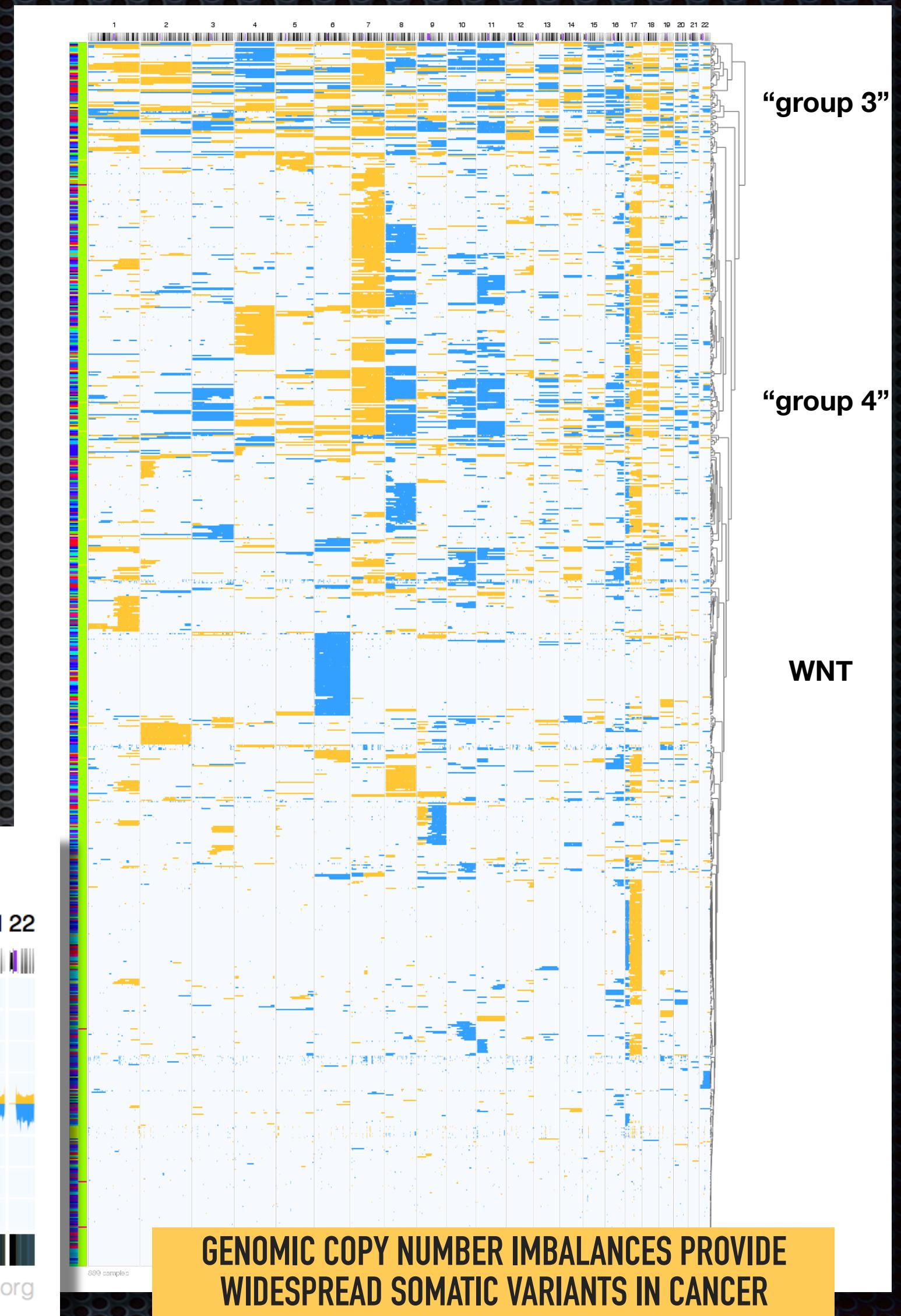
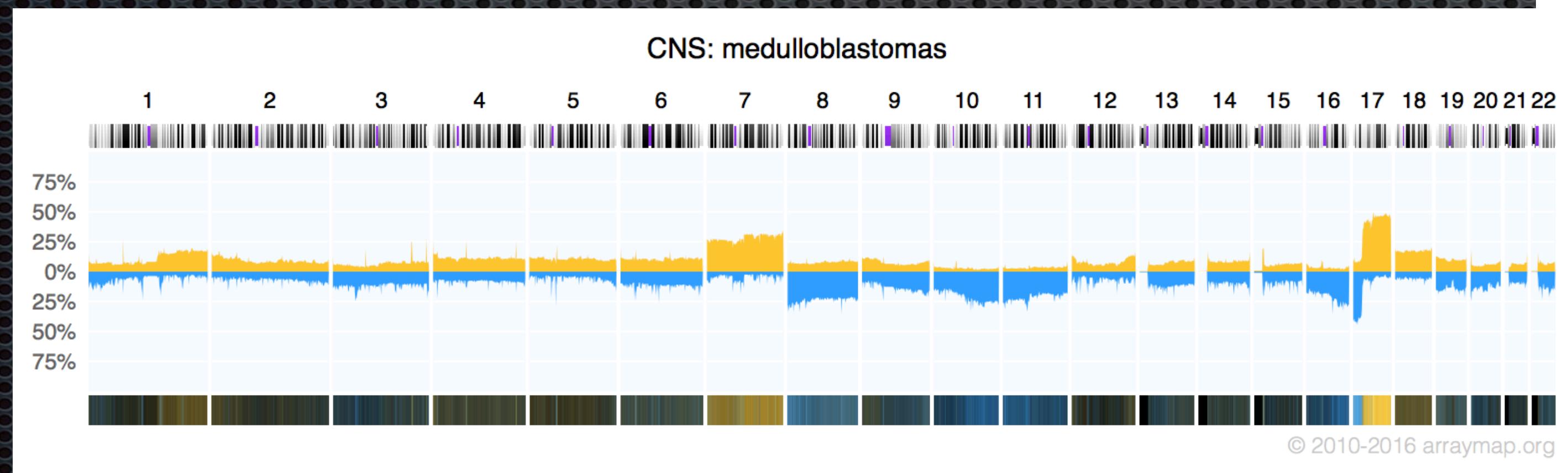
BIOINFORMATICS & CURATION

ARRAYMAP DATA “PIPELINE”



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?

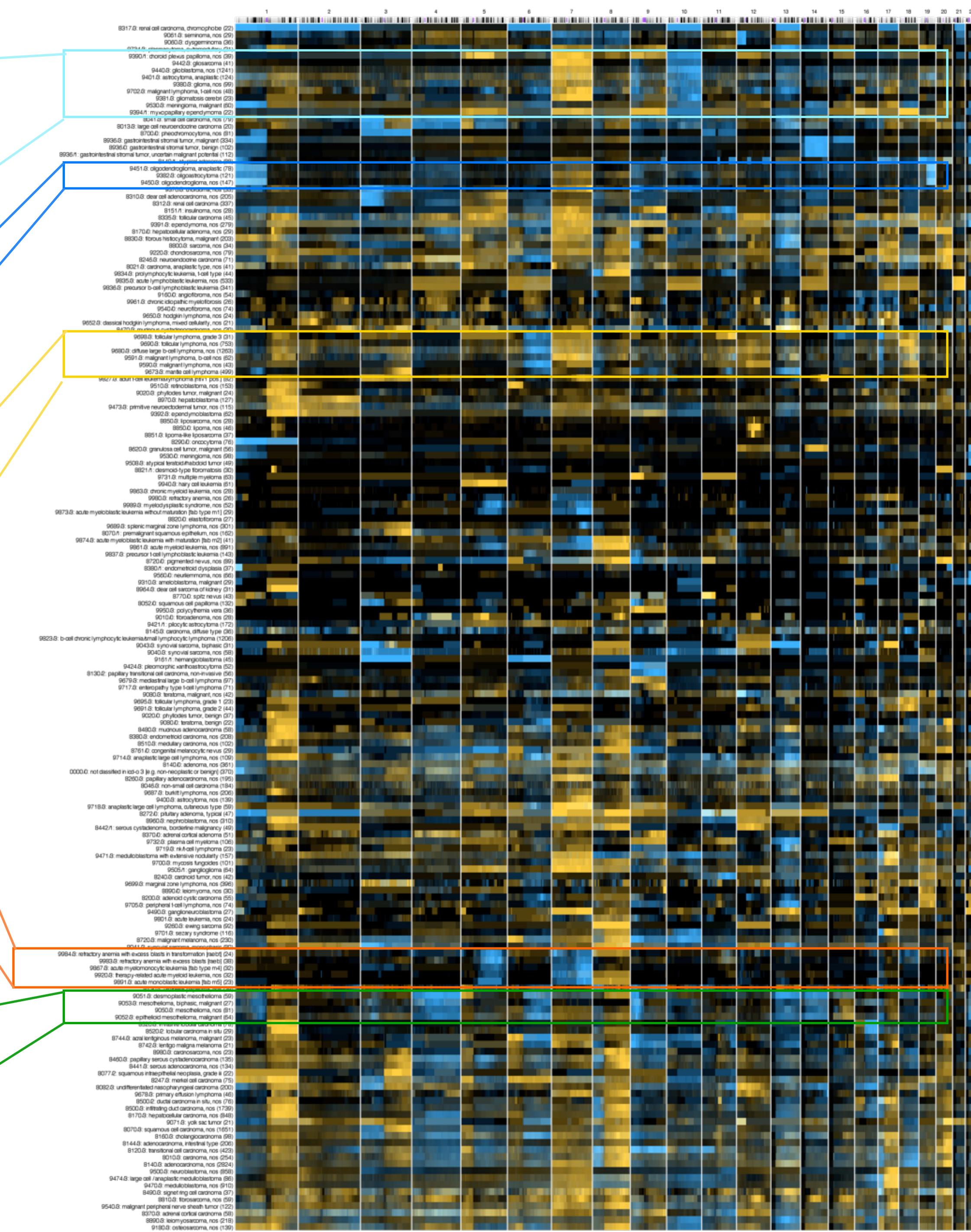
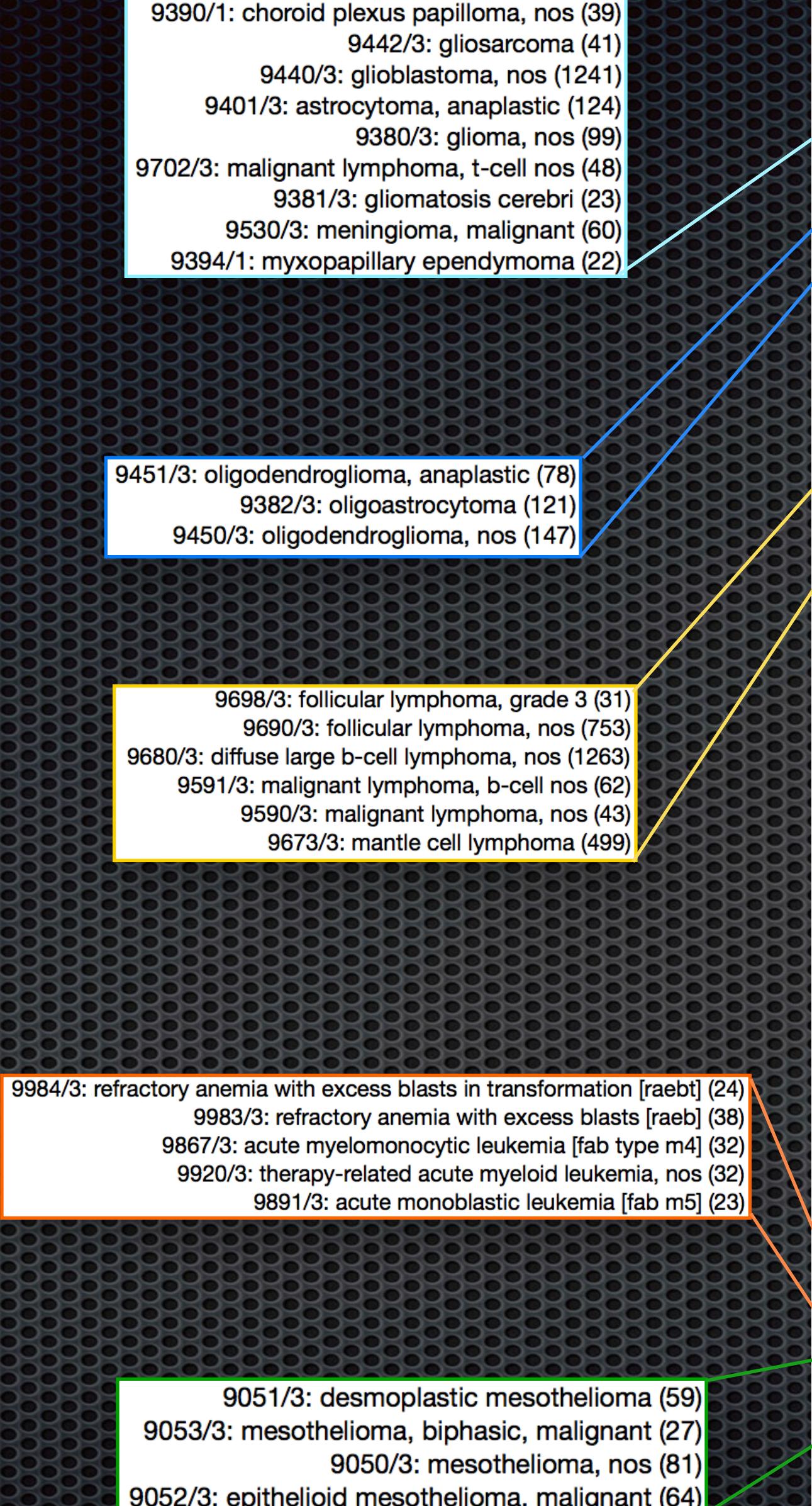


A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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

