

Cancer Genome (Data) Analysis

Why, Which, When, Whereto?

Michael Baudis
Computational Oncogenomics

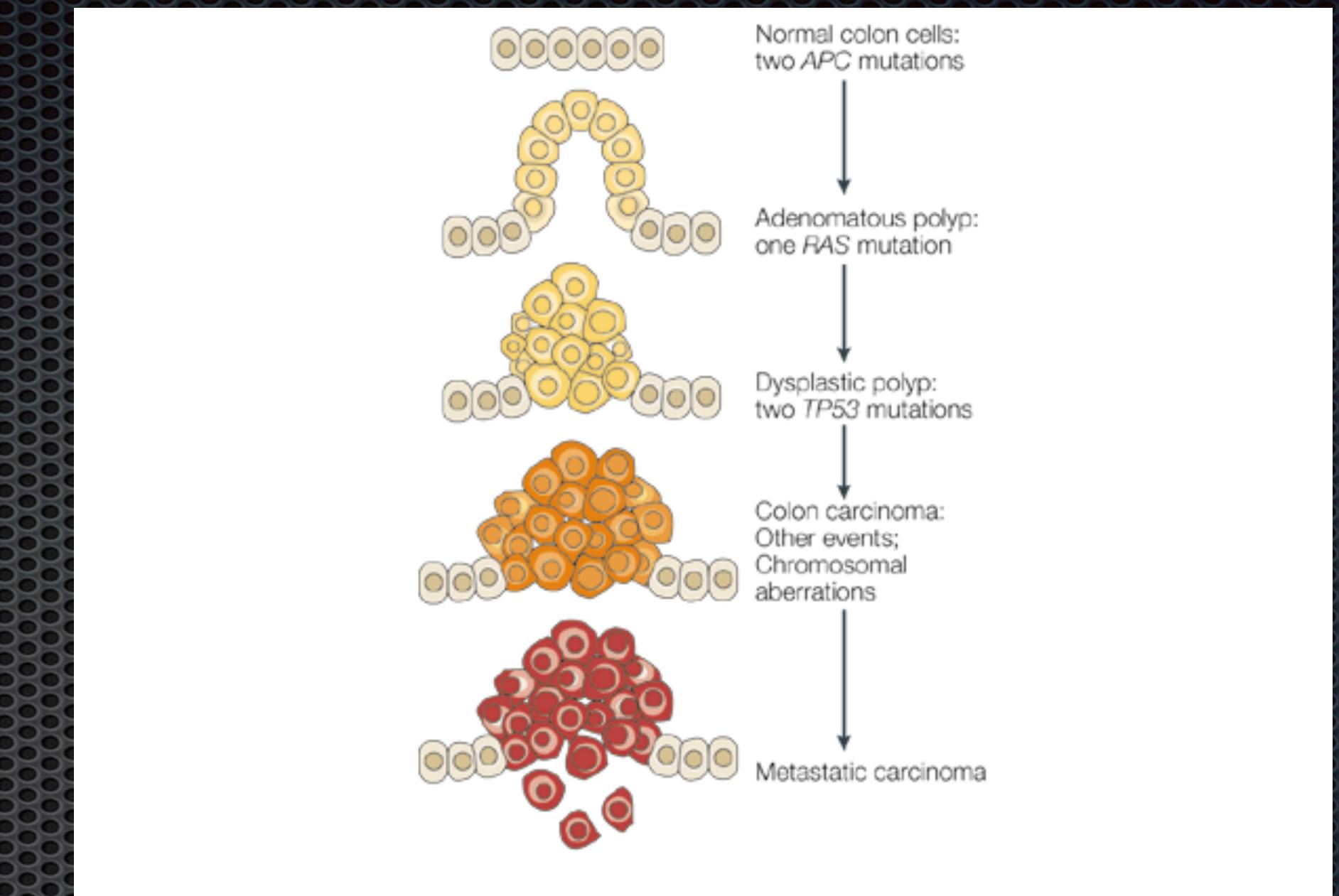
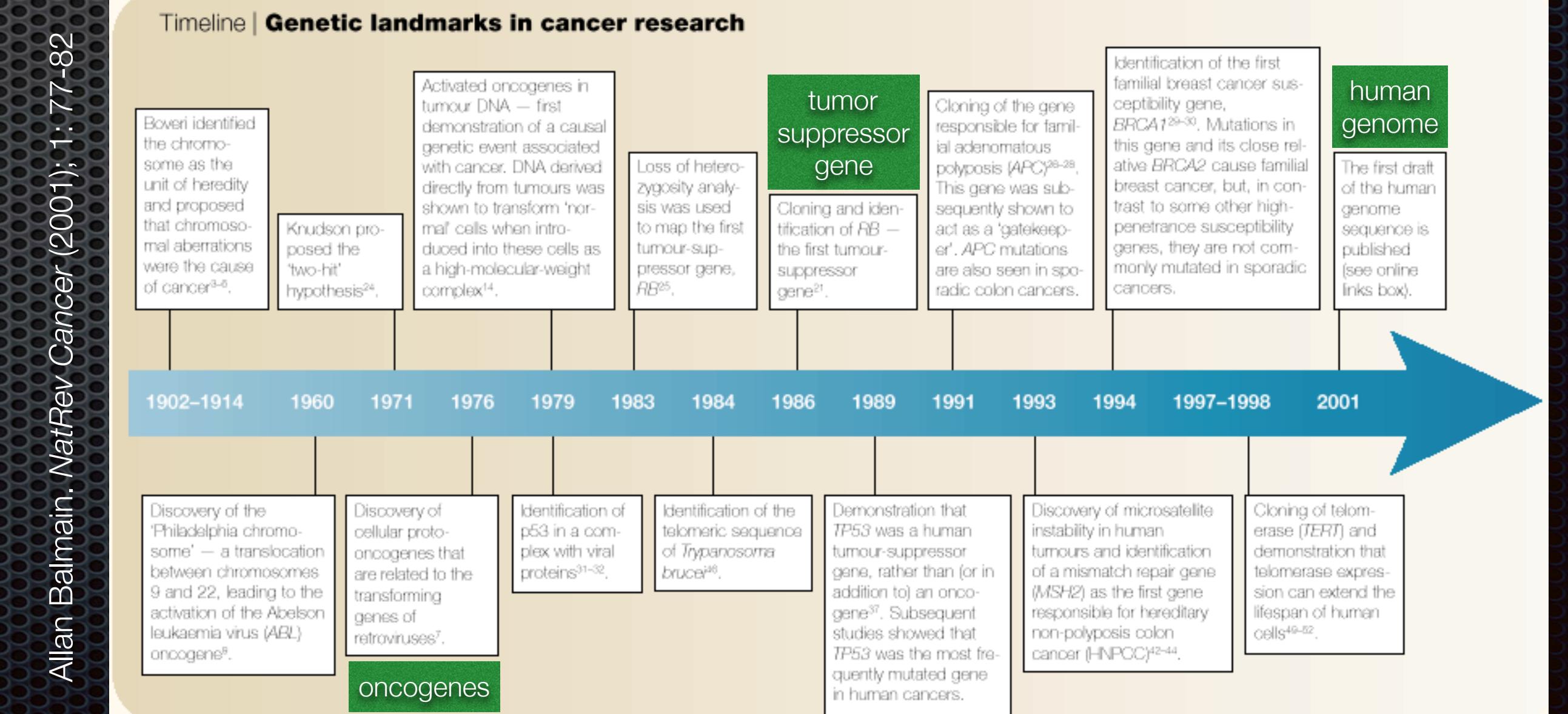


University of
Zurich ^{UZH}

Cancers are diseases of genomes

- every cancer's genome is different
- every individual's genome is different

Cancers arise from the clonal accumulation of **somatic genome mutations**, with varying but **limited** contribution of **inherited risk**



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



Theodor Boveri (1914)

(based on observations in sea urchin eggs)

- **Cell-cycle checkpoints** (“Hemmungseinrichtung”)
- **Tumour-suppressor genes** (“Teilungshemmende Chromosomen”); can be eliminated during tumour progression
- amplified **Oncogenes** (“Teilungsfoerdernde Chromosomen” ... “im permanenten Übergewicht”)
- sequential **Progression** (benign to malignant)
- Cancer **predisposition** through inheritance of less able suppressor “chromosomes”
- high-penetrance cancer syndromes - (e.g. xeroderma pigmentosum) through **homozygosity**
- Clonal origin & Genetic mosaicism; wounding and inflammation in tumour promotion; loss of cell adhesion in metastasis; sensitivity of malignant cells to radiation therapy (based on Hertwig *et al.*)

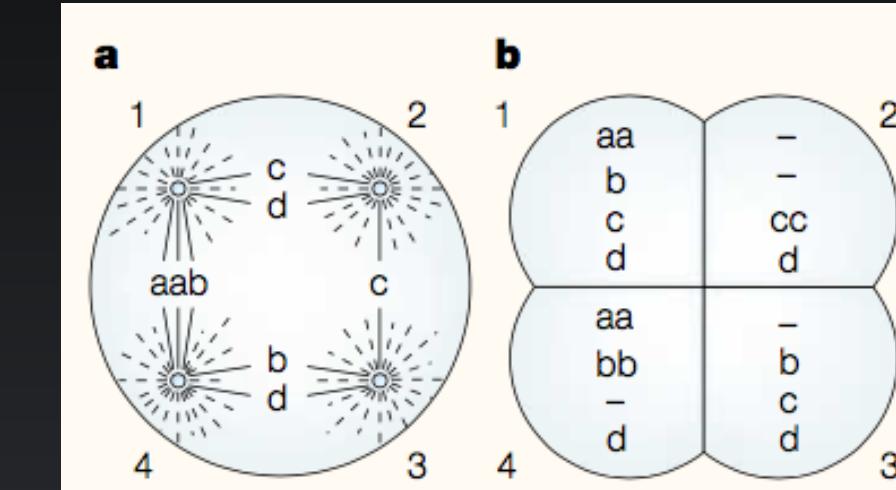
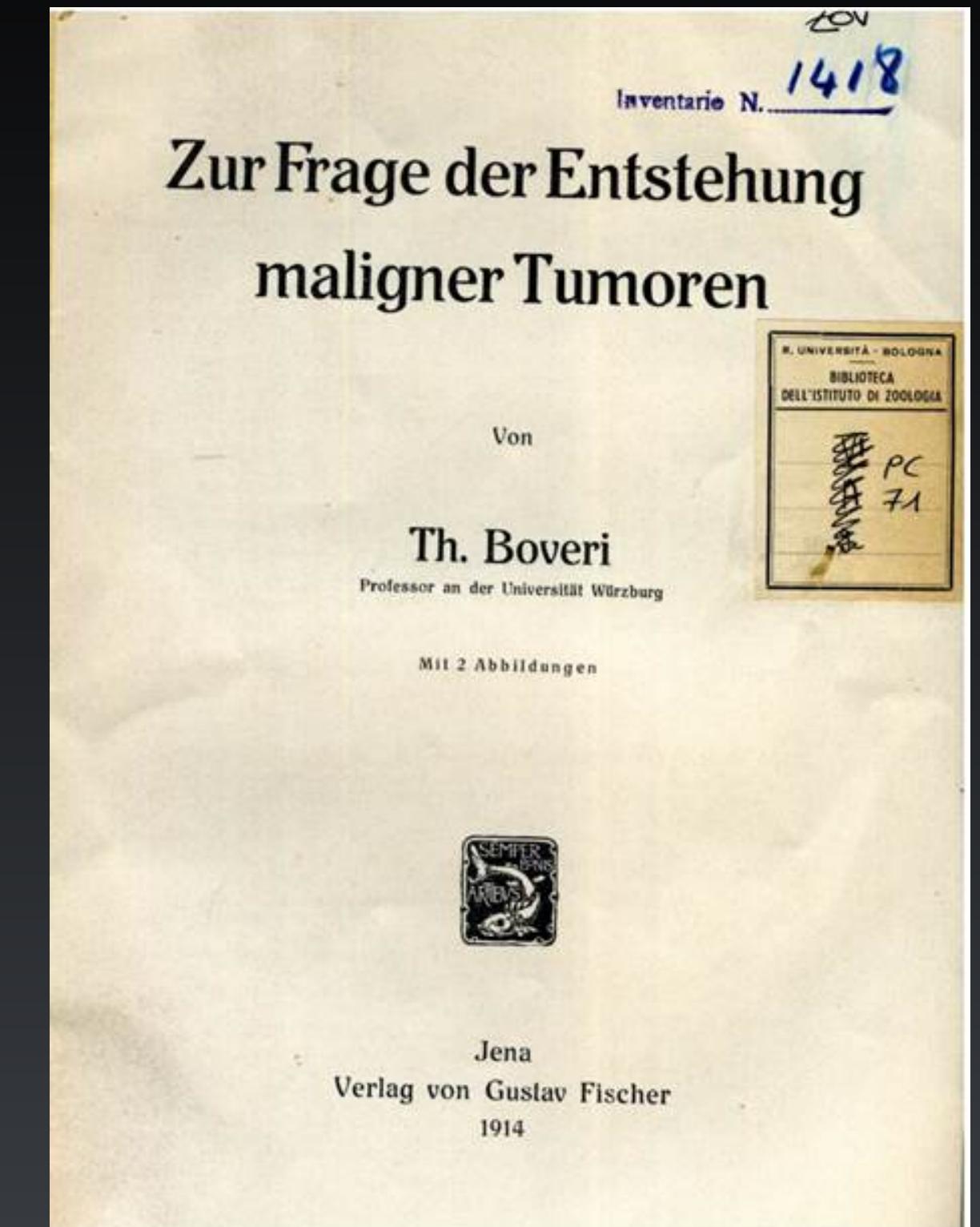


Figure 2 | Multiple cell poles cause unequal segregation of chromosomes. **a** | Boveri showed that fertilization of sea-urchin eggs by two sperm results in multiple cell poles. Individual chromosomes then attach to different combinations of poles — for example, one copy of chromosome c is attached to poles 1 and 2, and one copy is attached to poles 2 and 3. **b** | Chromosomes are segregated to the four poles at cell division, leaving some cells with too many copies of the chromosomes and some with too few — for example, cell 2 has two copies of chromosome c and cell 4 has none.