Somatic Mutations In Cancer: Patterns III

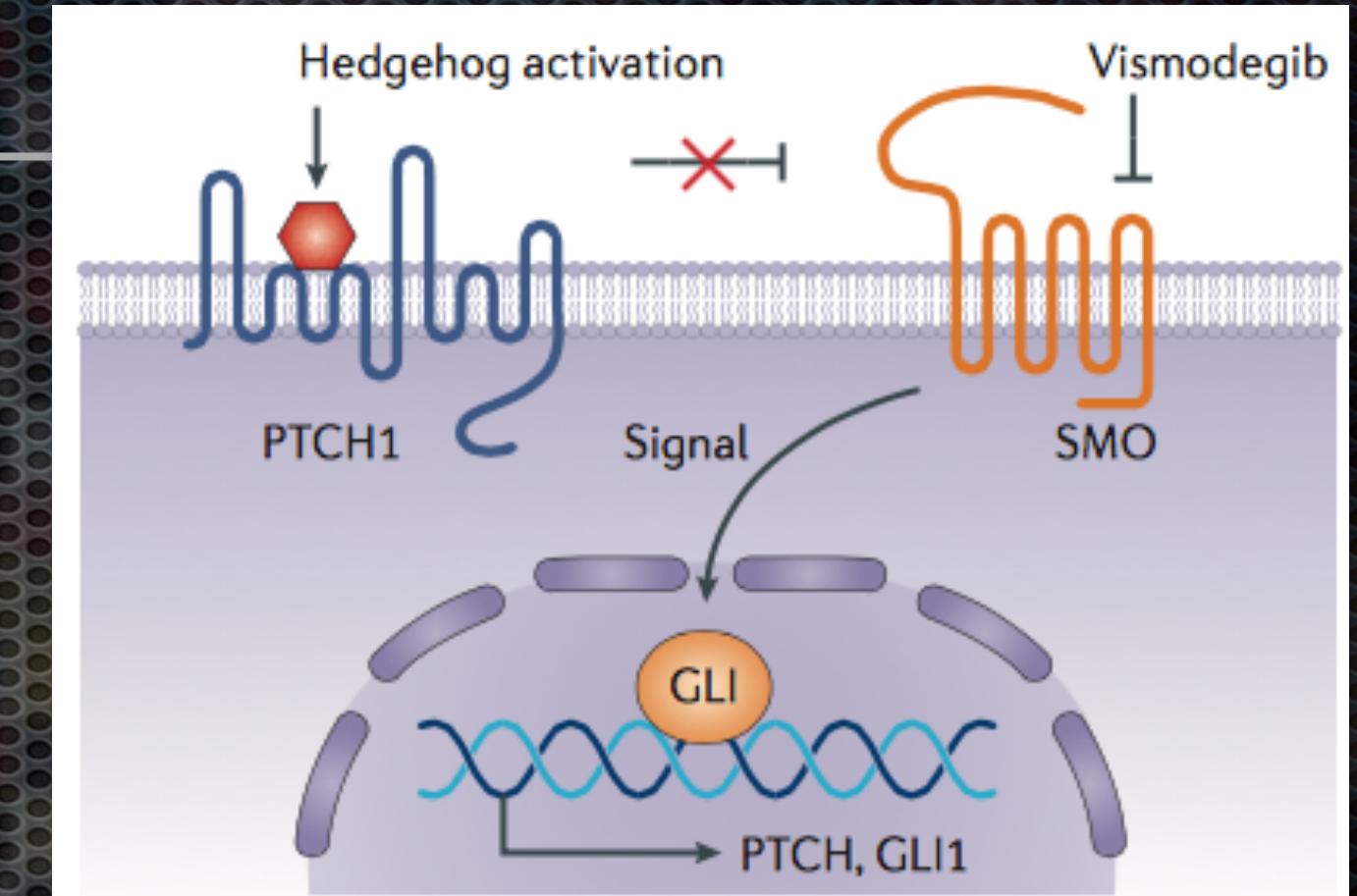
Making the case for genomic classifications

Some related cancer entities show similar copy number profiles

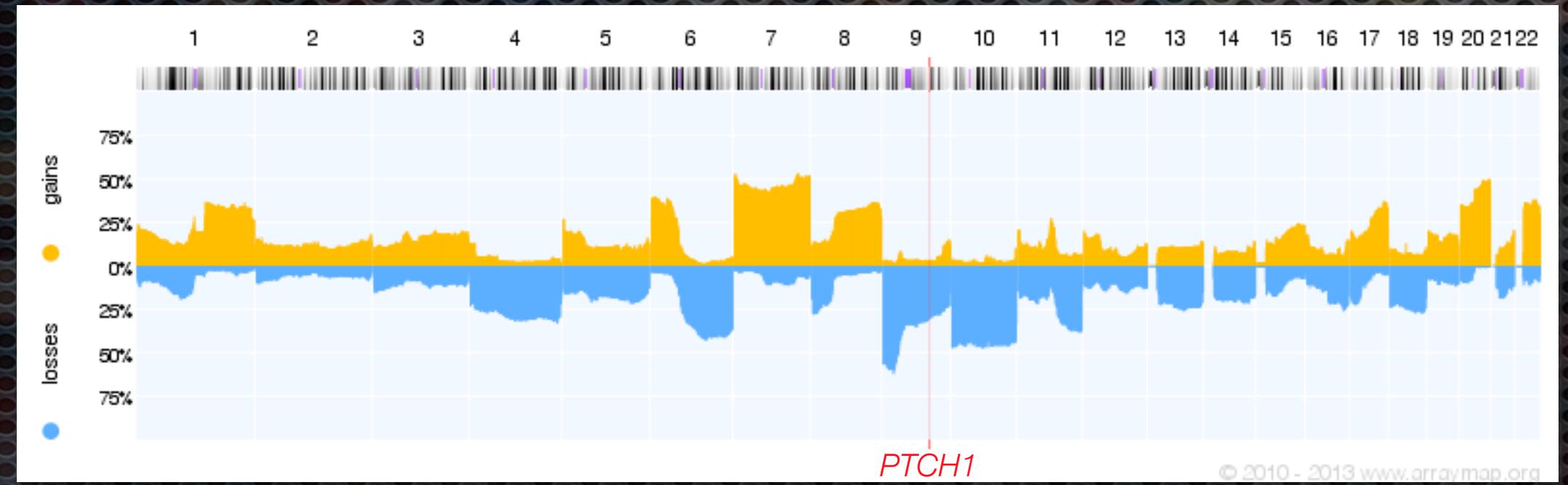


Large datasets for rare cancer gene events

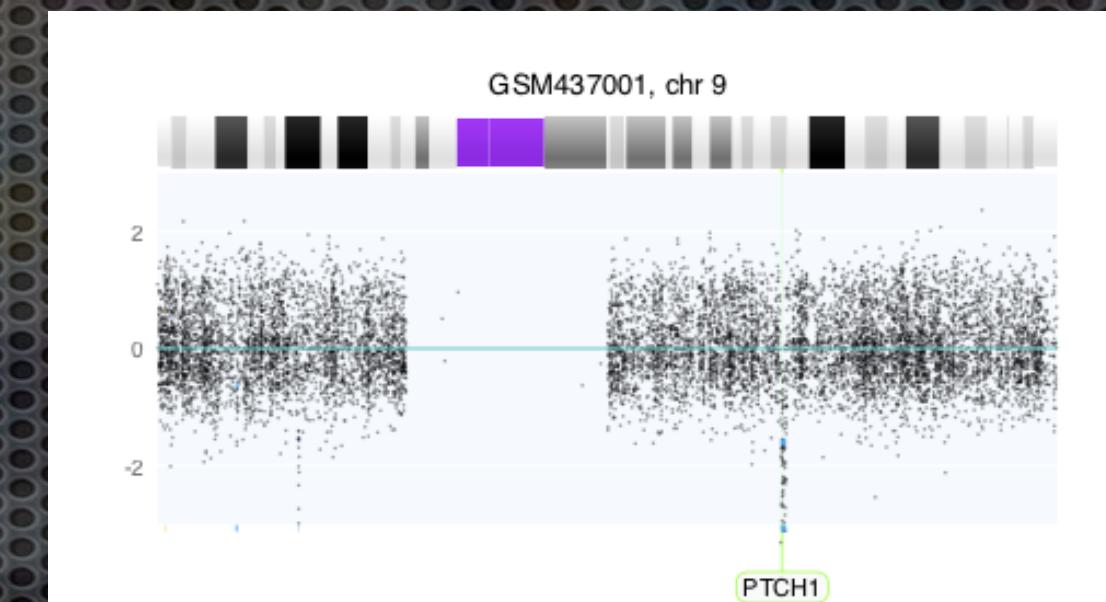
- The Sonic Hedgehog (SHH) pathway has become a “druggable” target in the therapy of syndromic and/or advanced basalomas (e.g. in Gorlin syndrome).
- In the pathway, PTCH1 acts as “tumor suppressor” counteracting SMO=>GLI mediated transcriptional activation.
- We were interested if the gene also could be involved in subsets of malignant melanomas ...



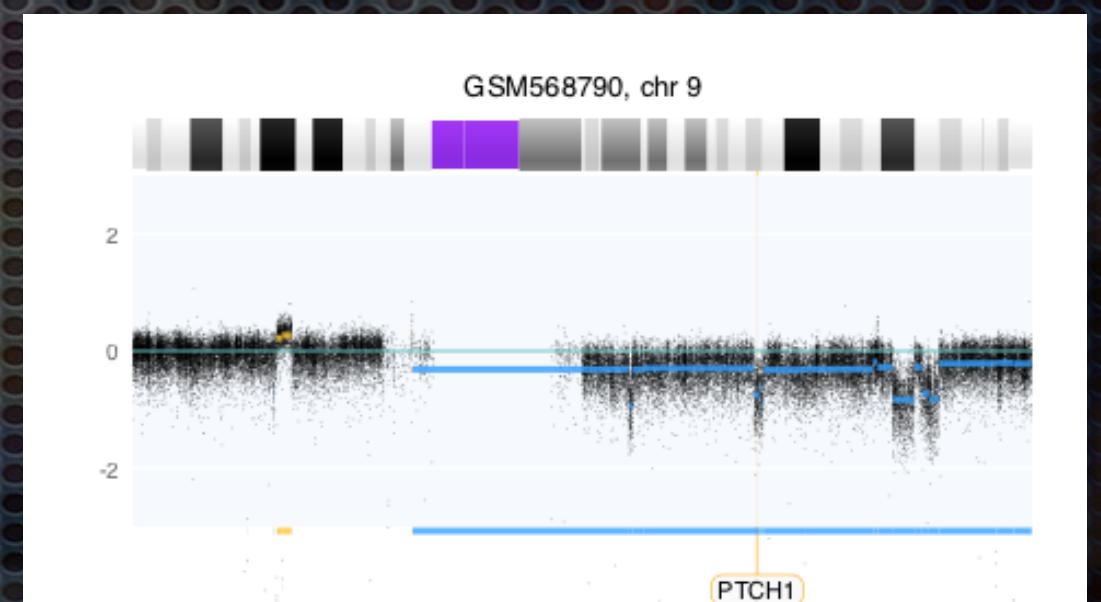
Dlugosz, A., Agrawal, S., & Kirkpatrick, P. (2012, June). Nature Reviews Drug Discovery, pp. 437–438



no “hot spot” (but 30% deletions)

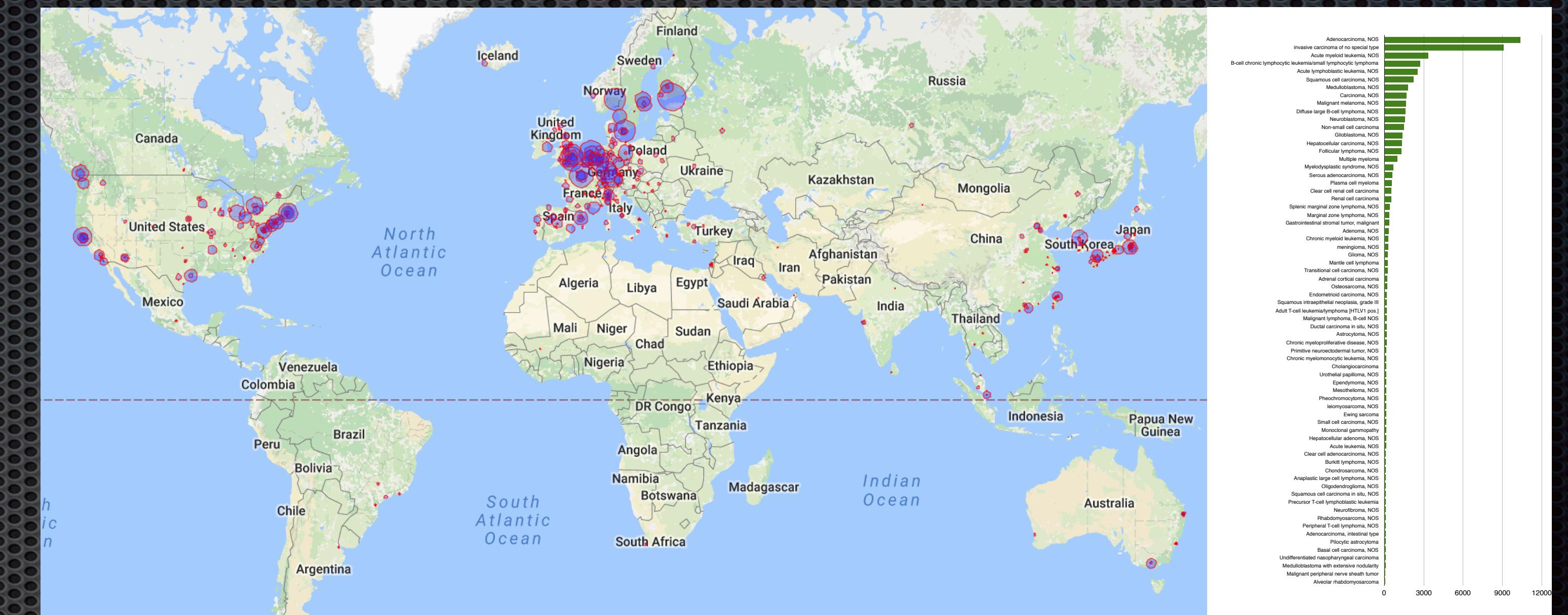


probably pathogenic homozygous deletions in few cases
(3/~700): large datasets needed