Allan Balmain
Cancer genetics: from Boveri and
Mendel to microarrays.
NatRev Cancer (2001); 1: 77-82

Anna Di Leonardo , Sergio Nasi , Simonetta Pulciani
Cancer: We Should Not Forget The Past
Journal of Cancer (2015), Vol. 6: 29-39
(for book cover & summary)



Janet Rowley (1972/73)

Chromosomal translocations in cancer

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

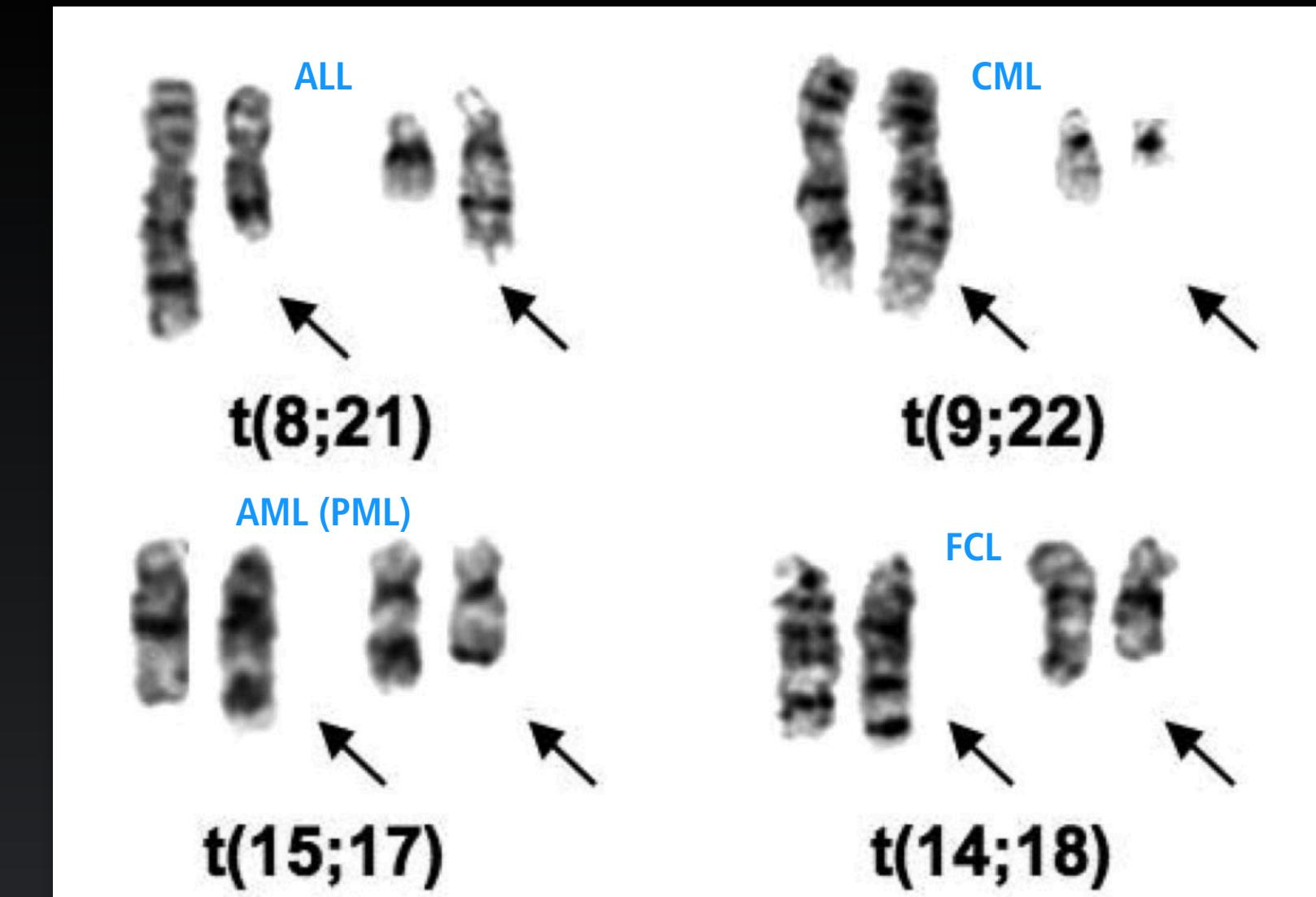
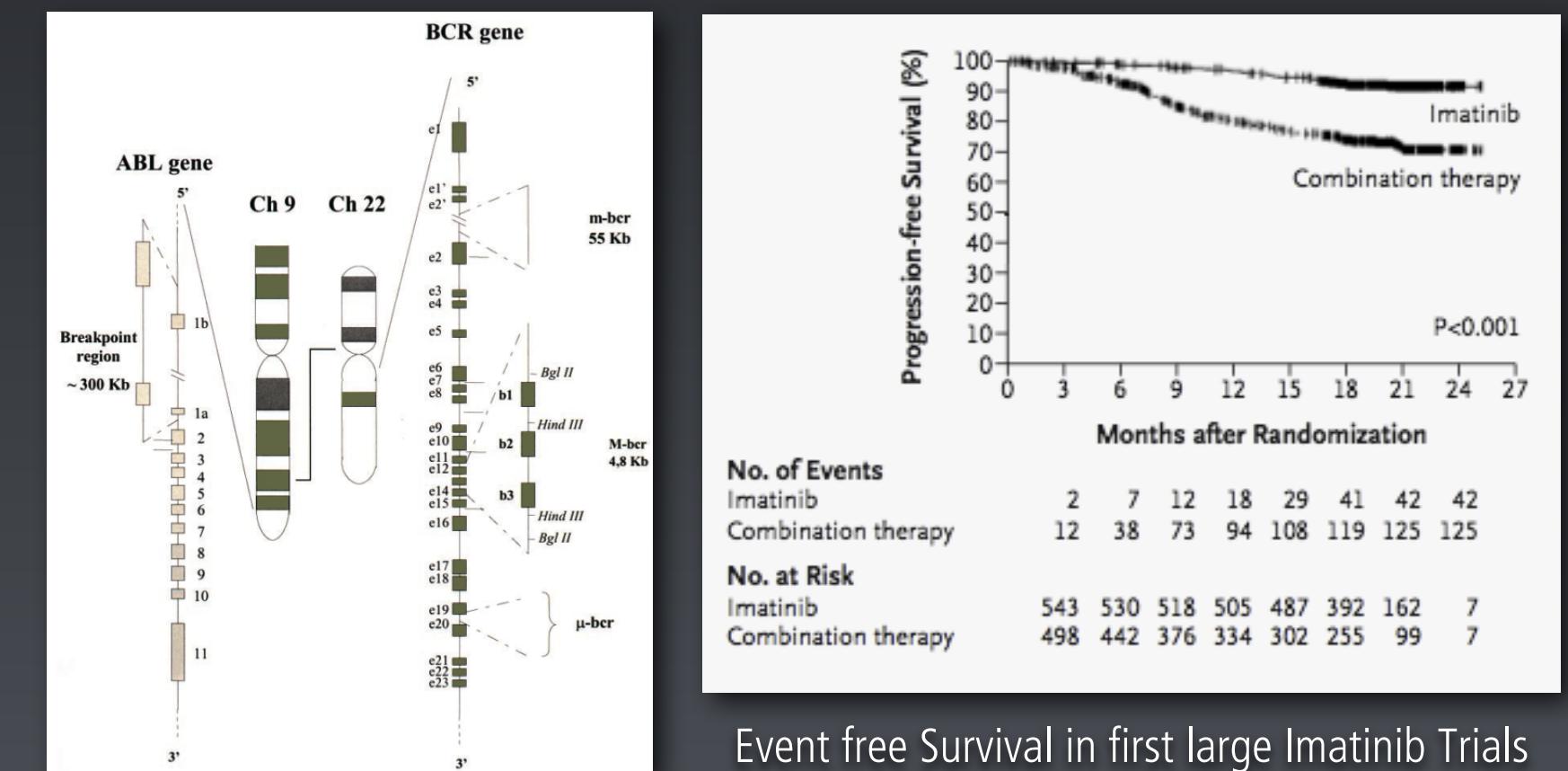


Figure 1. Partial karyotypes of common translocations discovered by Rowley.
The translocations appear in the order in which they were discovered.

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



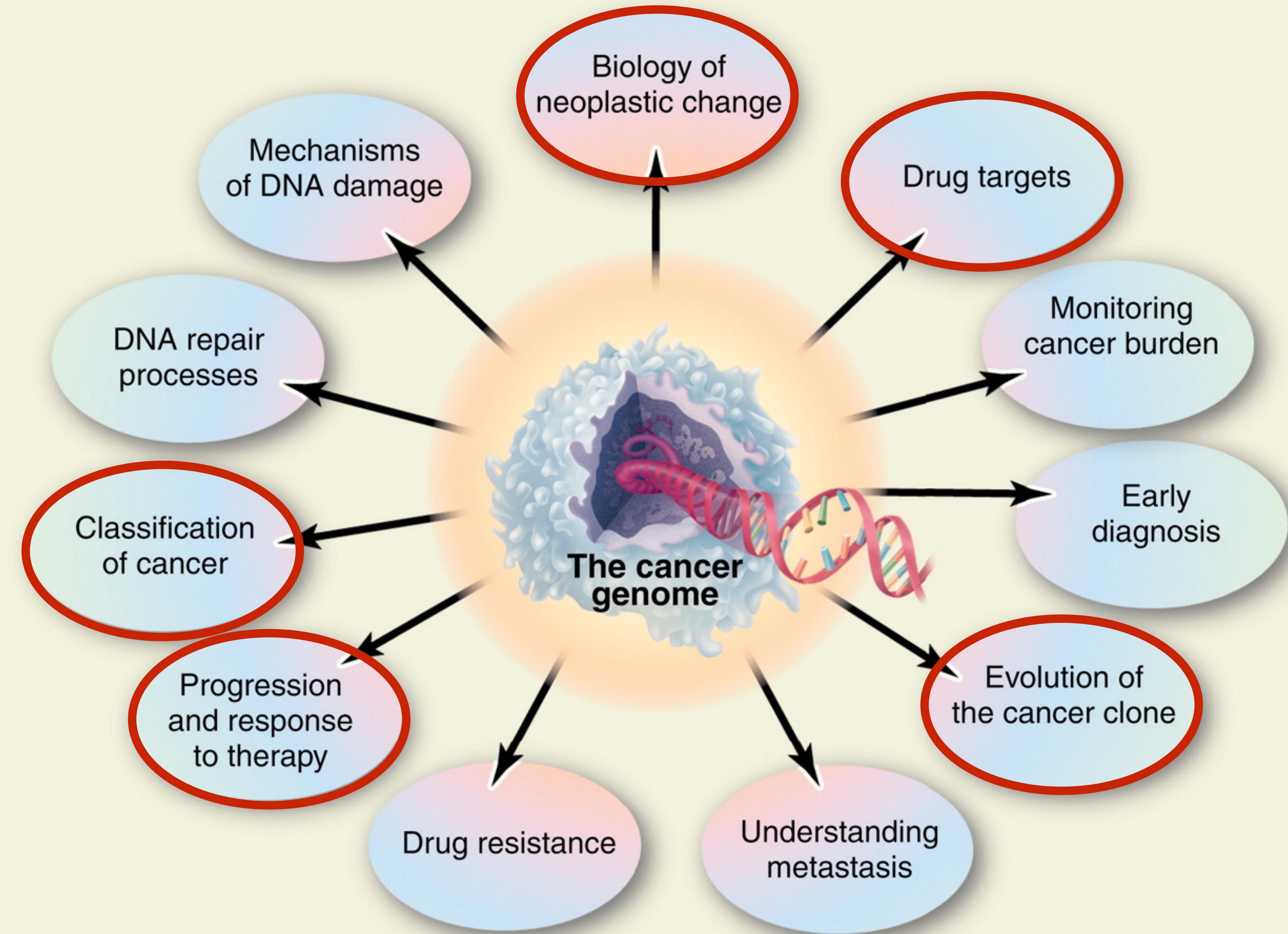
Event free Survival in first large Imatinib Trials

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

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

Why?

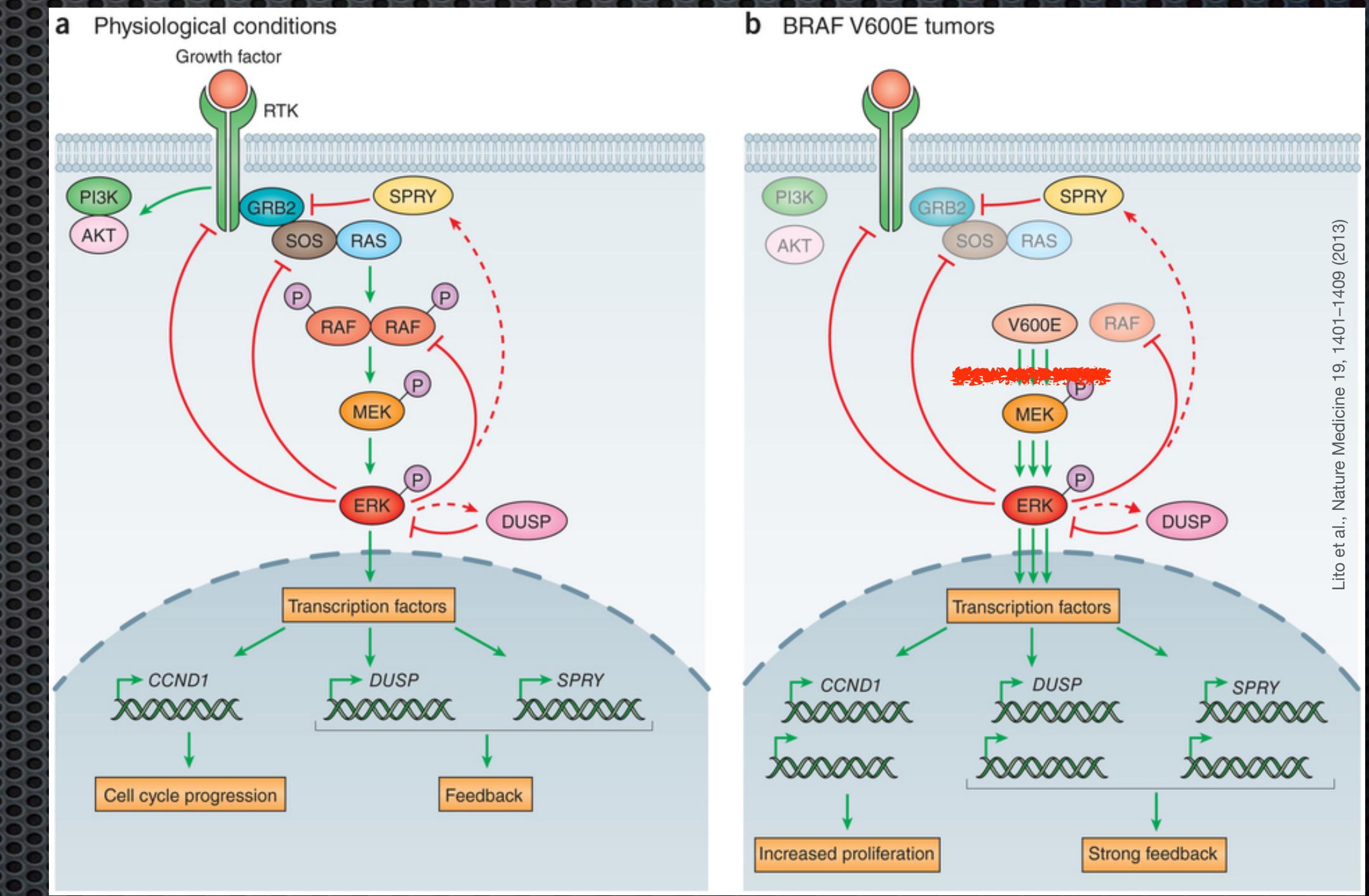
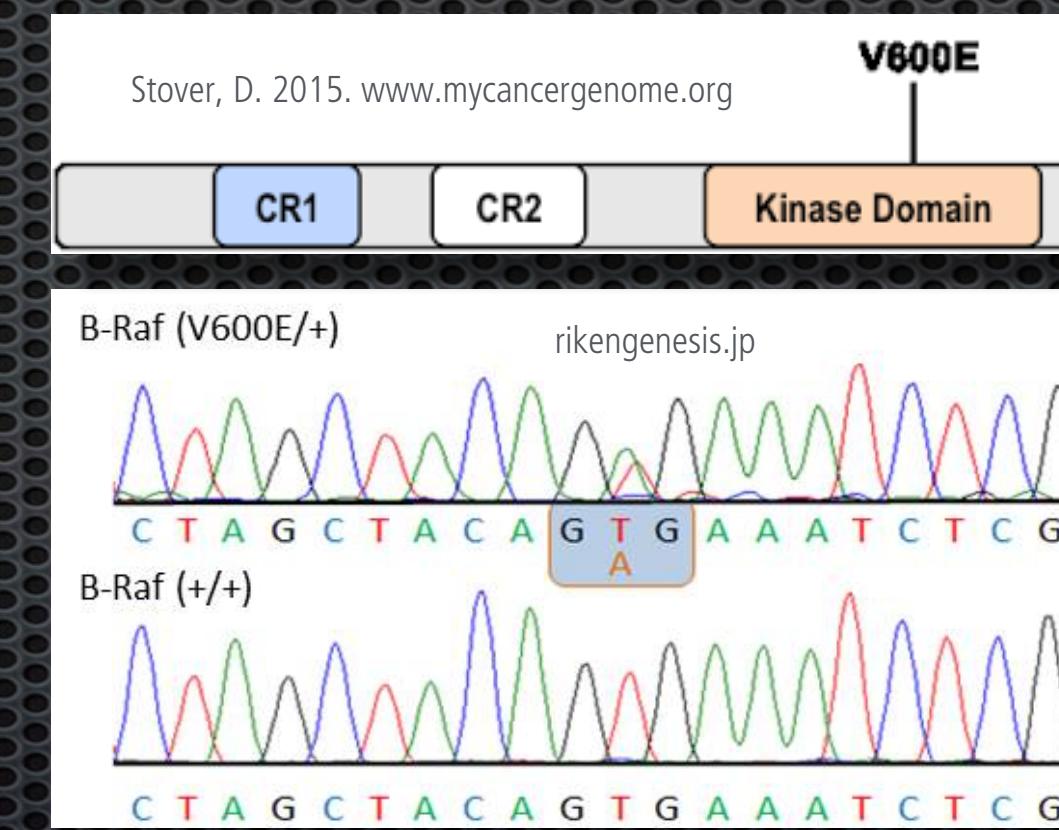
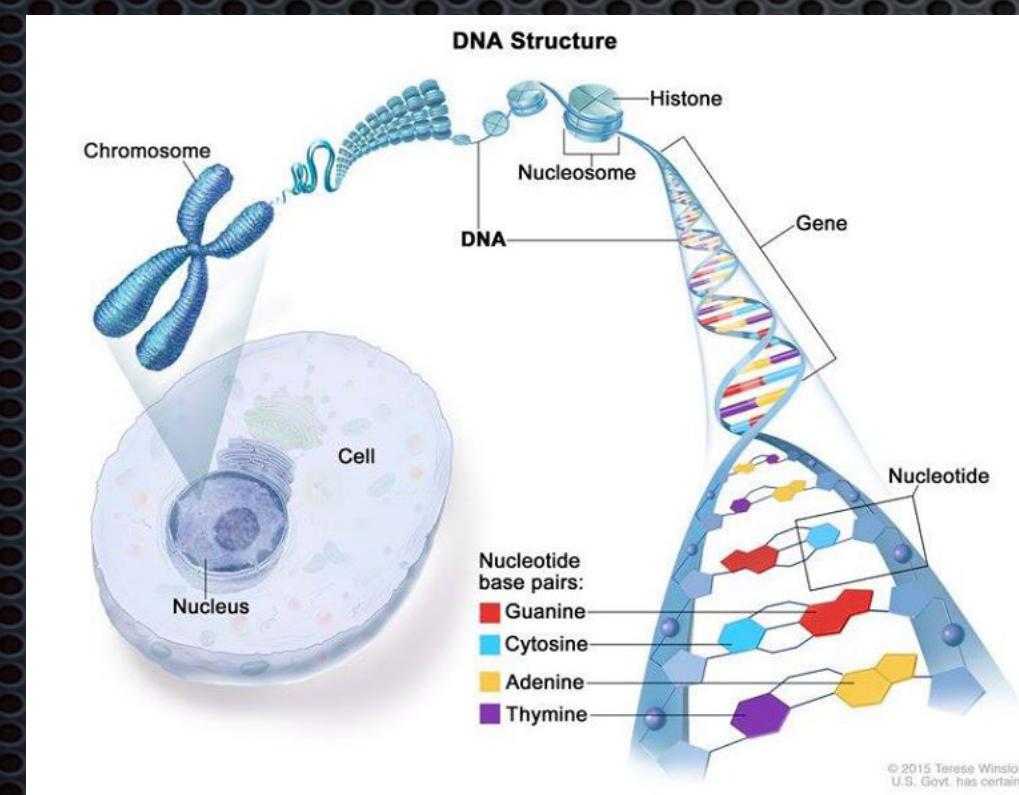
Michael R. Stratton.
Exploring the Genomes
of Cancer Cells:
Progress and Promise.
Science (2011)



BRAF V600E (c.1799T>A) Mutation

Oncogene Activation by Single Nucleotide Alteration

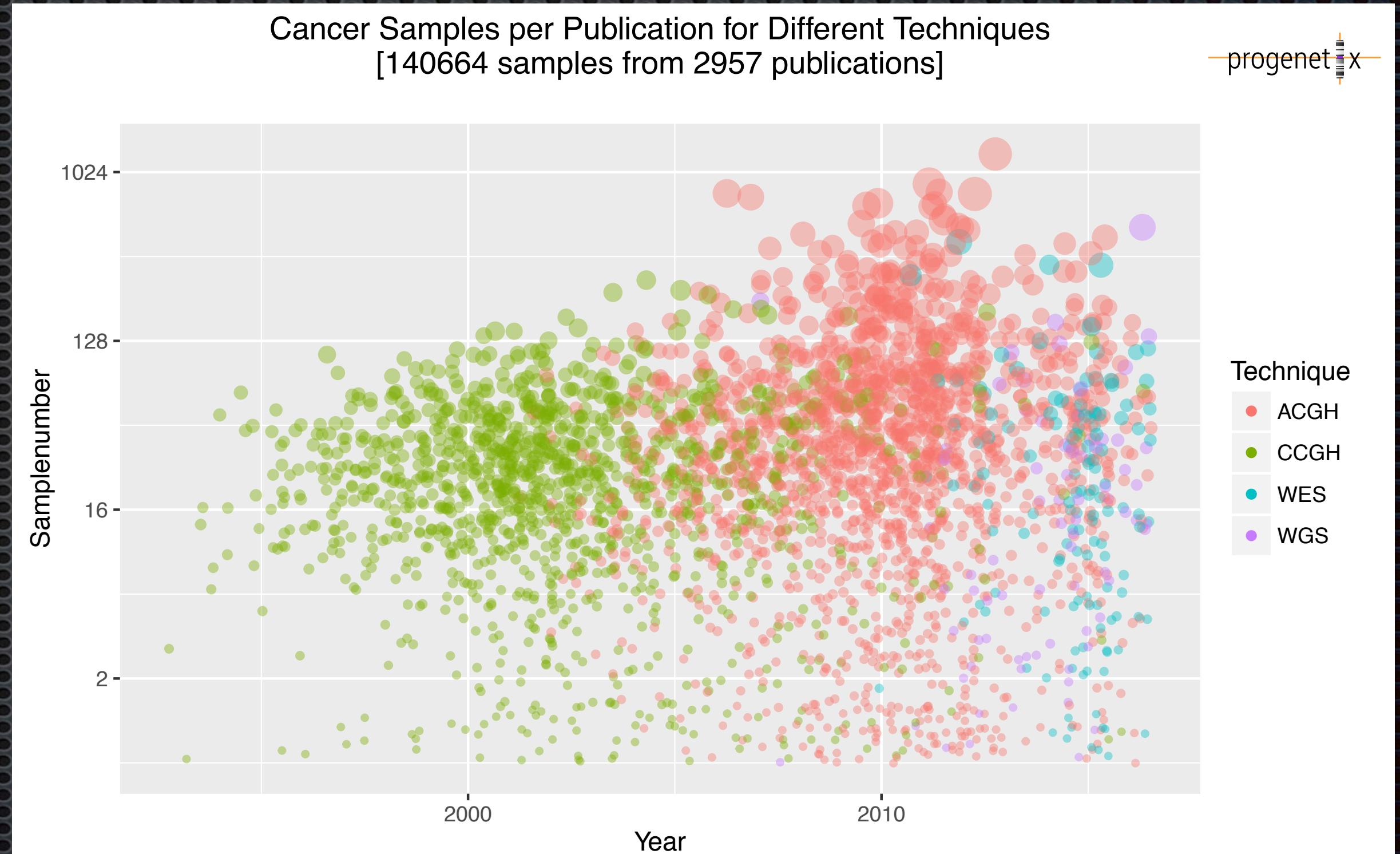
- a single nucleotide exchange Thymidine > Adenine leads to continuous RAF based activation of the MEK-ERK pathway
- BRAF V600E is a frequent mutation in >50% of malignant melanomas, but also CRC, lung ADC ...
- pharmacologic block of B-Raf (e.g. through **Vemurafenib**)