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



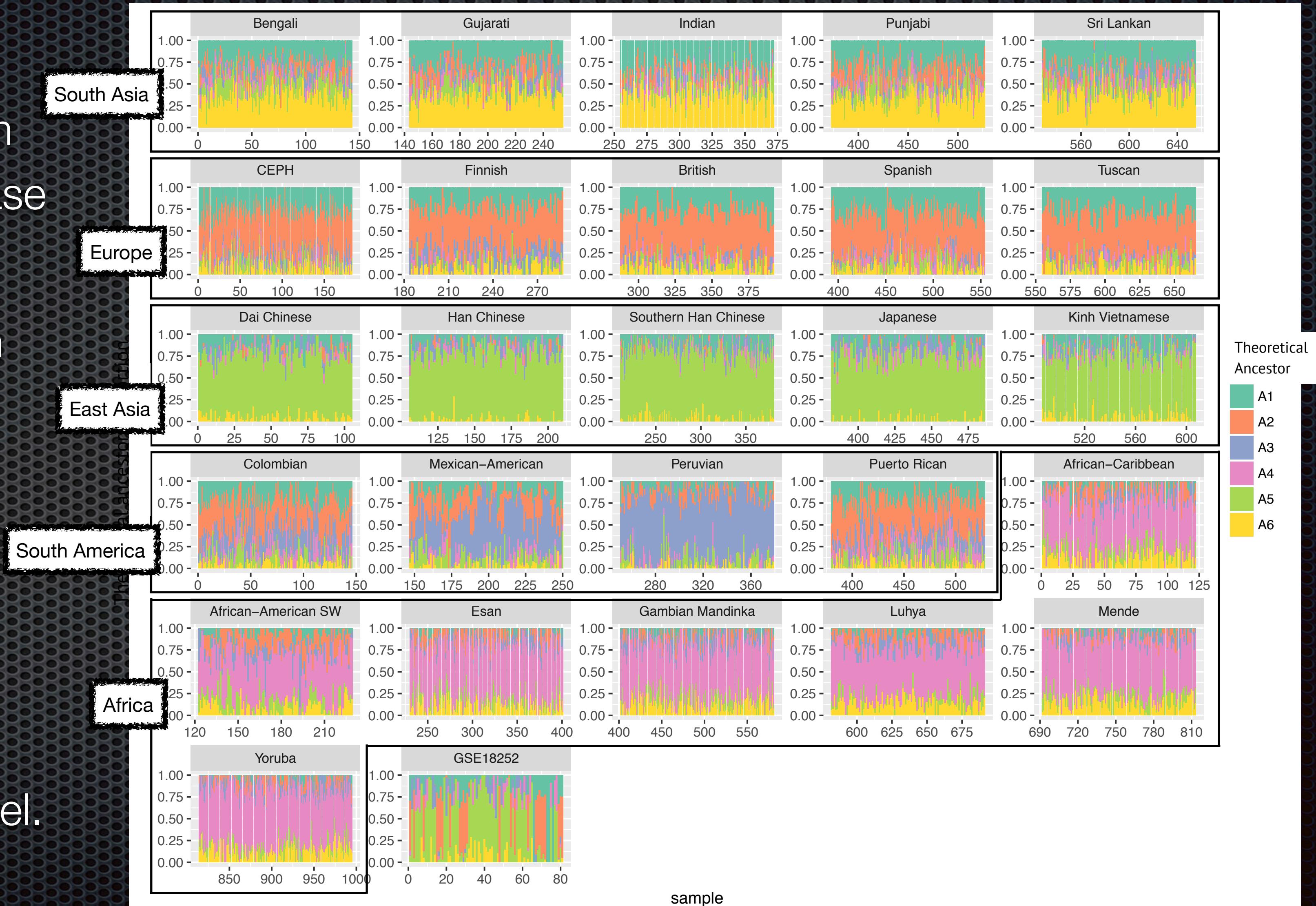
Geographic distribution of >140'000 cancer genome profiles reported in the literature. The numbers are derived from the 2947 publications registered in the Progenetix database.

Population stratification in cancer samples based on SNP array data

- 2504 genome profiles from 1000 Genome project phase 1 as reference

- 5 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.



GA4GH to solve accessibility...



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**



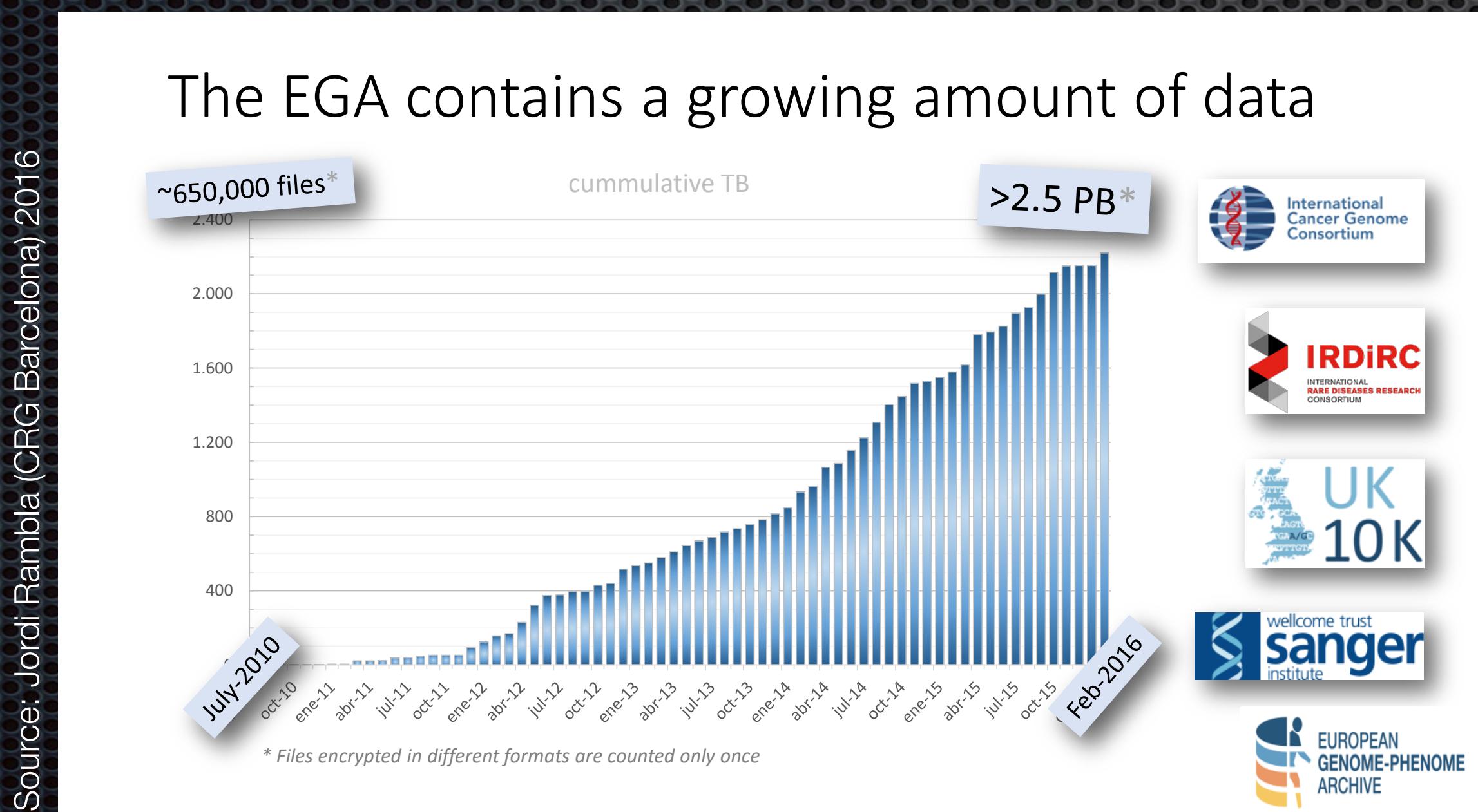
[VIEW OUR LEADERSHIP](#)

[MORE ABOUT US](#)

[BECOME A MEMBER](#)

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 API promotes sharing

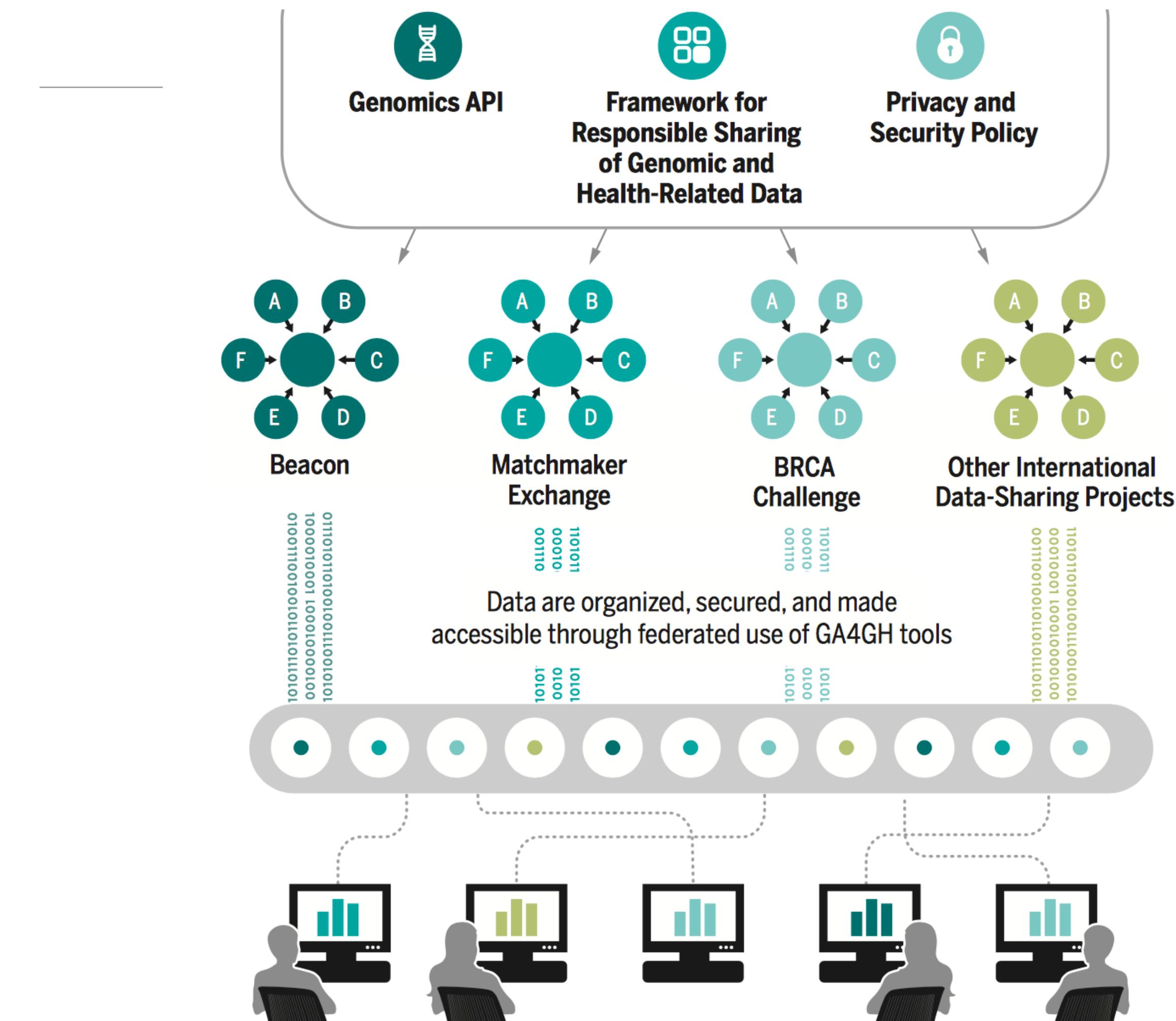
A federated data ecosystem. To share genomic data globally, this approach furthers medical research without requiring compatible data sets or compromising patient identity.



GENOMICS

A federated ecosystem for sharing genomic, clinical data

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



- 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): new structure

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

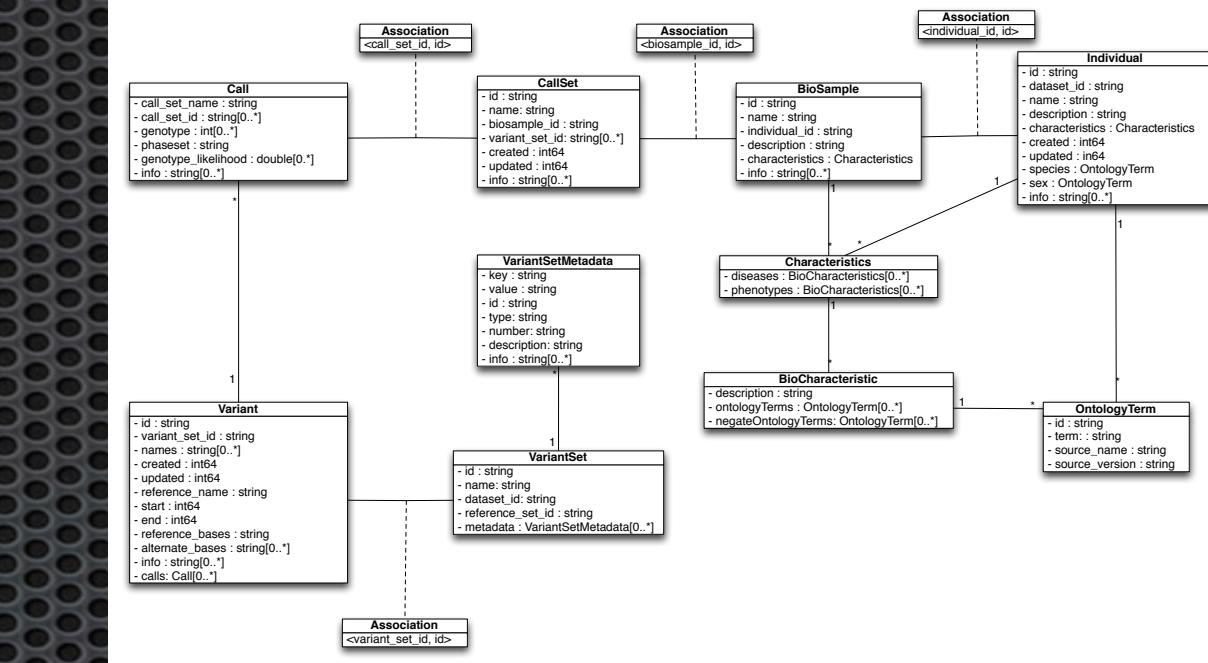


Global Alliance
for Genomics & Health

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 Project

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



Beacon Network

Search Beacons

A global search engine for genetic mutations.