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

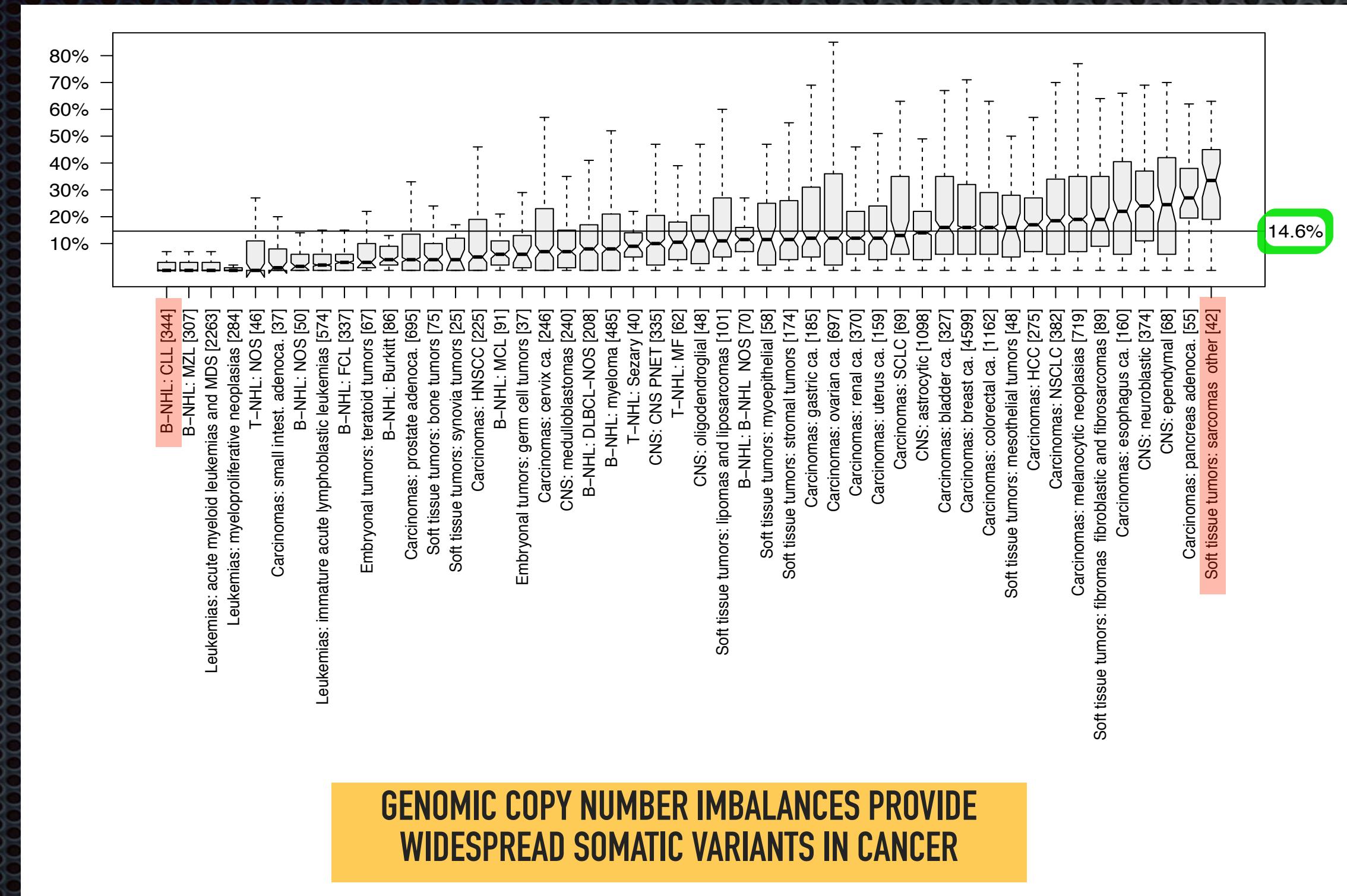
Molecular Cytogenetic & Sequencing Studies for **Whole Genome Profiling** in Cancer

- 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)



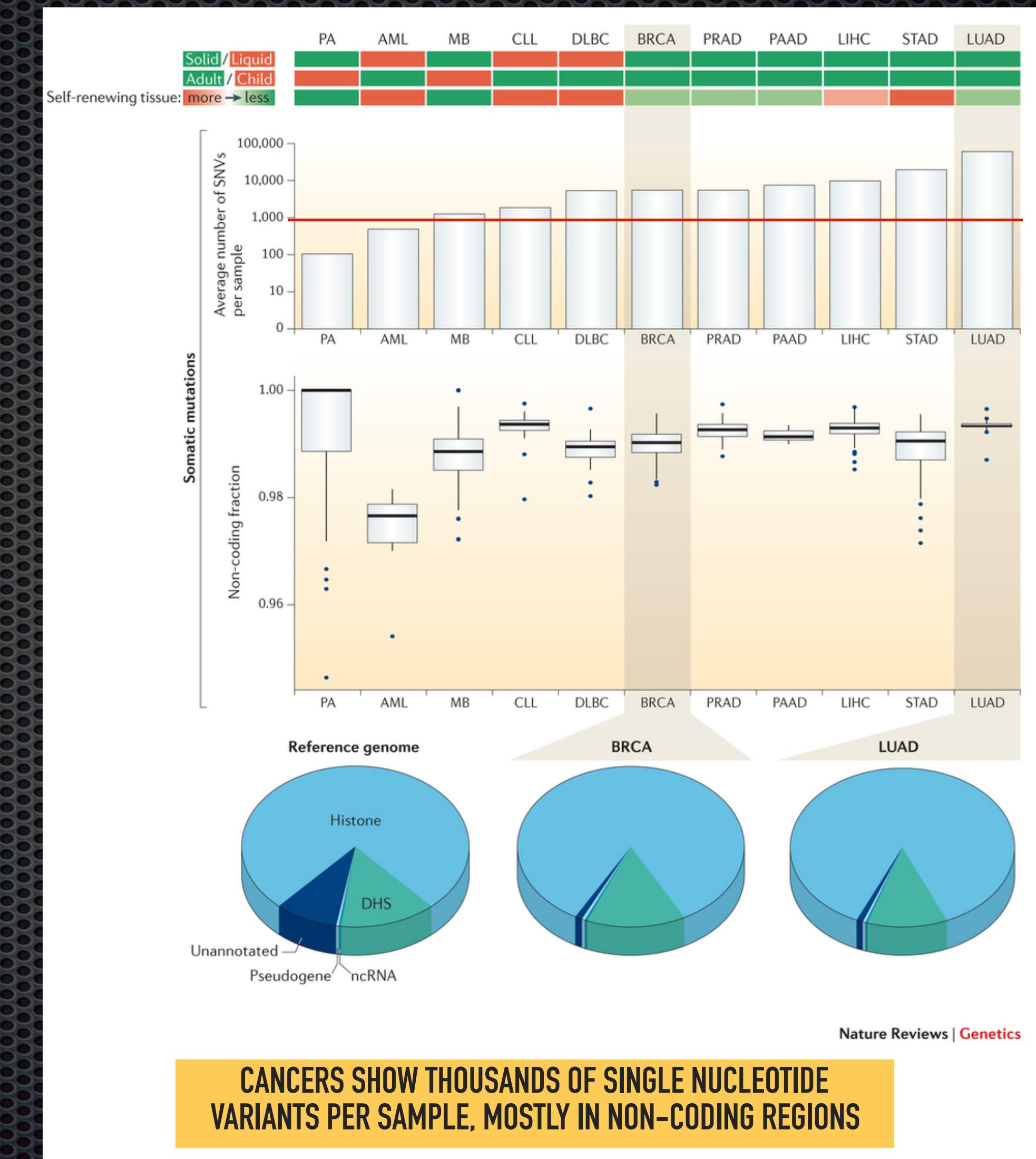
Overview of publications reporting whole-genome screening analysis of cancer samples, by molecular-cytogenetic or genome sequencing methods. The data represents articles assessed for the progenetix.org cancer genome data resource (M. Baudis, 2001-2016)

Quantifying Somatic Mutations In Cancer



GENOMIC COPY NUMBER IMBALANCES PROVIDE WIDESPREAD SOMATIC VARIANTS 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



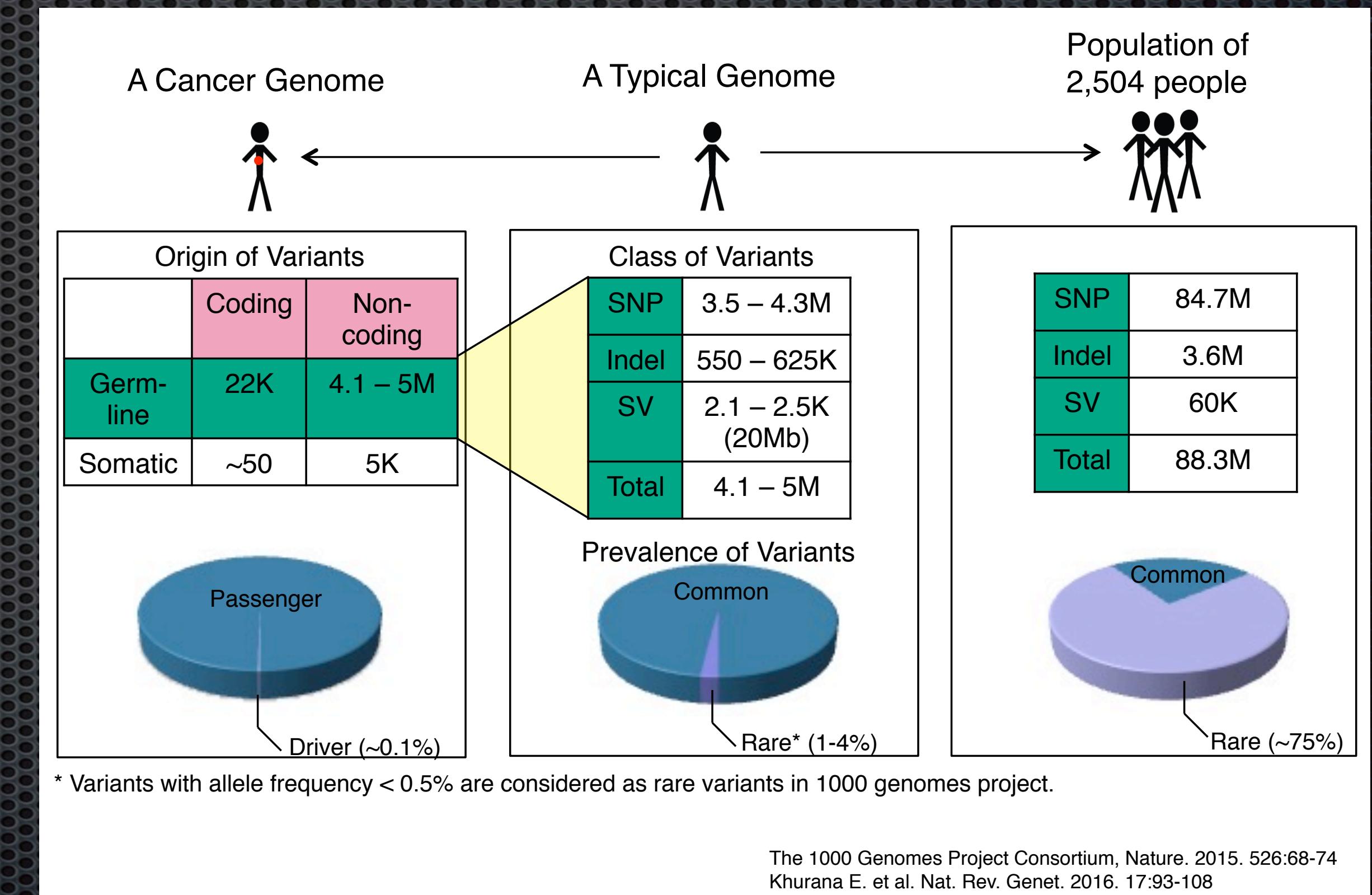
Pan-Cancer Analysis of Whole Genomes (PCAWG) data show widespread mutations in non-coding regions of cancer genomes (Khurana et al., Nat. Rev. Genet. (2016))