GRCh37 ▾ e.g. 1: 100,000 A>C Search

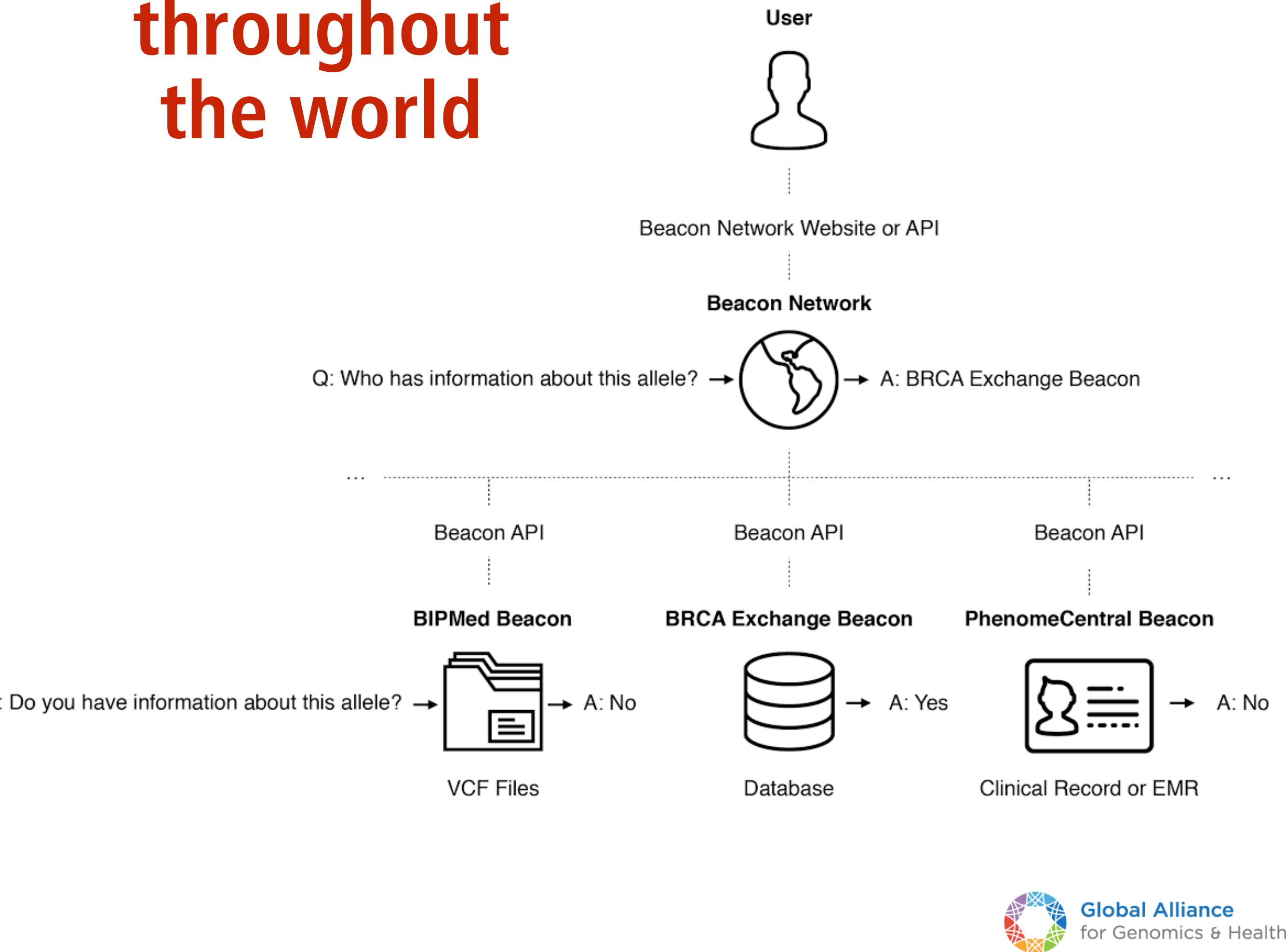
Quickstart: Search for a BRCA2 variant

Find genetic mutations shared by these organizations

- Global Gene Corp
- BRCA EXCHANGE
- Google
- BIPMed Beacon
- PC
- PhenomeCentral Beacon
- Clinical Record or EMR

Browse Beacons »

> 50 Beacons throughout the world



Beacon+ Concept

- Implementation of cancer beacon prototype, backed by arrayMap and DIPG data set
(MacKay *et al.*, Cancer Cell 2017, in print)
- structural variations (DUP, DEL) in addition to SNV
- diagnosis queries using ontology codes (NCIT, ICD-O)
- quantitative responses
- current version uses **GA4GH schema compatible** database

Beacon+

This forward looking Beacon interface implements additional, planned features beyond the current GA4GH specifications. [Info](#)

Query

Dataset: DIPG (CNV + selected SNV)

Reference name*: 17

Genome Assembly*: GRCh36 / hg18

Variant type*: SNV / indel

Position*: 7577121

Ref. Base(s)*: G

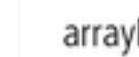
Alt. Base(s)*: A

Bio-ontology: pgx:icdom:9380_3

[Beacon Query](#)

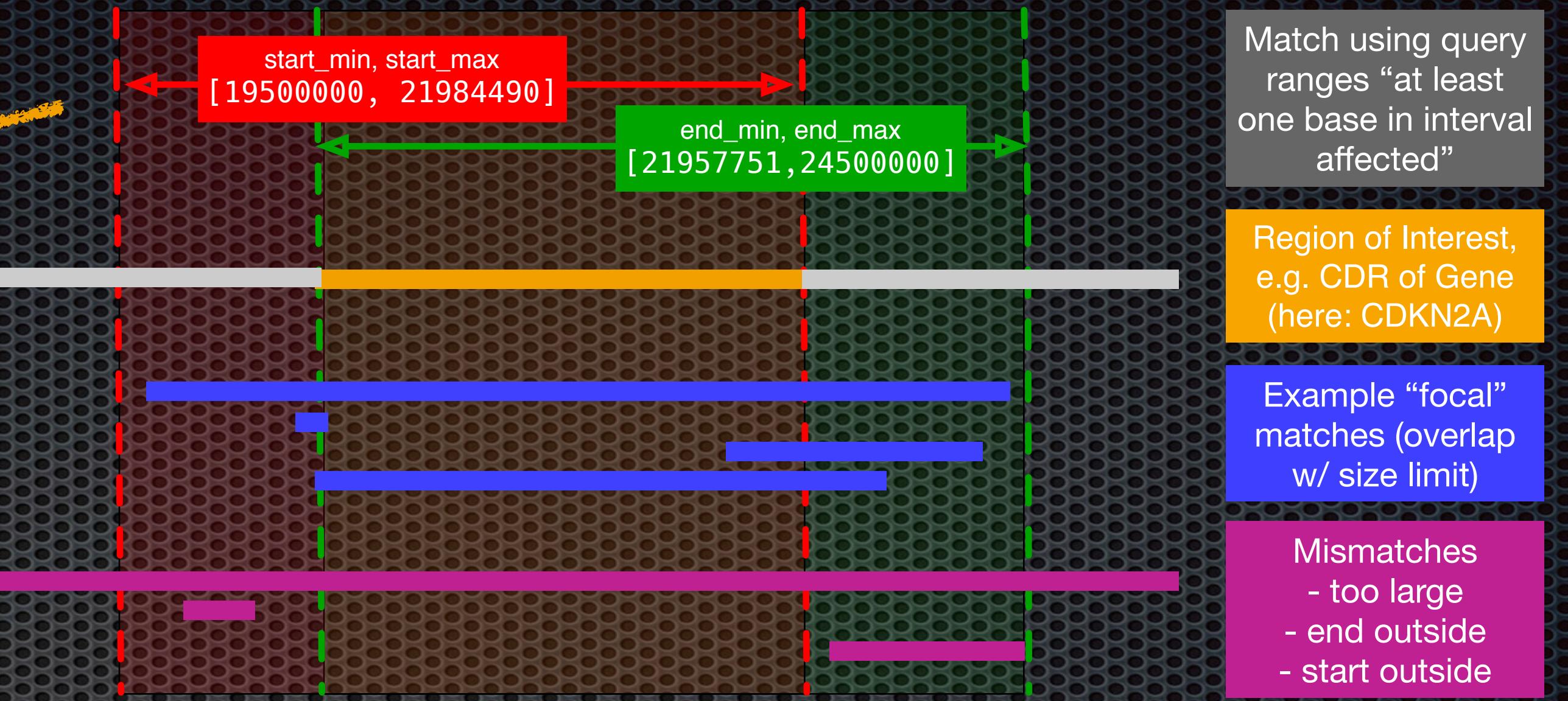
Response

Dataset	Chro.	Assembly	Var. Type	Start Min	Start Max	End Min	End Max	Pos.	Ref.	Alt.	Bio Query	Call Count	Samples	f	Query
arraymap	9	GRCh36	DEL	19000000	21984490	21900000	25000000				pgx:icdom:8140_3	3781	403	0.0065	show JSON
dipg	17	GRCh36	SNV			7577121		G	A	pgx:icdom:9380_3	21	20	0.0187	show JSON	

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

```
{
  "allele_request" : {
    "$and" : [
      { "reference_name" : "9" },
      { "variant_type" : "DEL" },
      { "start" : { "$gte" : 19500000 } },
      { "start" : { "$lte" : 21984490 } },
      { "end" : { "$gte" : 21957751 } },
      { "end" : { "$lte" : 24500000 } }
    ]
  },
  "api_version" : "0.4",
  "beacon_id" : "org.progenetix:progenetix-beacon",
  "exists" : true,
  "info" : {
    "url" : "http://progenetix.org/beacon/info/",
    "dataset_allele_responses" : [
      {
        "dataset_id" : "arraymap",
        "error" : null,
        "exists" : true,
        "external_url" : "http://arraymap.org",
        "sample_count" : 584,
        "call_count" : 3781,
        "variant_count" : 3244,
        "frequency" : 0.0094,
        "info" : {
          "description" : "The query was against database \"arraymap_ga4gh\", variant collection \"variants_cnv_grch36\". 3781 / 59428 matched callsets for 3602919 variants. Out of 62105 biosamples in the database, 2047 matched the biosample query; of those, 584 had the variant."
        },
        "ontology_ids" : [
          "ncit:C3058",
          "pgx:icdom:9440_3",
          "pgx:icdot:C71.9",
          "pgx:icdot:C71.0"
        ]
      }
    ]
  }
}
```

Metadata



- 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



Bioinformatics: **Ontologies**

- ontologies in information sciences describe concrete and abstract **objects**, there precisely defined **hierarchies** and **relationships**
- ontologies in bioinformatics support the move from a descriptive towards an **analytical science** in describing biological data and relations among it

"The widest use of ontologies within biology is for conceptual annotation – a representation of stored knowledge more computationally amenable than natural language."*

- Gene ontology (GO)
- NCI Neoplasm Core
- Uberon anatomical structures
- Experimental Factor Ontology (EFO)
- Disease Ontology (DO)



```
id: GO:0000118
name: histone deacetylase complex
namespace: cellular_component
def: "A protein complex that possesses histone deacetylase activity." [GOC:mah]
comment: Note that this term represents a location, not a function; the activity possessed by this complex is mentioned in the definition for the purpose of describing and distinguishing the complex. The function of this complex is represented by the molecular function term 'histone deacetylase activity'.
synonym: "HDAC complex" EXACT [
is_a: GO:0044451 ! nucleoplasm ]
is_a: GO:1902494 ! catalytic complex
```

- □ Neoplasm by Morphology
 - □ Epithelial Neoplasm [C3709](#)
 - □ Germ Cell Tumor [C3708](#)
 - □ Giant Cell Neoplasm [C7069](#)
 - □ Hematopoietic and Lymphoid Cell Neoplasm [C27134](#)
 - □ Melanocytic Neoplasm [C7058](#)
 - □ Benign Melanocytic Skin Nevus [C7571](#)
 - □ Dysplastic Nevus [C3694](#)
 - □ Melanoma [C3224](#)
 - □ Amelanotic Melanoma [C3802](#)
 - □ Cutaneous Melanoma [C3510](#)
 - □ Epithelioid Cell Melanoma [C4236](#)
 - □ Mixed Epithelioid and Spindle Cell Melanoma [C66756](#)
 - □ Non-Cutaneous Melanoma [C8711](#)
 - □ Spindle Cell Melanoma [C4237](#)
 - □ Meningothelial Cell Neoplasm [C6971](#)

ONTOLOGIES ARE RARELY USED IN CASE REPORTING

- ▶ Medical practice relies on established, slow moving classification systems.
- ▶ Medical diagnoses consist of an abundance of observations and classification items.
- ▶ We do not have (never will?) enough ontology concepts for detailed disease descriptions (Where to stop?)
- ▶ Relationships may help - but how to do them uniformly?

ONTOLOGIES ARE RARELY USED IN CASE REPORTING

- ▶ Medical practice relies on established, slow-moving classification systems.
- ▶ Medical diagnoses consist of a abundance of observations and classification items
- ▶ We do not have (or ever will?) enough ontology concepts for detailed case descriptions (Where to stop?)
- ▶ Relationships may help - but how to do them uniformly?