Personalized health & The trouble with human genetic variation



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

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

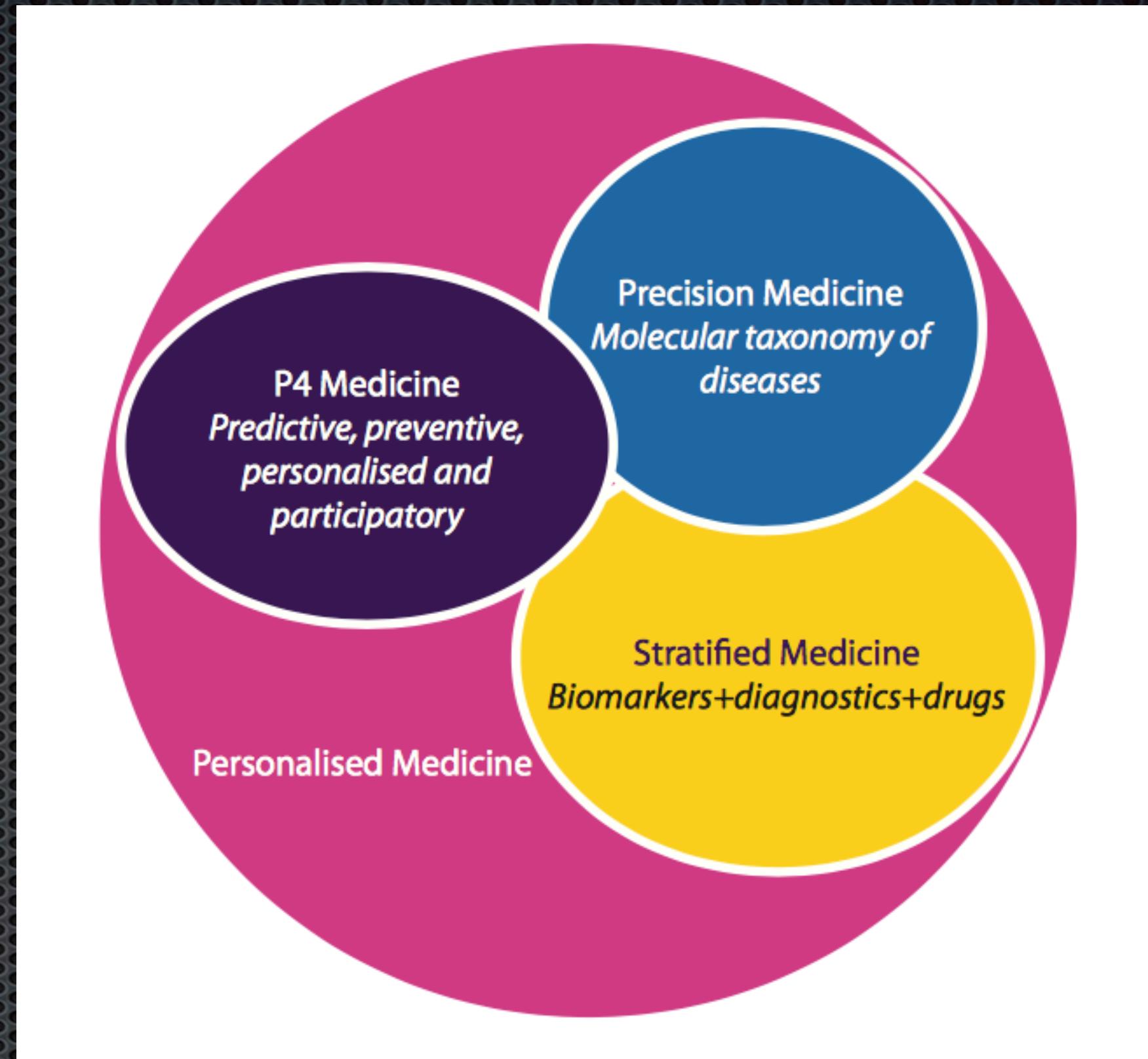
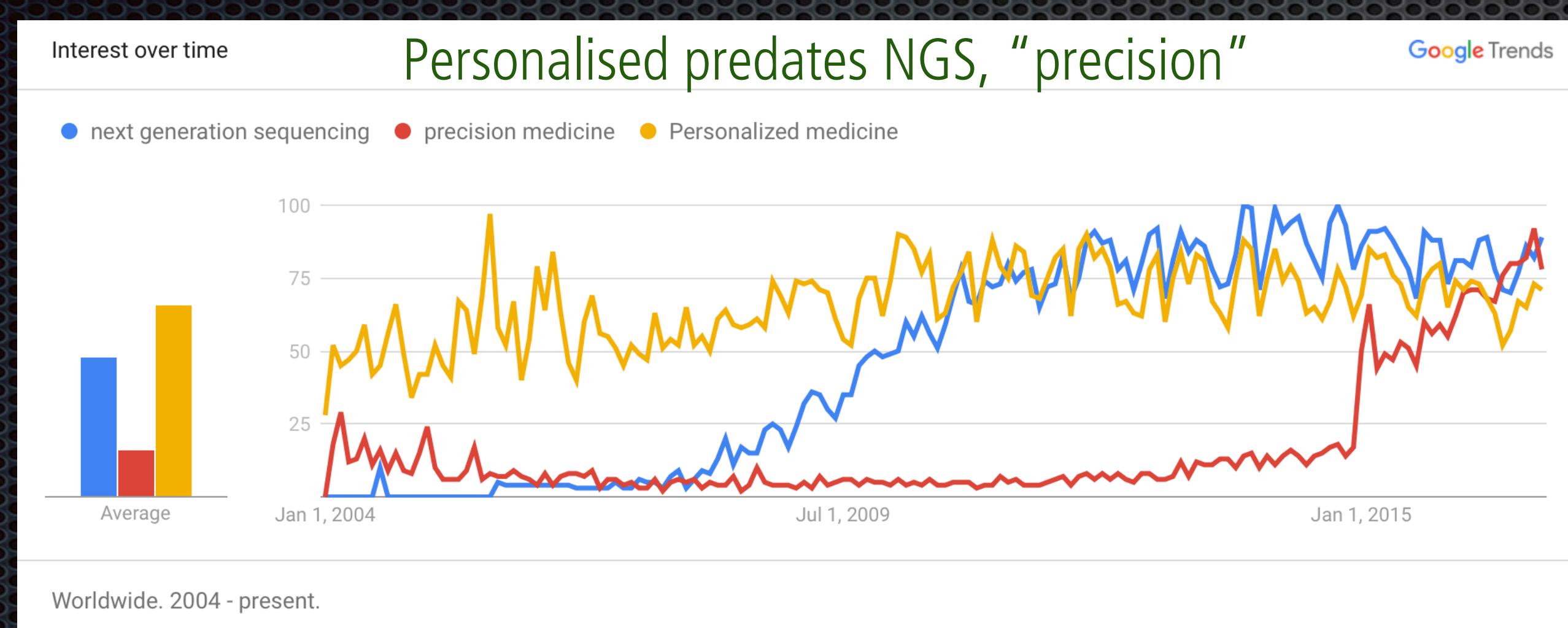


Graphic adapted from Mark Gerstein (GersteinLab.org; @markgerstein)

Genomic Background + Disease Parameters
Personalised Medicine

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 and individually targeted therapies.

Genome analyses at the core of Personalised Medicine™

Susceptibility, Pharmacogenomics, Classification, 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}, 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



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: iwatson@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>

doi:10.1111/pcn.12128

PCN Frontier Review

Psychiatry and Clinical Neurosciences

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*}

¹National Institute for Physiological Sciences, ²Restart Postdoctoral Fellowship of Japan Society for the Promotion of Science, Tokyo, and ³Department of Psychiatry, Nagoya University Graduate School of Medicine, Nagoya, Japan

RESEARCH ARTICLE

Open Access

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,*}

NEURODEVELOPMENT

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

Mustafa Sahin* and Mriganka 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

RESEARCH ARTICLE

Open Access

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

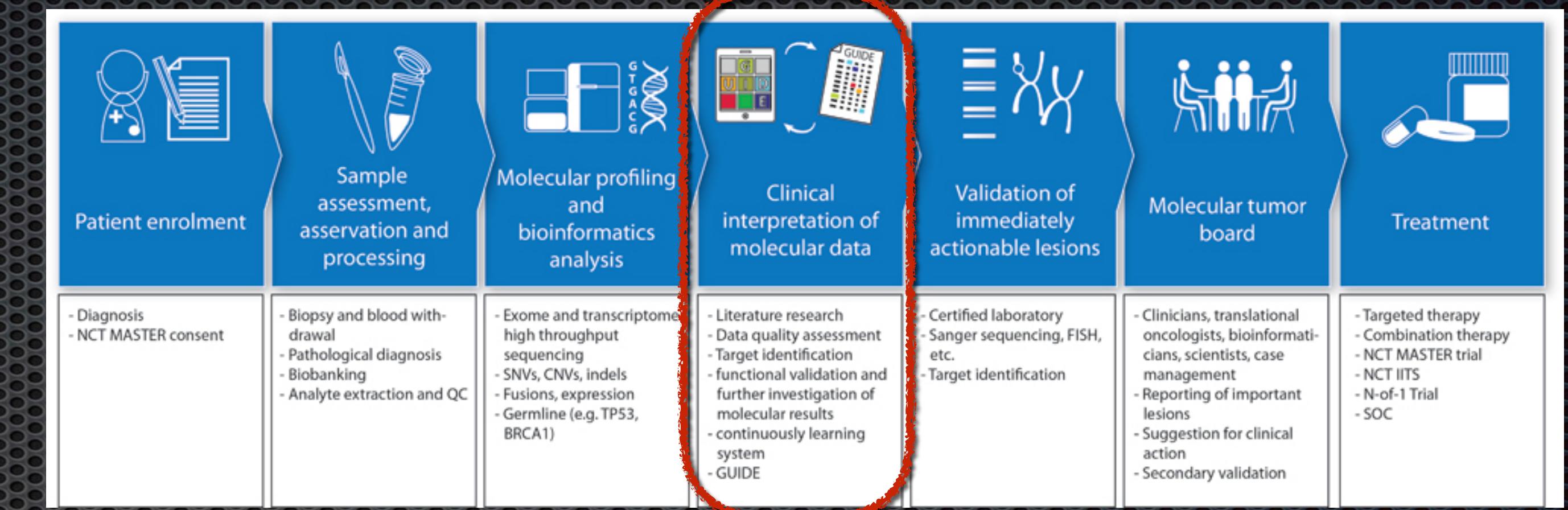
Manfred Beutel^{1,5*}, Philip Zimmermann², Michael Baudis³, Nicole Bruni⁴, Peter Bühlmann⁴, Oliver Laule², Van-Duc Luu¹, Wilhelm Gruissem², Peter Schraml^{1,*} and Holger Moch¹

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. Boehml¹, 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 Hoffner^{1,3}, John Prensner^{1,3}, Ted Liefeld¹, Qing Gao¹, Derek Yecies³, Sabina Signoretto^{3,4}, Elizabeth Maher¹⁰, Frederic J. Kaye¹¹, Hidefumi Sasaki¹², Joel E. Tepper¹³, Jonathan A. Fletcher⁴, Josep Tabernero¹⁴, Jose 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 Okamoto³², 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}