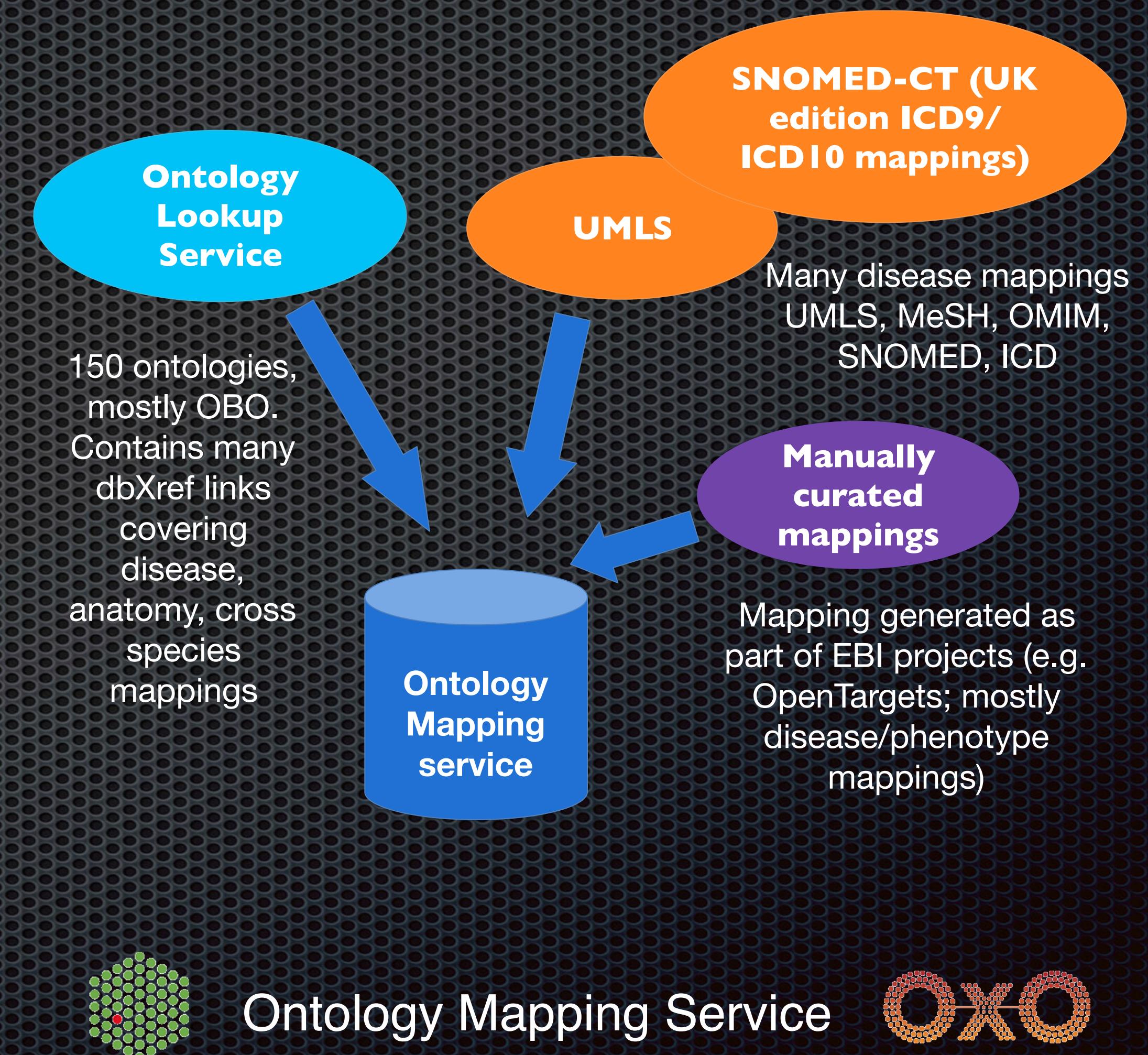
LHE
the core

THE GOLD STANDARD FOR A MEDICAL DIAGNOSIS IS STILL WELL WRITTEN PROSE.

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



Working towards ontologies w/ arrayMap: Mapping >55'000 samples from ICD-O to NCIt neoplasm core

ICDM	ICDMORPHOLOGY
8021/3	Carcinoma anaplastic type
9451/3	Oligodendrogloma anaplastic
9051/3	Desmoplastic mesothelioma
9732/3	Plasma cell myeloma
8070/3	Squamous cell carcinoma
8380/3	Endometrioid adenocarcinoma
8070/3	Squamous cell carcinoma
8430/3	Mucoepidermoid carcinoma
9680/3	Diffuse large B-cell lymphoma
8800/3	Sarcoma
8441/3	Serous adenocarcinoma
9689/3	splenic marginal zone lymphoma nos
8077/2	Squamous intraepithelial neoplasia grade III
8140/0	Adenoma
8272/3	Pituitary carcinoma
8500/2	Ductal carcinoma in situ
8200/3	Adenoid cystic carcinoma
9370/3	Chordoma
9717/3	Enteropathy type T-cell lymphoma
9698/3	Follicular lymphoma grade 3
9863/3	Chronic myeloid leukemia
8852/3	Liposarcoma myxoid
9080/3	Teratoma malignant
8530/3	Inflammatory carcinoma
8140/3	Adenocarcinoma
8200/3	Adenoid cystic carcinoma

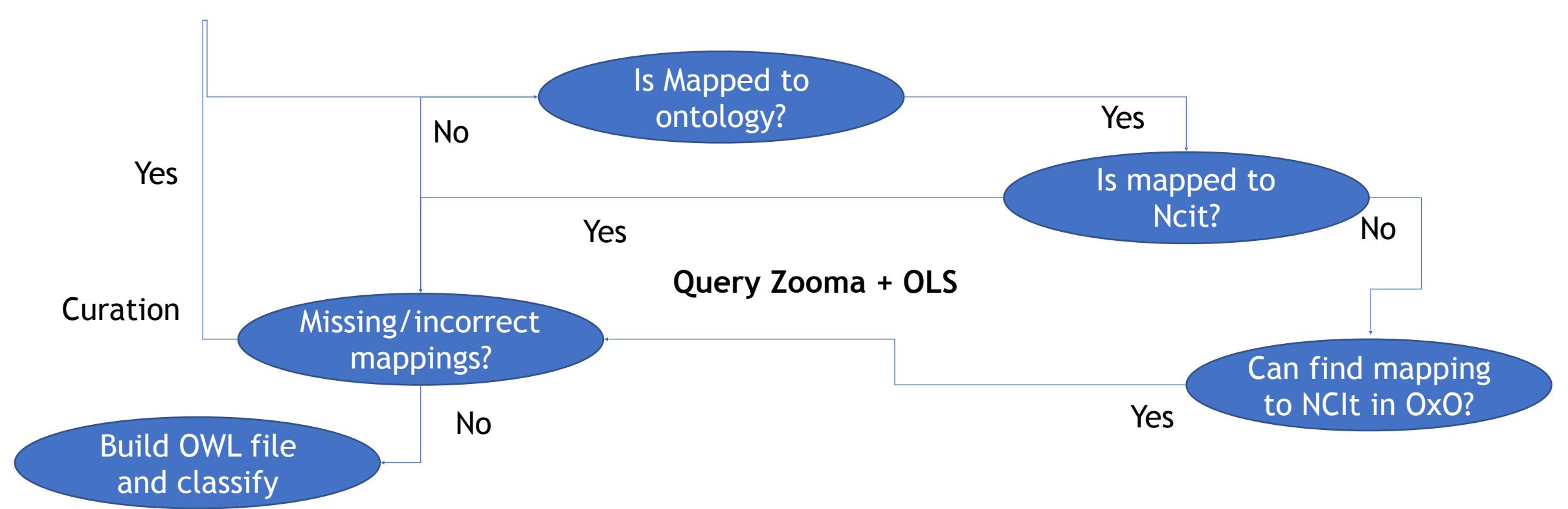
NCItcode	NCItlabel
C4326	anaplastic oligodendrogloma
C6747	
C3242	multiple myeloma
C2926	non-small cell lung carcinoma
C3769	endometrioid carcinoma
C2926	non-small cell lung carcinoma
C45544	pulmonary mucoepidermoid carcinoma
C8851	diffuse large B-cell lymphoma
C9118	sarcoma
C7550	ovarian serous adenocarcinoma
C4196	adenoma
C4536	Pituitary carcinoma
C3641	ductal carcinoma in situ
C2970	adenoid cystic carcinoma
C2947	Chordoma
C3177	chronic myelogenous leukemia
C3735	myxoid liposarcoma
C4872	breast carcinoma
C27745	lung adenocarcinoma
C2670	

ICDT	ICDTOPOGRAPHY
C739	thyroid gland
C719	Brain
C499	connective and soft tissue
C42	hematopoietic and reticuloendothelial systems
C140	pharynx
C54	corpus uteri
C44	skin
C089	salivary gland
C42	hematopoietic and reticuloendothelial systems
C559	uterus nos
C570	fallopian tube
C422	spleen
C53	cervix uteri
C189	large intestine excl. rectum and rectosigmoid junction
C751	pituitary gland
C50	breast
C32	larynx
C419	bone
C17	small intestine
C42	hematopoietic and reticuloendothelial systems
C42	hematopoietic and reticuloendothelial systems
C499	connective and soft tissue
C809	unknown
C50	breast
C809	unknown
C12	uterus nos

NCItcode	NCItlabel
C12400	thyroid gland
C12439	brain
C12316	
C12470	zone of skin
C12426	saliva-secreting gland
C12403	fallopian tube
C12432	spleen
C12311	
C12399	pituitary gland
C12971	breast
C12420	larynx
C13076	bone tissue
C12386	small intestine
C35882	Hereditary elliptocytosis
C12971	breast
C35882	Hereditary elliptocytosis
C12762	oropharynx
C12415	kidney
C12499	internal ear
C12683	bronchus
C12343	retina
C12393	pancreas
C12422	tongue
C12390	rectum
C12404	female gonad
C12391	

NCIt_mapped	NCIt_mapped_ICDM_T_label
C3878	Thyroid Gland Undifferentiated (Anaplastic) Carcinoma
C4326	Anaplastic Oligodendrogloma
C6747	Desmoplastic Mesothelioma
C3242	Plasma Cell Myeloma
C102872	Pharyngeal Squamous Cell Carcinoma
C6287	Endometrial Endometrioid Adenocarcinoma
C4819	Skin Squamous Cell Carcinoma
C5953	Minor Salivary Gland Mucoepidermoid Carcinoma
C8851	Diffuse Large B-Cell Lymphoma
C9306	Soft Tissue Sarcoma
C40101	Serous Adenocarcinoma
C4663	Splenic Marginal Zone Lymphoma
C89476	Grade III Vaginal Intraepithelial Neoplasia
C4349	Colon Adenocarcinoma
C4536	Pituitary Gland Carcinoma
C2924	Ductal Breast Carcinoma In Situ
C2970	Adenoid Cystic Carcinoma
C2947	Chordoma
C4737	Enteropathy-Associated T-Cell Lymphoma
C3460	Grade 3 Follicular Lymphoma
C3174	Chronic Myelogenous Leukemia BCR-ABL1 Positive
C27781	Myxoid Liposarcoma
C3403	Tetrotoma
C4001	Inflammatory Breast Carcinoma
C2852	Adenocarcinoma
C2970	Adenoid Cystic Carcinoma
C3158	Leiomyosarcoma
C2970	Adenoid Cystic Carcinoma
C2923	Bronchioloalveolar Carcinoma
C3224	Melanoma
C8459	Hepatosplenic T-Cell Lymphoma
C8294	Pancreatic Adenocarcinoma
C3996	Monoclonal gammopathy of Undetermined Significance
C4817	Ewing Sarcoma
C3288	Oligodendrogloma
C4648	Tongue Squamous Cell Carcinoma
C2862	Primary Myelofibrosis
C4833	Oral Cavity Squamous Cell Carcinoma
C9383	Rectal Adenocarcinoma
C3158	Leiomyosarcoma
C3898	Extranodal Marginal Zone Lymphoma of Mucosa-Associated Lymphoid Tissue
C4512	Ovarian Mucinous Cystadenoma
C5519	Other

- From 456 pairs of ICD-O terms Morphology and Topography representative of cancer entities in arrayMap
- Develop Python script to take ICD-O Morphology and Topography labels separately QUERY ZOOMA, Oxo and OLS to find mapping to NCIt



From 456 pairs of ICD-O
70% ICD-O Morphology - NCIt
65% ICD-O Topography - NCIt

45% ICD-O-3 Pairs mapped to NCIt terms

=> MANUAL CURATION of >50%

Beacon+ Concept

- Implementation of cancer beacon prototype, backed by arrayMap and DIPG data set
(MacKay *et al.*, Cancer Cell 2017, in print)
- structural variations (DUP, DFI) in addition to SNV
- diagnosis queries using ontology codes (NCIT, ICD-O)
- quantitative responses
- current version uses **GA4GH schema compatible** database

Beacon+

This forward looking Beacon interface implements additional, planned features beyond the current GA4GH specifications. [Info](#)

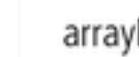
Query

Dataset: DIPG (CNV + selected SNV)
Reference name*: 17
Genome Assembly*: GRCh36 / hg18
Variant type*: SNV / indel
Position*: 7577121
Ref. Base(s)*: G
Alt. Base(s)*: A
Bio-ontology: pgx:icdom:9380_3

Beacon Query

Response

Dataset	Chro.	Assembly	Var. Type	Start Min	Start Max	End Min	End Max	Pos.	Ref.	Alt.	Bio Query	Call Count	Samples	f	Query
arraymap	9	GRCh36	DEL	19000000	21984490	21900000	25000000				pgx:icdom:8140_3	3781	403	0.0065	show JSON
dipg	17	GRCh36	SNV			7577121		G	A	pgx:icdom:9380_3	21	20	0.0187	show JSON	

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



This forward looking Beacon interface implements additional, planned features beyond the current GA4GH specifications.

[Info](#)

Query

[SNV Example](#)

[DGV Example](#)

[CNV Example](#)

Dataset

arraymap

Reference name*

9

Genome Assembly*

GRCh36 / hg18

Variant type*

DEL (Deletion)

Start min Position*

19,500,000

Start max Position

21,964,826

End min Position

21,958,228

End max Position

24,500,000

Bio-ontology

ncit:c3224: Melanoma (1098)

Beacon Query

Response

Dataset	Assembly	Chro	Var Type	Start Range	End Range	Pos	Ref Alt	Bio Query	Variants	Calls	Samples	f_alleles	Response Context
arraymap	hg18	9	DEL	19,500,000 21,964,826	21,958,228 24,500,000			ncit:c3224	157 171 171			0.1557	JSON UCSC Handover

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

Variants	Calls	Samples	f_alleles	Response Context
157				JSON
171				UCSC
171			0.1557	Handover



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
- ```
"dataset_allele_responses" : [{ "dataset_id" : "arraymap", "call_count" : "171", "sample_count" : "171", "variant_count" : "157", "error" : null, "exists" : true, "external_url" : "http://beacon.arraymap.org", "frequency" : "0.1557", "info" : { "callset_access_handle" : "d5850347-d411-11e7-8c89-ec436516cb41", "description" : "The query was against database \"arraymap_ga4gh\", variant collection \"variants_cnv_grch36\". 171 matched calls for 157 distinct variants. Out of 62033 biosamples in the database, 1098 matched the biosample query; of those, 171 had the variant.", }, "note" : "", }],
```





This is an implementation of a Beacon "handover" concept, in which a Beacon query response additionally delivers an "accessid" value. This value represents a pointer to an internal representation of the query results (i.e. callsets, biosamples, metadata ...), which can then be accessed after authentication. The "handover" scenario separates the standard qualitative ("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 variation data
- returning data of the associated callsets which matched the Beacon query (this is for feature demonstration only...)
- 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**

This Beacon implementation is developed by the Computational  
support from the SIB Technology group and ELIXIR.