Personalised Medicine in Cancer - A Genome Based Approach

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



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

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

Curated Variant Data Resources as Backbone of Personalised Cancer Therapy

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

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

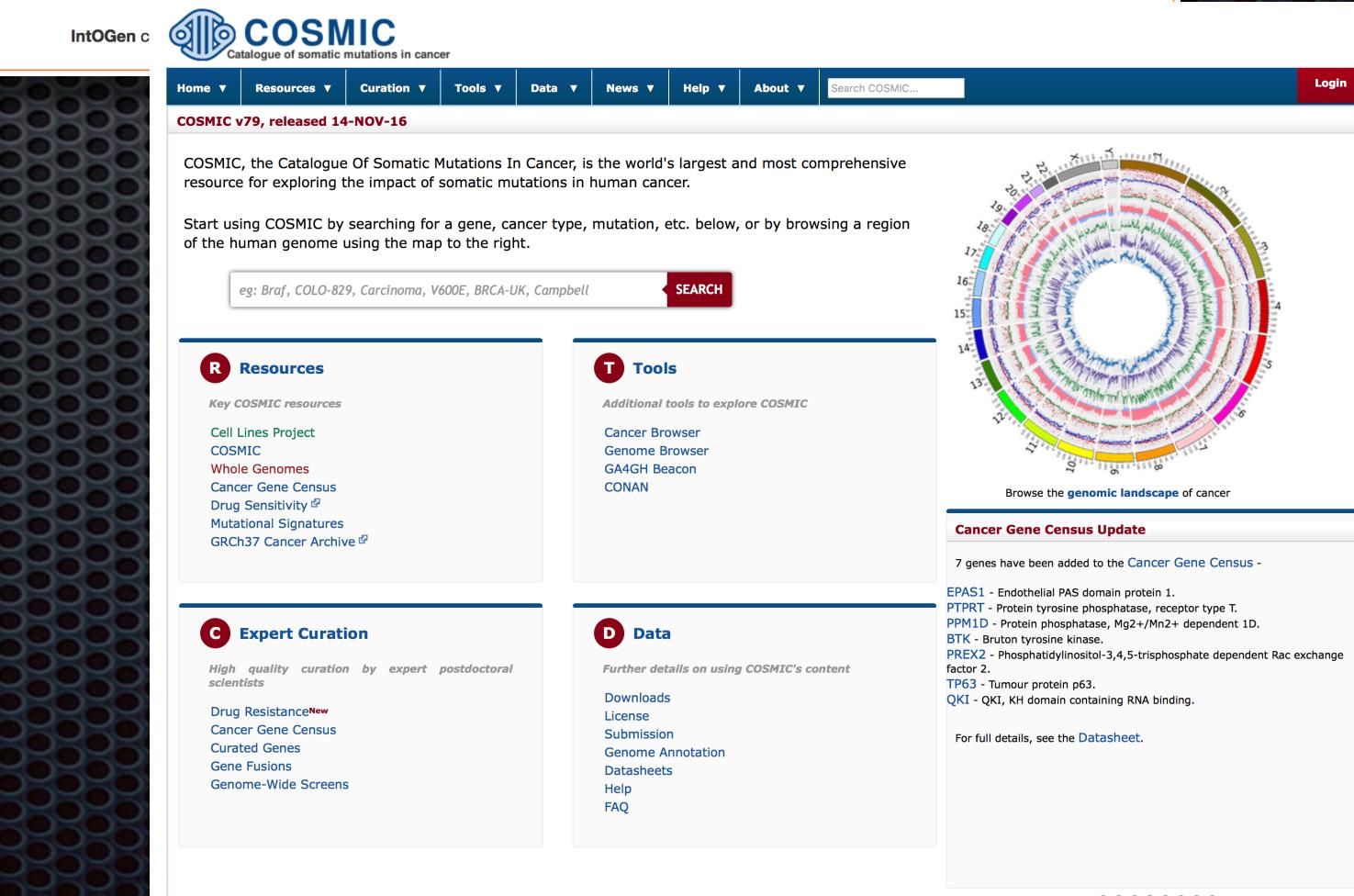
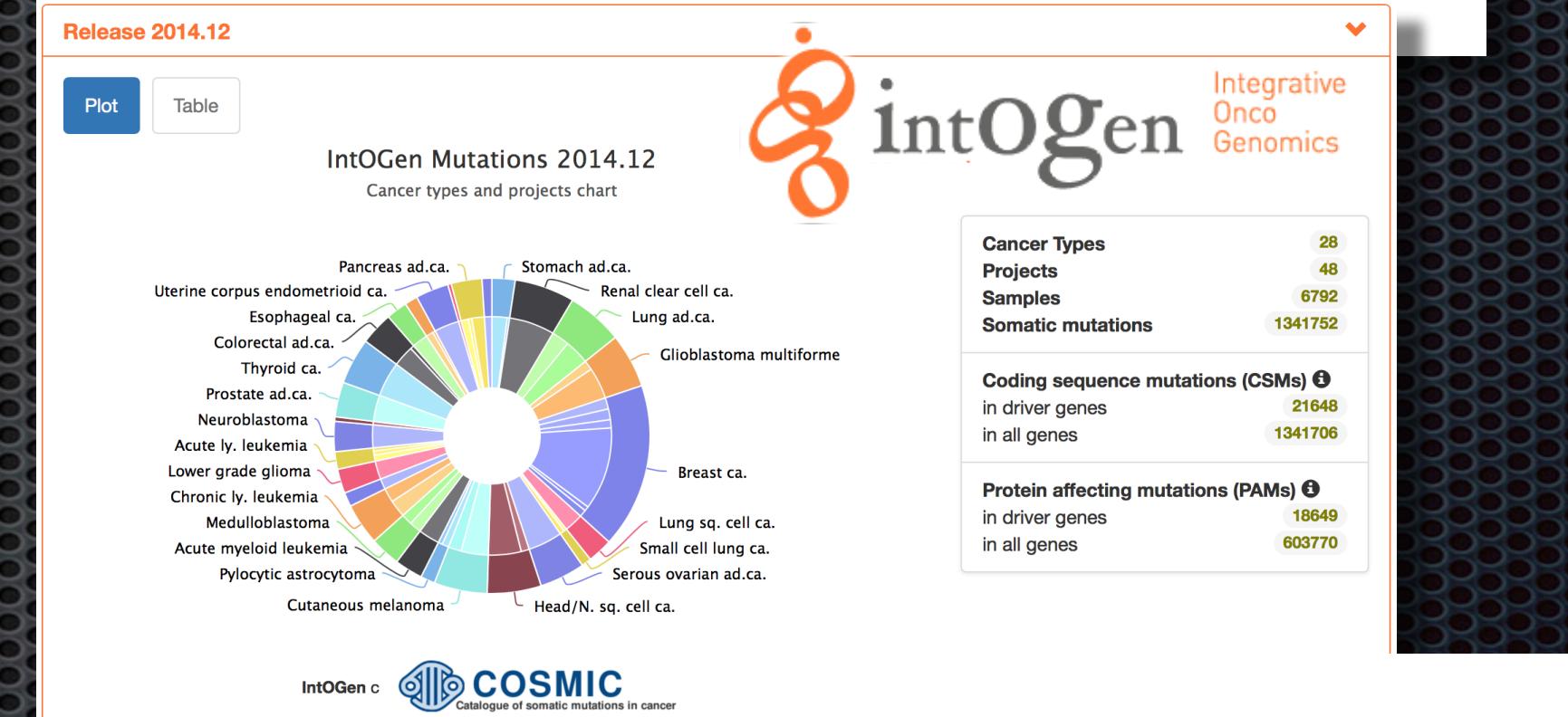
OncoKB Home About Team Levels of Evidence Actionable Genes Data Access News

OncoKB
Precision Oncology Knowledge Base
Annotation of Somatic Mutations in Cancer

418 Genes 3332 Variants 50 Tumor Types 71 Drugs

Search Gene

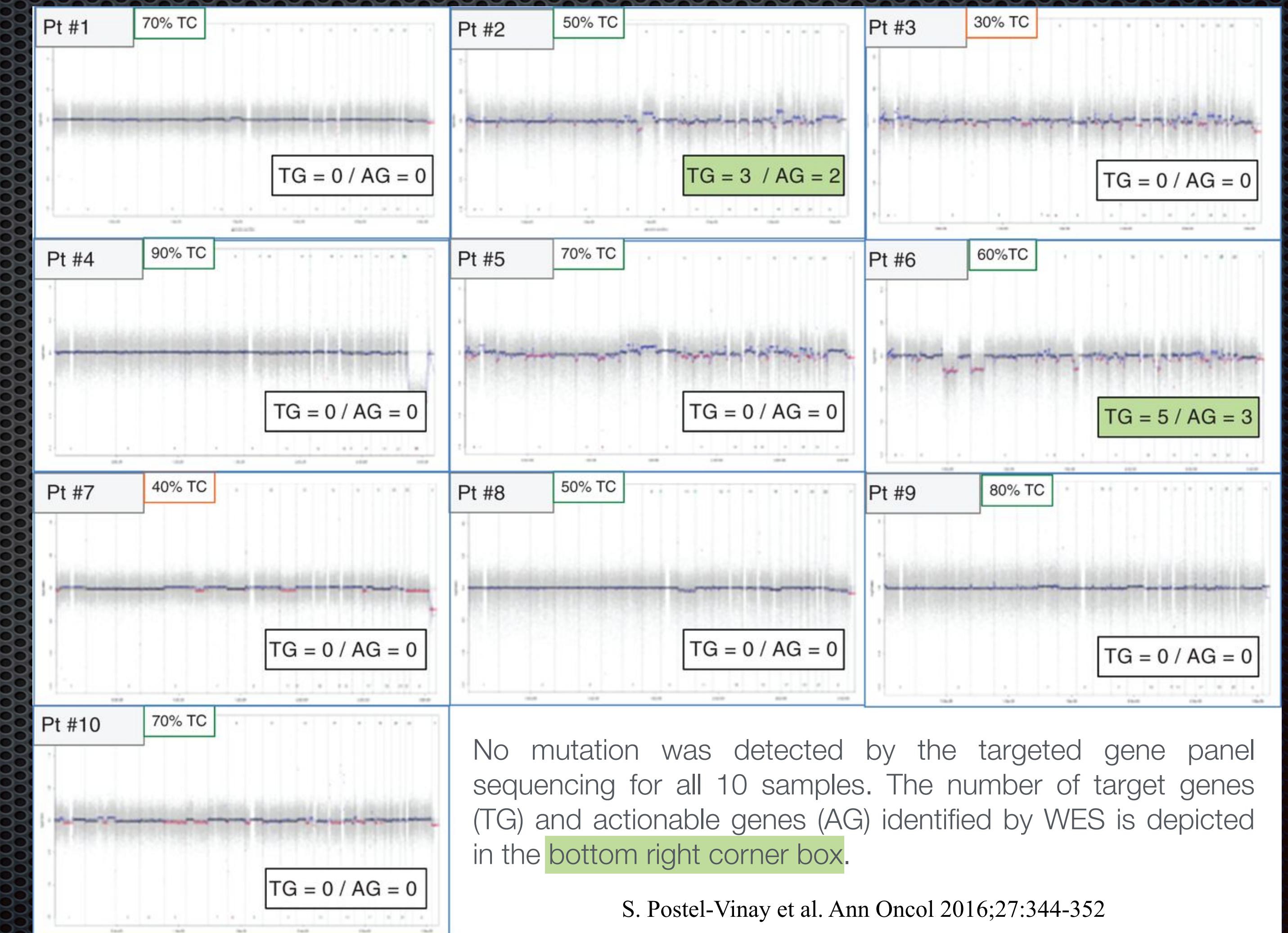
Level 1 FDA-approved Level 2 Standard-of-care Level 3 Clinical evidence



Drivers in tumours with apparent normal molecular profile on CGH and **targeted gene panel** sequencing: what is the added value of **WES**?

French “MOSCATO-01” trial (2011-2013)

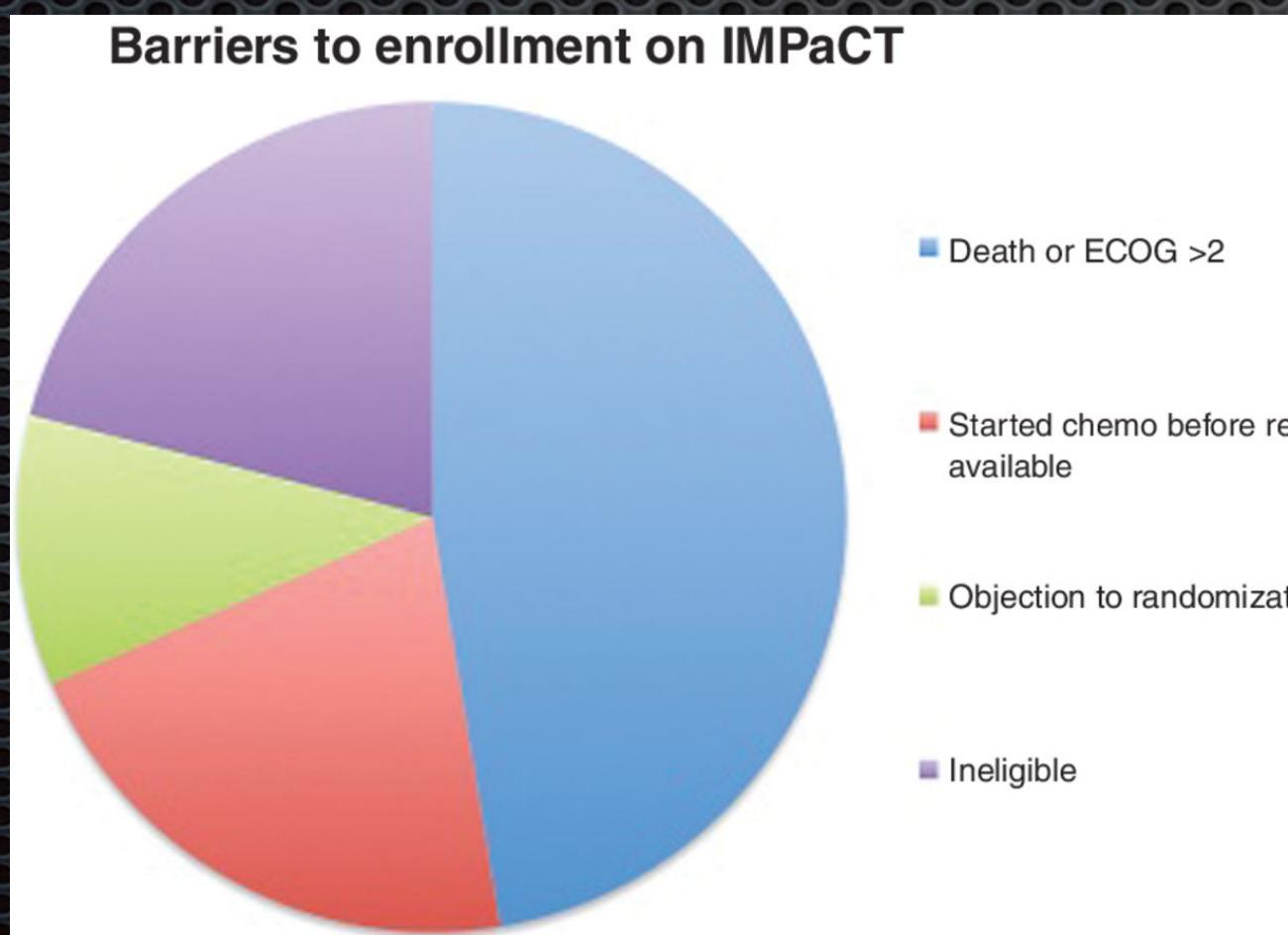
- samples **without** mutation in hot spot regions of a **panel** of 46 critical cancer-related genes were submitted to WES
- actionable mutations were identified in **2/10** randomly selected patients



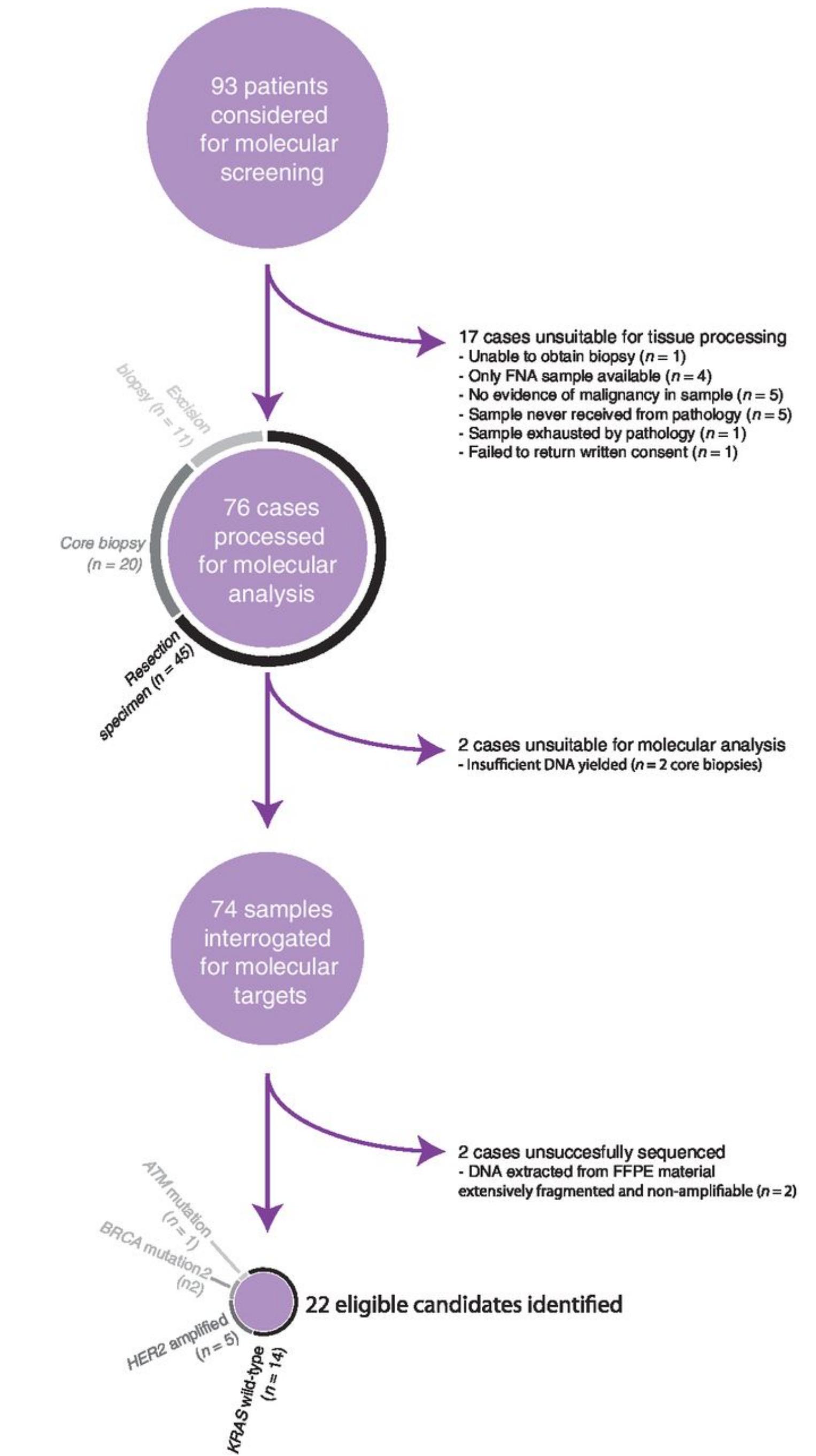
No mutation was detected by the targeted gene panel sequencing for all 10 samples. The number of target genes (TG) and actionable genes (AG) identified by WES is depicted in the bottom right corner box.

Personalised Cancer Therapy Needs More Targets

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



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



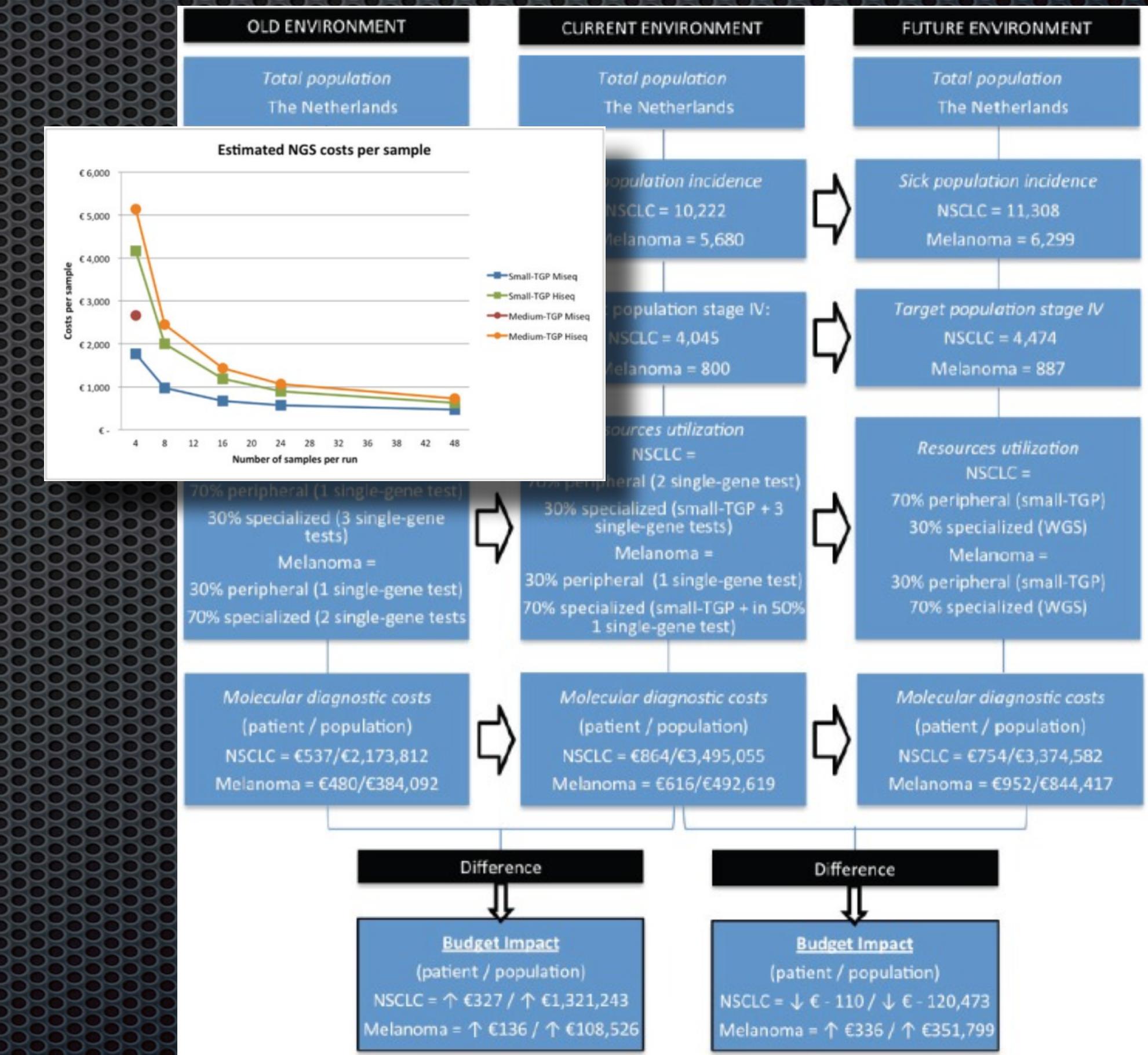
Next-generation sequencing strategies in malignancies

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

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

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

WGS use is expected to reduce the cost of additional tests

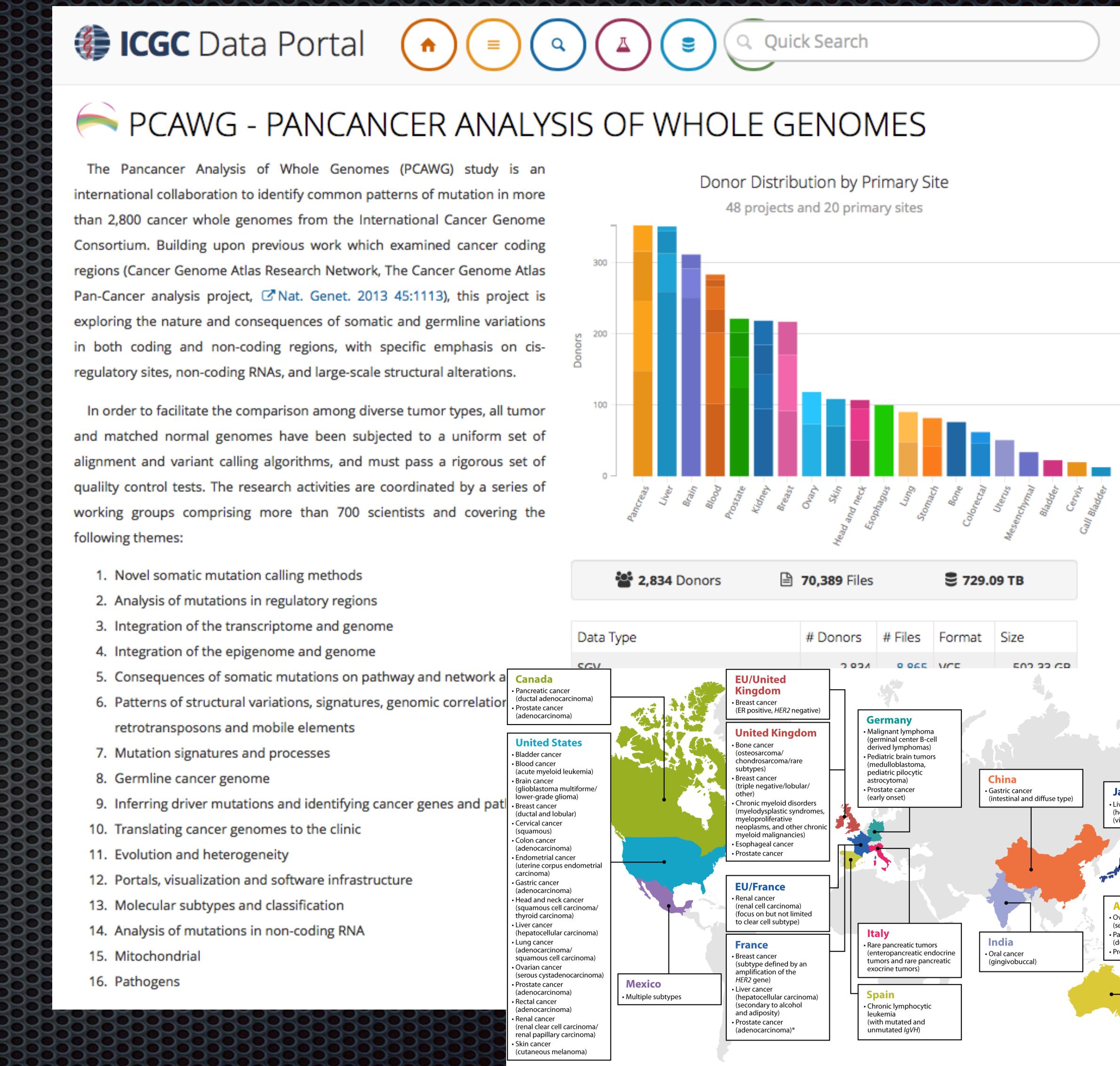


Estimated costs per sample for small- and medium TGP and WGS; van Amerongen et al., Ecancermedicalscience. 2016