## 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 not managed by Beacon protocol / system => "Discovery" protocol?





This is an implementation of a Beacon "handover" concept, in which a Beacon query response additionally delivers an "accessid" representation of the query results (i.e. callsets, biosamples, metadata ...), which can then be accessed after authentication. The ("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 variation data
- returning data of the associated callsets which matched the Beacon query (this is for feature demonstration only...)
- 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 hist

Export Callset Dat

Export Biosample

## Login

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

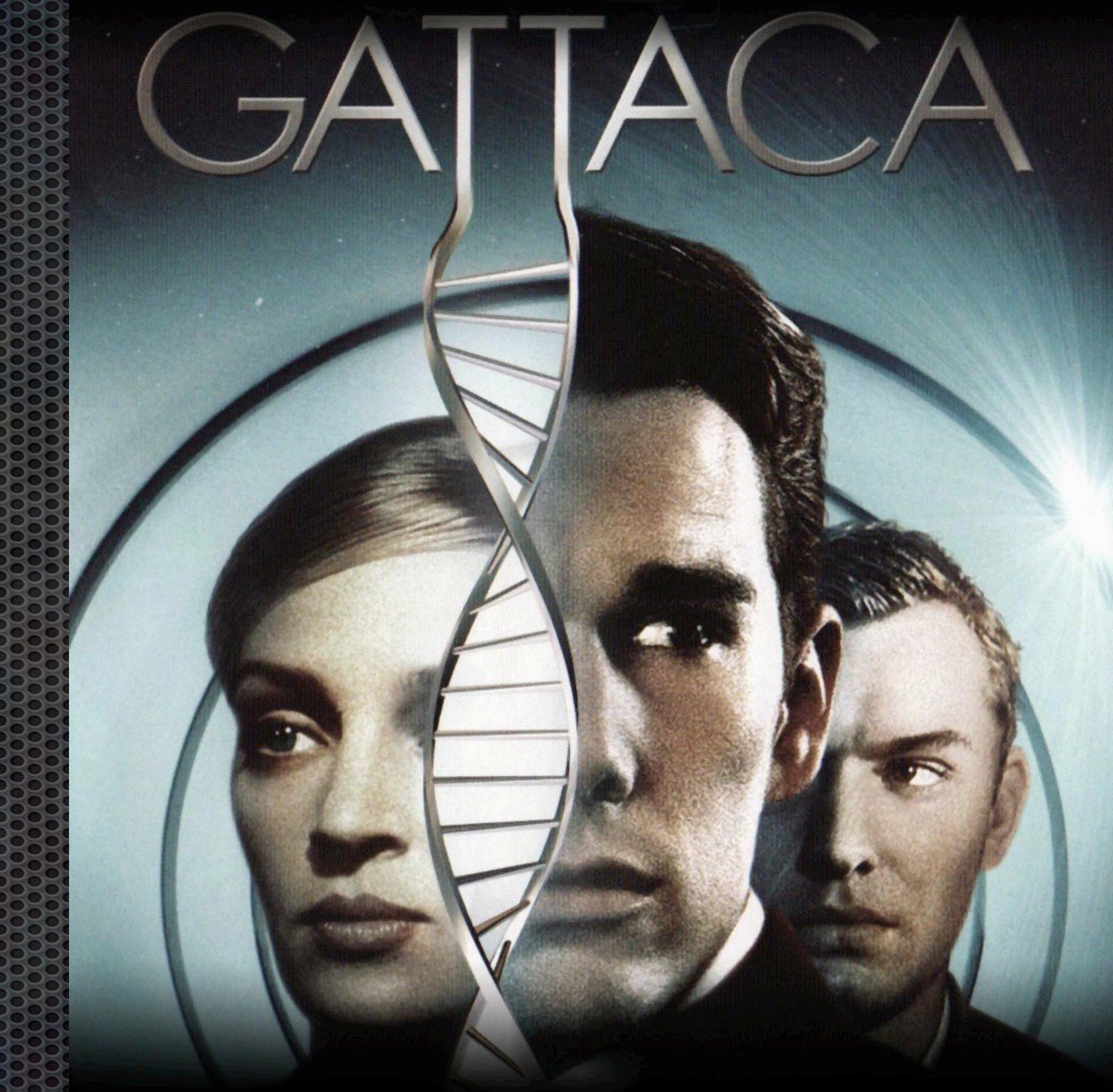
.....

.....

.....

.....

# Genomes & Privacy



Generalkonsent

PRIVACY

HACKERS

Health  
Insurance  
Portability and  
Accountability  
Act

BENEFIT

CONSENT

LAWS

SAFETY

BLOCKCHAIN

SECURITY

Right to Research

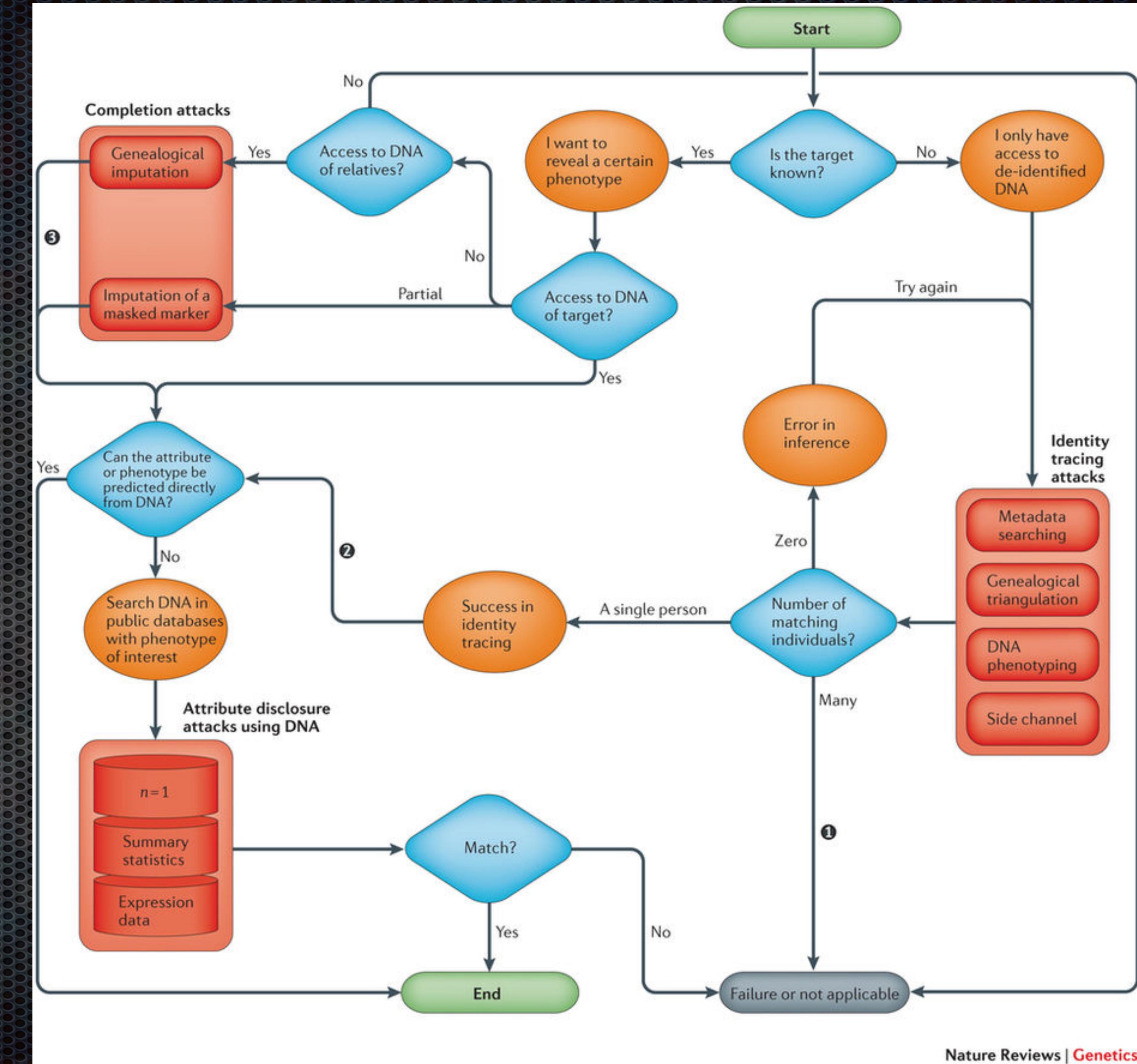
Genetic  
Information  
Nondiscrimination  
Act

CRYPTOGRAPHY

# Routes for breaching and protecting genetic privacy

The map contrasts different scenarios, such as identifying de-identified genetic data sets, revealing an attribute from genetic data and unmasking of data. It also shows the interdependencies between the techniques and suggests potential routes to exploit further information after the completion of one attack. There are several simplifying assumptions (black circles). In certain scenarios (such as insurance decisions), uncertainty about the target's identity within a small group of people could still be considered a success (assumption 1). For certain privacy harms (such as surveillance), identity tracing can be considered a success and the end point of the process (assumption 2). The complete DNA sequence is not always necessary (assumption 3).

Yaniv Erlich & Arvind Narayanan. *Nature Reviews Genetics* 15, 409–421 (2014)



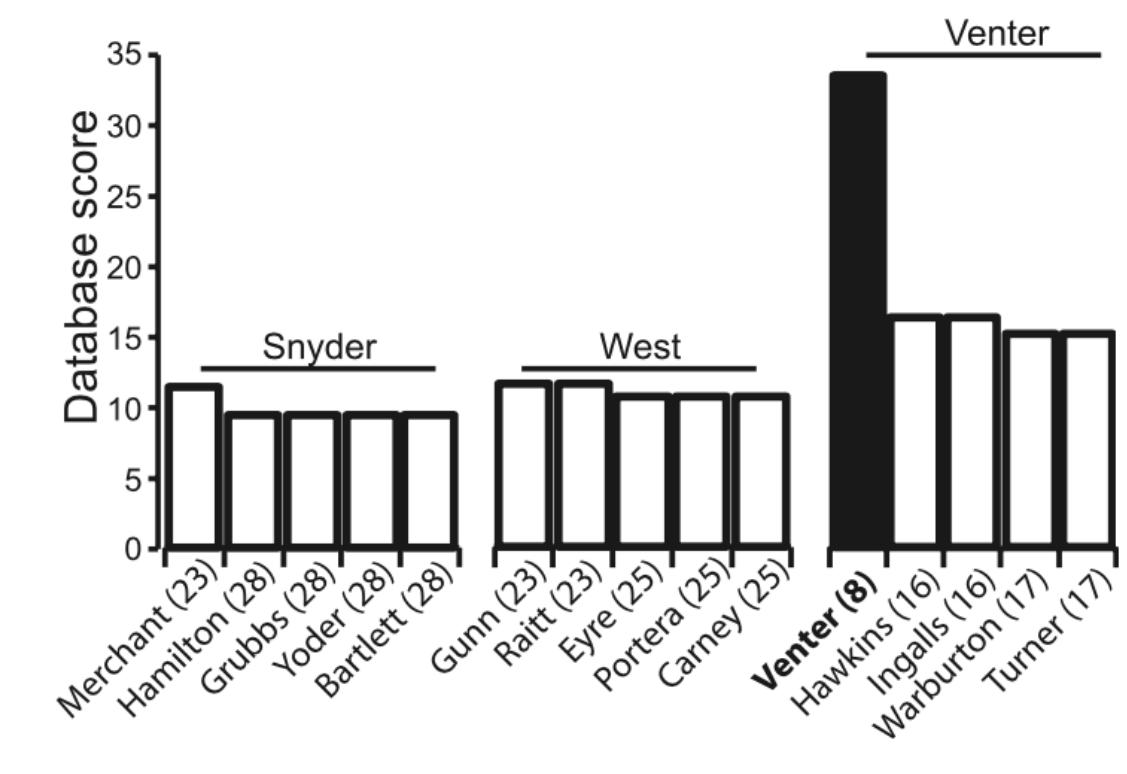
# IDENTIFICATION OF INDIVIDUALS BASED ON "GENOMIC FINGERPRINTS"

## Identifying Personal Genomes by Surname Inference

Melissa Gymrek,<sup>1,2,3,4</sup> Amy L. McGuire,<sup>5</sup> David Golan,<sup>6</sup> Eran Halperin,<sup>7,8,9</sup> Yaniv Erlich<sup>1\*</sup>

Sharing sequencing data sets without identifiers has become a common practice in genomics. Here, we report that surnames can be recovered from personal genomes by profiling short tandem repeats on the Y chromosome (Y-STRs) and querying recreational genetic genealogy databases. We show that a combination of a surname with other types of metadata, such as age and state, can be used to triangulate the identity of the target. A key feature of this technique is that it entirely relies on free, publicly accessible Internet resources. We quantitatively analyze the probability of identification for U.S. males. We further demonstrate the feasibility of this technique by tracing back with high probability the identities of multiple participants in public sequencing projects.

**Fig. 2.** The top five records retrieved after searching Ysearch with the Y-STR haplotypes of Michael Snyder, John West, and Craig Venter. The expected number of generations to the MRCA is given in parentheses for each record. Searching with Craig Venter returned a "Venter" record (closed bar) as the top match.



- ▶ Genomic data of many types can be used to re-identify individuals in data collections



# Genome *Beacons* Compromise Security?

## Querying for thousands of specific SNV occurrences in a genomic data pool can identify individuals

Stanford researchers identify potential security hole in genomic data-sharing network

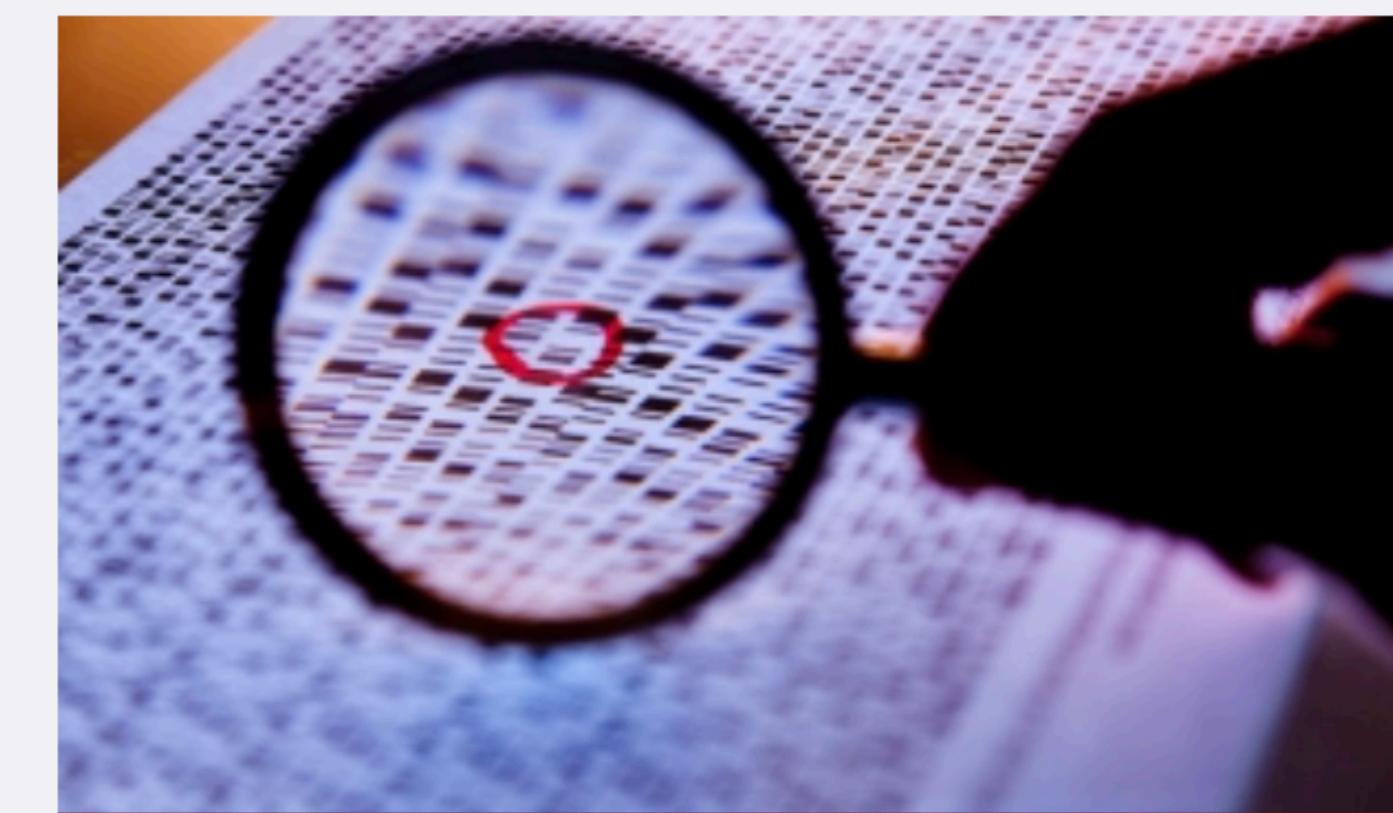
Hackers with access to a person's genome might find out if that genome is in an international network of disease databases.

OCT 29  
2015

Sharing genomic information among researchers is critical to the advance of biomedical research. Yet genomic data contains identifiable information and, in the wrong hands, poses a risk to individual privacy. If someone had access to your genome sequence — either directly from your saliva or other tissues, or from a popular genomic information service — they could check to see if you appear in a database of people with certain medical conditions, such as heart disease, lung cancer or autism.

Work by a pair of researchers at the [Stanford University School of Medicine](#) makes that genomic data more secure. [Suyash Shringarpure](#), PhD, a postdoctoral scholar in genetics, and [Carlos Bustamante](#), PhD, a professor of genetics, have demonstrated a technique for hacking a network of global genomic databases and how to prevent it. They are working with investigators from the Global Alliance for Genomics and Health on implementing preventive measures.

The work, published Oct. 29 in *The American Journal of Human Genetics*, also bears importantly on the larger question of how to analyze mixtures of genomes, such as those from different people at a crime scene.



Stanford researchers are working with the Global Alliance for Genomics and Health to make genomic information in the Beacon Project more secure.

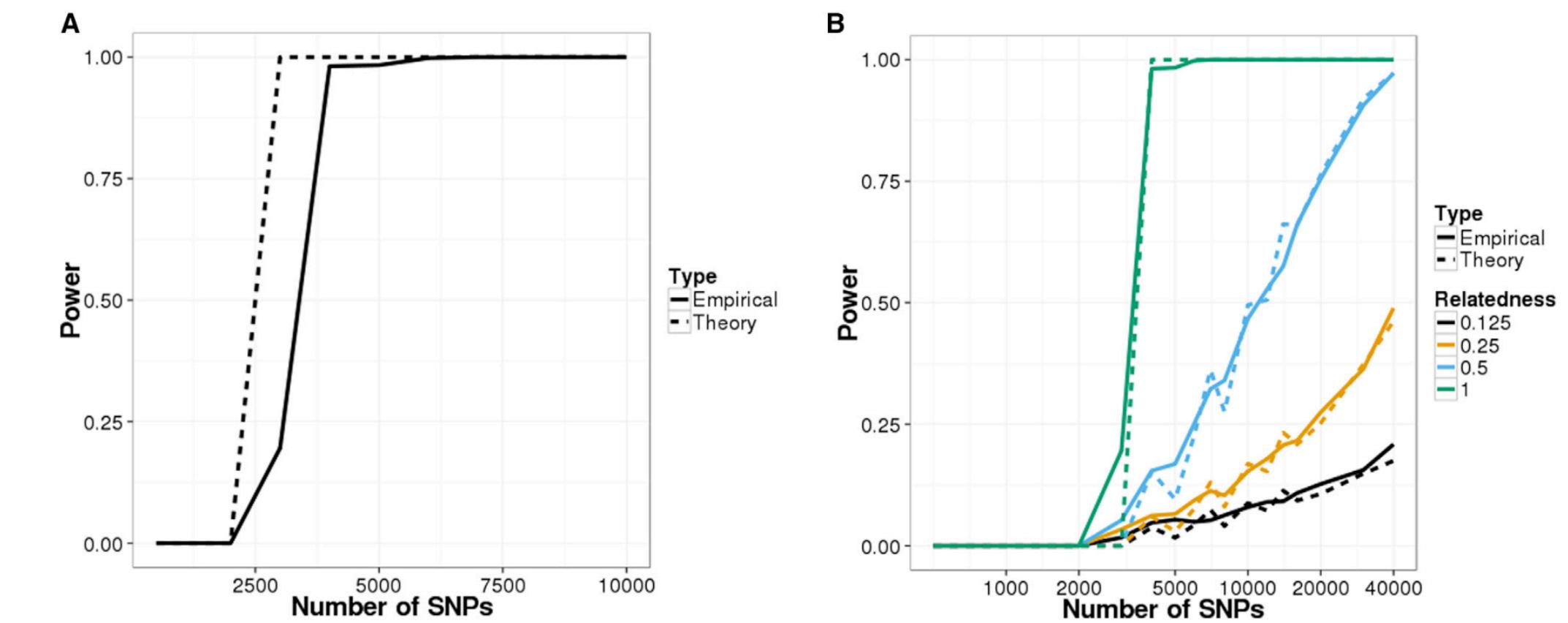
*Science photo/Shutterstock*

# IDENTIFICATION OF INDIVIDUALS FROM MIXED COLLECTIONS USING RARE ALLELES

## Privacy Risks from Genomic Data-Sharing Beacons

Suyash S. Shringarpure<sup>1,\*</sup> and Carlos D. Bustamante<sup>1,\*</sup>

The human genetics community needs robust protocols that enable secure sharing of genomic data from participants in genetic research. Beacons are web servers that answer allele-presence queries—such as “Do you have a genome that has a specific nucleotide (e.g., A) at a specific genomic position (e.g., position 11,272 on chromosome 1)?”—with either “yes” or “no.” Here, we show that individuals in a beacon are susceptible to re-identification even if the only data shared include presence or absence information about alleles in a beacon. Specifically, we propose a likelihood-ratio test of whether a given individual is present in a given genetic beacon. Our test is not dependent on allele frequencies and is the most powerful test for a specified false-positive rate. Through simulations, we showed that in a beacon with 1,000 individuals, re-identification is possible with just 5,000 queries. Relatives can also be identified in the beacon. Re-identification is possible even in the presence of sequencing errors and variant-calling differences. In a beacon constructed with 65 European individuals from the 1000 Genomes Project, we demonstrated that it is possible to detect membership in the beacon with just 250 SNPs. With just 1,000 SNP queries, we were able to detect the presence of an individual genome from the Personal Genome Project in an existing beacon. Our results show that beacons can disclose membership and implied phenotypic information about participants and do not protect privacy *a priori*. We discuss risk mitigation through policies and standards such as not allowing anonymous pings of genetic beacons and requiring minimum beacon sizes.



**Figure 1. Power of Re-identification Attacks on Beacons Constructed with Simulated Data**  
Power curves for the likelihood-ratio test (LRT) on (A) a simulated beacon with 1,000 individuals and (B) detecting relatives in the simulated beacon. The false-positive rate was set to 0.05 for all scenarios.

- ▶ rare allelic variants can be used to identify an individual (or her relatives) in a genome collection without having access to individual datasets
- ▶ however, such an approach requires previous knowledge about the individual's SNPs

# The Right to Scientific Knowledge

In 1948, the General assembly of the United nations adopted the Universal Declaration of Human Rights (UDHR) to guarantee the rights of every individual in the world. Included were twin rights “to share in scientific advancement and its benefits” and “to the protection of the moral and material interests resulting from any scientific...production of which [a person] is the author” (art. 27, United nations 1948).

from Knoppers et al, 2014

Hum Genet (2014) 133:895–903  
DOI 10.1007/s00439-014-1432-6

ORIGINAL INVESTIGATION

## A human rights approach to an international code of conduct for genomic and clinical data sharing

Bartha M. Knoppers · Jennifer R. Harris ·  
Isabelle Budin-Ljøsne · Edward S. Dove

Received: 9 December 2013 / Accepted: 16 February 2014 / Published online: 27 February 2014  
© The Author(s) 2014. This article is published with open access at Springerlink.com

**Abstract** Fostering data sharing is a scientific and ethical imperative. Health gains can be achieved more comprehensively and quickly by combining large, information-rich datasets from across conventionally siloed disciplines and geographic areas. While collaboration for data sharing is increasingly embraced by policymakers and the international biomedical community, we lack a common ethical and legal framework to connect regulators, funders, consortia, and research projects so as to facilitate genomic and clinical data linkage, global science collaboration, and responsible research conduct. Governance tools can be used to responsibly steer the sharing of data for proper stewardship of research discovery, genomics research resources, and their clinical applications. In this article, we propose that an international code of conduct be designed to enable global genomic and clinical data sharing for biomedical research. To give this proposed code universal application and accountability, however, we propose to position it within a human rights framework. This proposition is not without precedent: international treaties have long recognized that everyone has a right to the benefits of scientific

progress and its applications, and a right to the protection of the moral and material interests resulting from scientific productions. It is time to apply these twin rights to internationally collaborative genomic and clinical data sharing.

### Introduction

In 1948, the General Assembly of the United Nations adopted the *Universal Declaration of Human Rights* (UDHR) to guarantee the rights of every individual in the world. Included were twin rights “to share in scientific advancement and its benefits” and “to the protection of the moral and material interests resulting from any scientific...production of which [a person] is the author” (Art. 27, United Nations 1948). In the 21st century, where are we in realizing the sharing of scientific advancement and its benefits, and the importance of protecting a scientific producer’s moral and material interests? In this article, we argue that these little-developed twin rights, what we call the right “to benefit from” and “to be recognized for”, have direct application to internationally collaborative genomic and clinical data sharing, and can be activated through an international code of conduct.

Sharing genomic and clinical data is critical to achieve precision medicine (National Research Council 2011), that is, more accurate disease classification based on molecular profiles to enable tailored effective treatments, interventions, and models for prevention. Better communication flow across borders and research teams, encompassing data from clinical and population research, enables researchers to connect the diverse types of datasets and expertise needed to elucidate the genomic basis and complexities of disease etiology. Such data integration can make it possible to reveal the genetic basis of cancer, inherited diseases,

B. M. Knoppers (✉) · E. S. Dove  
Centre of Genomics and Policy, McGill University, 740 Dr.  
Penfield Avenue, Suite 5200, Montreal H3A 0G1, Canada  
e-mail: bartha.knoppers@mcgill.ca

E. S. Dove  
e-mail: edward.dove@mcgill.ca

J. R. Harris · I. Budin-Ljøsne  
Division of Epidemiology, Department of Genes  
and Environment, Norwegian Institute of Public Health,  
PO Box 4404, Nydalen 0403, Oslo, Norway  
e-mail: Jennifer.Harris@fhi.no

I. Budin-Ljøsne  
e-mail: Isabelle.Budin.Ljosne@fhi.no

# Modernizing Patient Consent

- forward looking, transparent and technically feasible regulations for enabling access to research material and data while empowering patients

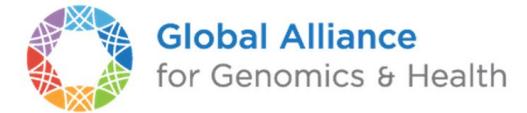
## Generalkonsent: Eine einheitliche Vorlage soll schweizweite Forschung erleichtern

| Art des Forschungs-materials<br>Personenbezug | Biologisches Material und genetische Daten                                                                                                                                                                                                          | Nicht-genetische Daten                                                                                                                                                                        |              |
|-----------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------|
|                                               | Unverschlüsselt (identifizierend)                                                                                                                                                                                                                   | Verschlüsselt                                                                                                                                                                                 | Anonymisiert |
|                                               | Information + Einwilligung in jedes einzelne Forschungsprojekt                                                                                                                                                                                      | Information über Weiterverwendung für zukünftige noch unbestimmte Forschungsprojekte + Generalkonsent für Forschungszwecke                                                                    |              |
|                                               | Information über Weiterverwendung für zukünftige noch unbestimmte Forschungsprojekte + Generalkonsent für Forschungszwecke                                                                                                                          | Information über Weiterverwendung für zukünftige noch unbestimmte Forschungsprojekte + Generalkonsent für Forschungszwecke + über Möglichkeit Weiterverwendung abzulehnen > Widerspruchsrecht |              |
|                                               | <b>Genetische Daten:</b> Information über Weiterverwendung für zukünftige noch unbestimmte Forschungszwecke + über Möglichkeit Weiterverwendung abzulehnen > Widerspruchsrecht<br><b>Proben:</b> Information zur Anonymisierung > Widerspruchsrecht | Ausserhalb des Geltungsbereichs des HFG                                                                                                                                                       |              |

Switzerland: Definition of a unified "Generalkonsent", to provide a single framework to manage permissions for access to patient derived material and related data

## Consent Codes: Upholding Standard Data Use Conditions

Stephanie O. M. Dyke<sup>1\*</sup>, Anthony A. Philippakis<sup>2</sup>, Jordi Rambla De Argila<sup>3,4</sup>, Dina N. Paltoo<sup>5</sup>, Erin S. Luetkemeier<sup>5</sup>, Bartha M. Knoppers<sup>1</sup>, Anthony J. Brookes<sup>6</sup>, J. Dylan Spalding<sup>7</sup>, Mark Thompson<sup>8</sup>, Marco Roos<sup>8</sup>, Kym M. Boycott<sup>9</sup>, Michael Brudno<sup>10,11</sup>, Matthew Hurles<sup>12</sup>, Heidi L. Rehm<sup>2,13</sup>, Andreas Matern<sup>14</sup>, Marc Fiume<sup>15</sup>, Stephen T. Sherry<sup>16</sup>



| Consent Codes                                                                                                                |              |                                                                                                                                                           |
|------------------------------------------------------------------------------------------------------------------------------|--------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------|
| Name                                                                                                                         | Abbreviation | Description                                                                                                                                               |
| <b>Primary Categories (I<sup>IV</sup>)</b>                                                                                   |              |                                                                                                                                                           |
| no restrictions                                                                                                              | NRES         | No restrictions on data use.                                                                                                                              |
| general research use and clinical care                                                                                       | GRU(CC)      | For health/medical/biomedical purposes and other biological research, including the study of population origins or ancestry.                              |
| health/medical/biomedical research and clinical care                                                                         | HMB(CC)      | Use of the data is limited to health/medical/biomedical purposes, does not include the study of population origins or ancestry.                           |
| disease-specific research and clinical care                                                                                  | DS-[XX](CC)  | Use of the data must be related to [disease].                                                                                                             |
| population origins/ancestry research                                                                                         | POA          | Use of the data is limited to the study of population origins or ancestry.                                                                                |
| <b>Secondary Categories (II<sup>IV</sup>)</b> (can be one or more extra conditions, in addition to I <sup>IV</sup> category) |              |                                                                                                                                                           |
| other research-specific restrictions                                                                                         | RS-[XX]      | Use of the data is limited to studies of [research type] (e.g., pediatric research).                                                                      |
| research use only                                                                                                            | RUO          | Use of data is limited to research purposes (e.g., does not include its use in clinical care).                                                            |
| no “general methods” research                                                                                                | NMDS         | Use of the data includes methods development research (e.g., development of software or algorithms) ONLY within the bounds of other data use limitations. |
| genetic studies only                                                                                                         | GSO          | Use of the data is limited to genetic studies only (i.e., no research using only the phenotype data).                                                     |
| <b>Requirements</b>                                                                                                          |              |                                                                                                                                                           |
| not-for-profit use only                                                                                                      | NPU          | Use of the data is limited to not-for-profit organizations.                                                                                               |
| publication required                                                                                                         | PUB          | Requestor agrees to make results of studies using the data available to the larger scientific community.                                                  |
| collaboration required                                                                                                       | COL-[XX]     | Requestor must agree to collaboration with the primary study investigator(s).                                                                             |
| return data to database/resource                                                                                             | RTN          | Requestor must return derived/enriched data to the database/resource.                                                                                     |
| ethics approval required                                                                                                     | IRB          | Requestor must provide documentation of local IRB/REC approval.                                                                                           |
| geographical restrictions                                                                                                    | GS-[XX]      | Use of the data is limited to within [geographic region].                                                                                                 |
| publication moratorium/embargo                                                                                               | MOR-[XX]     | Requestor agrees not to publish results of studies until [date].                                                                                          |
| time limits on use                                                                                                           | TS-[XX]      | Use of data is approved for [x months].                                                                                                                   |
| user-specific restrictions                                                                                                   | US           | Use of data is limited to use by approved users.                                                                                                          |
| project-specific restrictions                                                                                                | PS           | Use of data is limited to use within an approved project.                                                                                                 |
| institution-specific restrictions                                                                                            | IS           | Use of data is limited to use within an approved institution.                                                                                             |

SOM Dyke, et al. Consent Codes: Upholding Standard Data Use Conditions. *PLoS Genetics* 12(1): e1005772. <http://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1005772>

Contact: Dr. Stephanie Dyke (stephanie.dyke@mcgill.ca)

# Health Related Data & Privacy

- Is the genetic condition **outwardly visible**?
- How **severe** is it? (serious disease, **penetrance**, age of onset)
- Is it associated with what could be considered to be **stigmatizing** health information (e.g., associated with **mental** health, **reproductive** care, **disability**)?
- Is it **familial** (i.e., potential carrier status/reproductive implications for family/relatives)?
- Does it provide information about the likely **geographical location** of individuals?
- Does it provide information about **ethnicity** that may be considered potentially stigmatizing information?

## Sharing health-related data: a privacy test?

Stephanie OM Dyke<sup>1</sup>, Edward S Dove<sup>2</sup> and Bartha M Knoppers<sup>1</sup>

Greater sharing of potentially sensitive data raises important ethical, legal and social issues (ELSI), which risk hindering and even preventing useful data sharing if not properly addressed. One such important issue is respecting the privacy-related interests of individuals whose data are used in genomic research and clinical care. As part of the Global Alliance for Genomics and Health (GA4GH), we examined the ELSI status of health-related data that are typically considered 'sensitive' in international policy and data protection laws. We propose that 'tiered protection' of such data could be implemented in contexts such as that of the GA4GH Beacon Project to facilitate responsible data sharing. To this end, we discuss a Data Sharing Privacy Test developed to distinguish degrees of sensitivity within categories of data recognised as 'sensitive'. Based on this, we propose guidance for determining the level of protection when sharing genomic and health-related data for the Beacon Project and in other international data sharing initiatives.

*npj Genomic Medicine* (2016) **1**, 16024; doi:10.1038/npjgenmed.2016.24; published online 17 August 2016

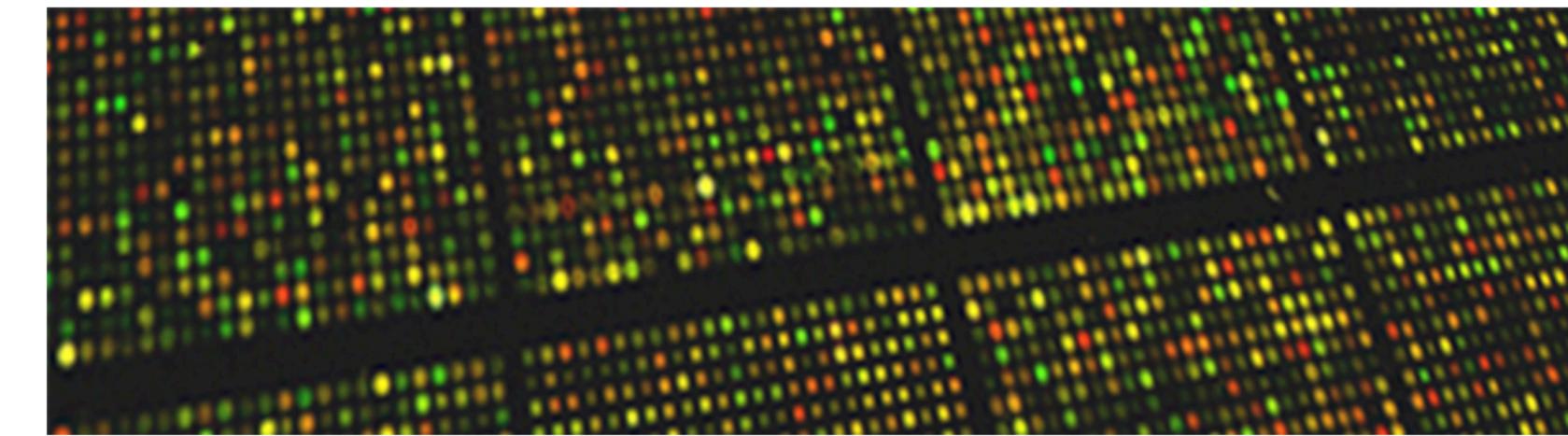


**Figure 1.** The three steps of a Data Sharing Privacy Test to distinguish degrees of data sensitivity within categories of data recognised as 'sensitive'.

# SHARE YOUR GENOME DATA?

- ▶ depositing genome data has the inherent risk of being identified and linked to your person
- ▶ What are the Risks?
- ▶ Would you contribute e.g. to OpenSNP?
- ▶ Discuss!

## Welcome to *openSNP*



*openSNP* lets customers of direct-to-customer genetic tests publish their test results, find others with similar genetic variations, learn more about their results by getting the latest primary literature on their variations, and helps scientists find new associations.

[Sign Up!](#)[Download the data!](#)

For Genotyping Users    For Scientists

### Upload Your Genotyping File



Upload your raw genotyping or exome data from [23andMe](#), [deCODEme](#) or [FamilyTreeDNA](#) to the *openSNP* database to make it available for everybody.

### Share Your Phenotypes & Traits



Phenotypes are the observable characteristics of your body, such as height, eye color or preference for coffee. Share your phenotype with other *openSNP* users, and find others with similar characteristics and traits. Your data may help scientists discover new genetic associations!

### Share your stories on variations & phenotypes



With *openSNP* you can share stories about your genetic variations and phenotypes, and discover the stories of other users.

### Find literature on genetic variation

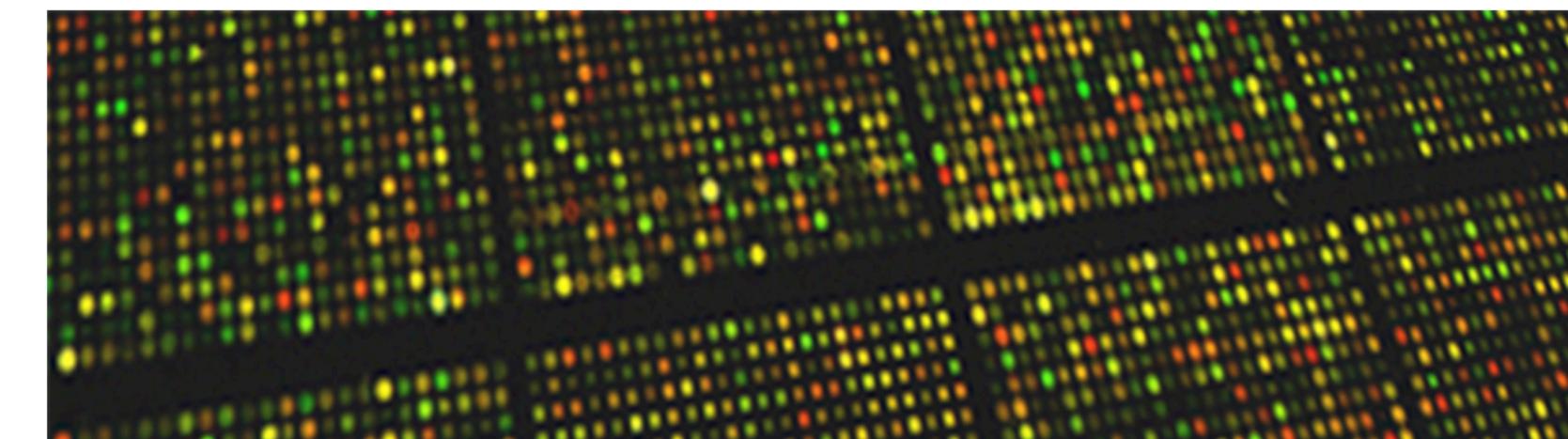


*openSNP* gets the latest open access journal articles on genetic variations from the [Public Library of Science](#). Popular articles are indexed via the social reference manager [Mendeley](#), and summaries are provided by [SNPedia](#).

# SHARE YOUR GENOME DATA?



## Welcome to *openSNP*



*openSNP* lets customers of direct-to-customer genetic tests publish their test results, find others with similar genetic variations, learn more about their results by getting the latest primary literature on their variations, and helps scientists find new associations.

[Sign Up!](#)[Download the data!](#)

For Genotyping Users    For Scientists

**Upload Your Genotyping File**



Upload your raw genotyping or exome data from *23andMe*, *deCODEme* or *FamilyTreeDNA* to the *openSNP* database to make it available for everybody.

**Share Your Phenotypes & Traits**



Phenotypes are the observable characteristics of your body, such as height, eye color or preference for coffee. Share your phenotype with other *openSNP* users, and find others with similar characteristics and traits. Your data may help scientists discover new genetic associations!

**Share your stories on variations & phenotypes**



With *openSNP* you can share stories about your genetic variations and phenotypes, and discover the stories of other users.

**Find literature on genetic variation**



*openSNP* gets the latest open access journal articles on genetic variations from the *Public Library of Science*. Popular articles are indexed via the social reference manager *Mendeley*, and summaries are provided by *SNPedia*.

Prof. Dr. Michael Baudis  
Institute 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*