Genome Data Access in Cancer Mining for New Knowledge

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

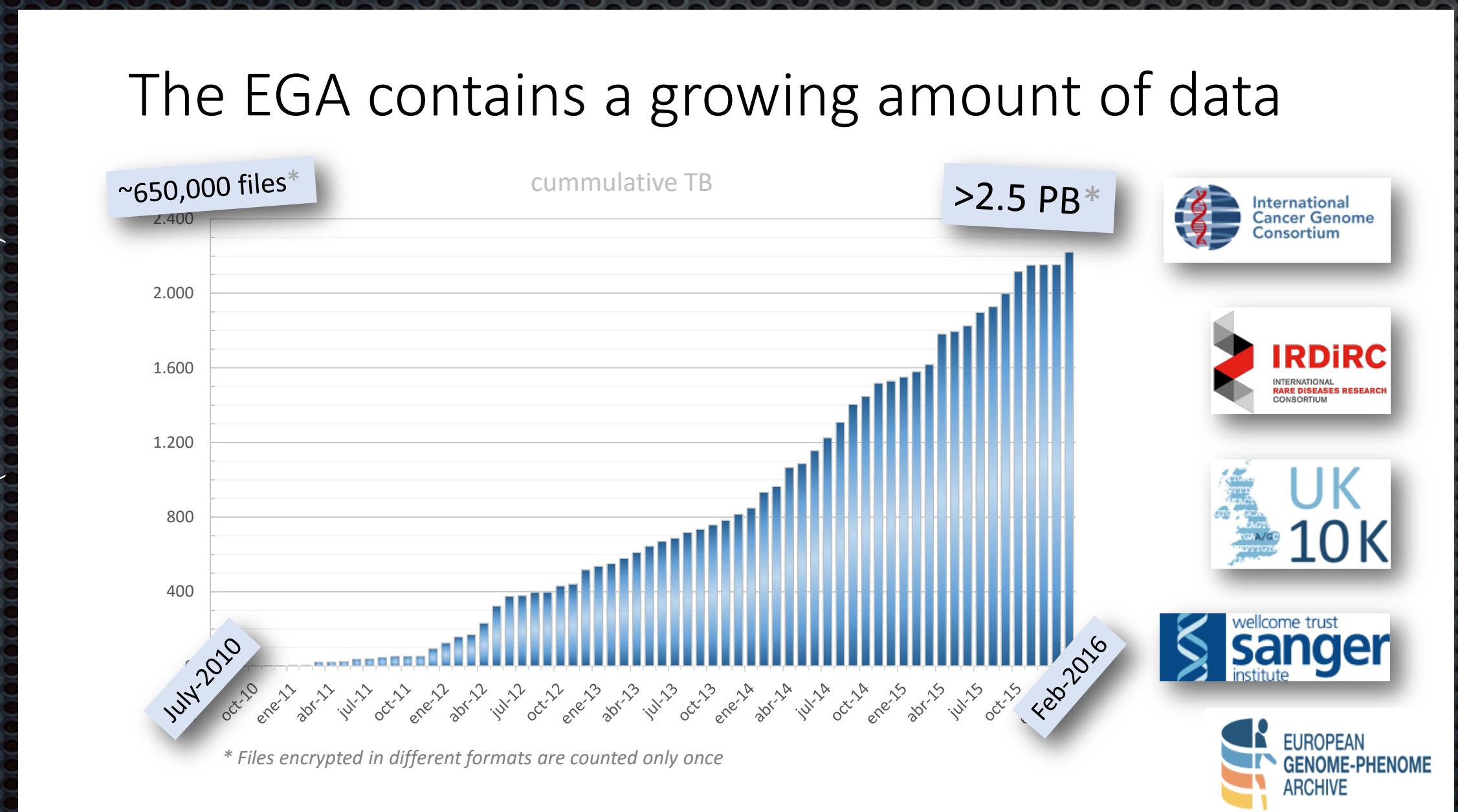
- the International Cancer Genome Consortium (ICGC) as leading example of deep analysis of multiple cancer entities
- international collaboration of leading research centers for each of ~20 tumor types
- limitations:
 - focus on prominent cancer types w/ limited representation of rare entities
 - data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



Wong KM, et al. 2011.
Annu. Rev. Genomics Hum. Genet. 12:407–30

Genome Datasets: Rapid Growth, Limited Access

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



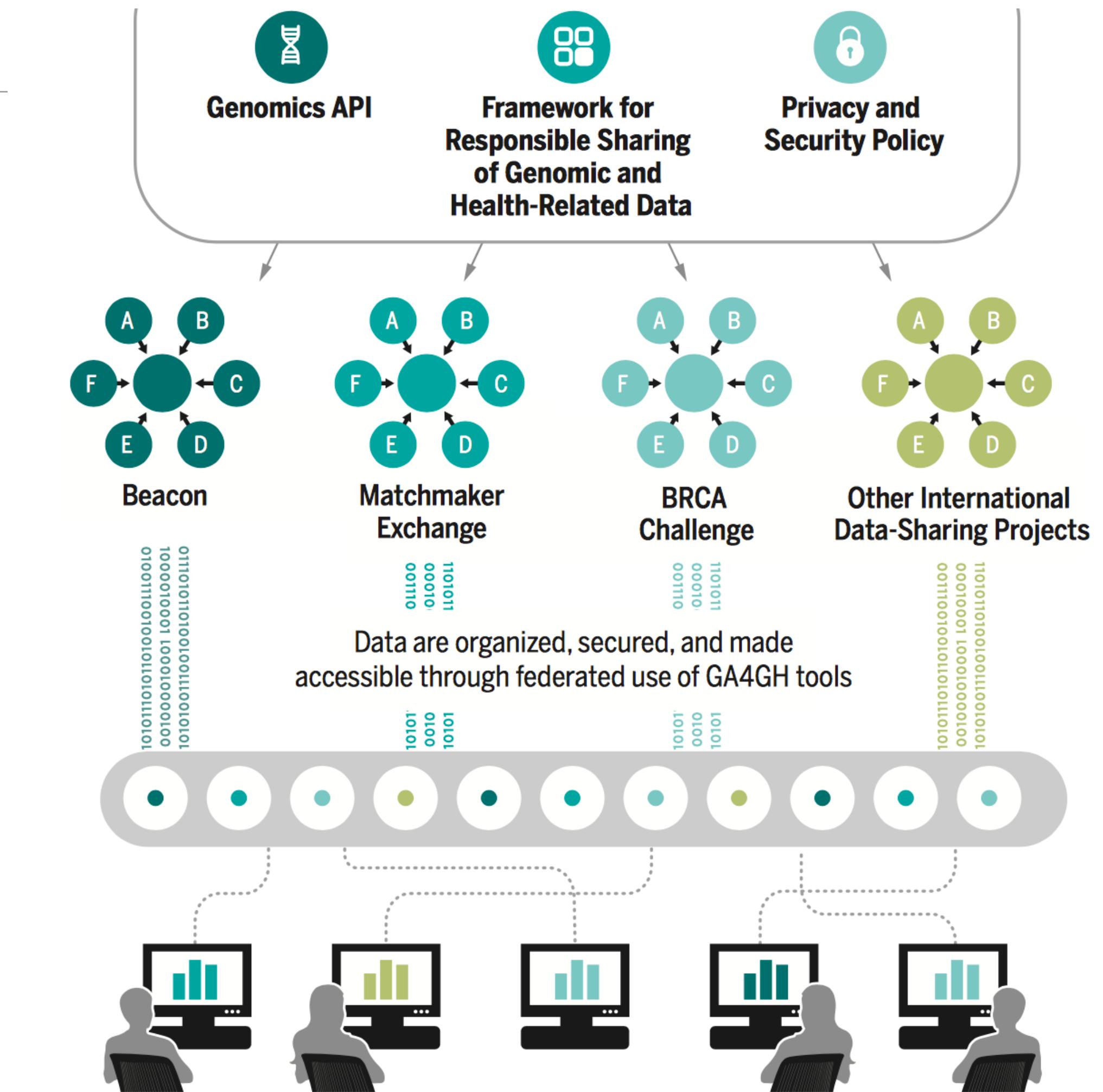
genome data is stored in an increasing number of institutional and core repositories, with **incompatible data** structures and **access** policies

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.



The mission of the Global Alliance for Genomics and Health is to accelerate progress in human health by helping to establish a **common framework** of harmonized approaches to enable **effective and responsible** sharing of **genomic and clinical data**, and by catalyzing data sharing projects that drive and demonstrate the value of data sharing.



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 Dyke¹, Edward S Dove² and Bartha M Knoppers¹

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.

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

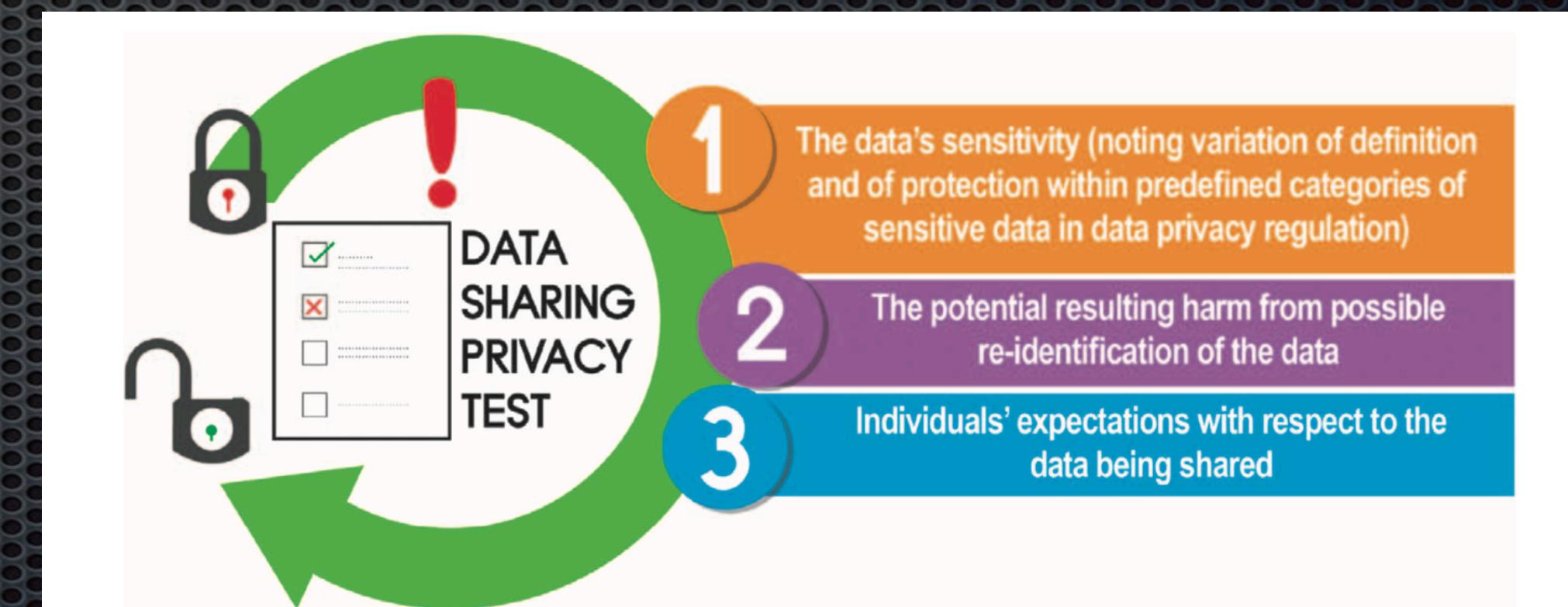


Figure 1. The three steps of a Data Sharing Privacy Test to distinguish degrees of data sensitivity within categories of data recognised as 'sensitive'.

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	Biologisches Material und genetische Daten	Nicht-genetische Daten
Personenbezug		
Unverschlüsselt (identifizierend)	Information + Einwilligung in jedes einzelne Forschungsprojekt	Information über Weiterverwendung für zukünftige noch unbestimmte Forschungsprojekte + Generalkonsent für Forschungszwecke
Verschlüsselt	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
Anonymisiert	Genetische Daten: Information über Weiterverwendung für zukünftige noch unbestimmte Forschungszwecke + über Möglichkeit Weiterverwendung abzulehnen > Widerspruchsrecht Proben: 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.
(Source: SAMW)

Consent Codes: Upholding Standard Data Use Conditions

Stephanie O. M. Dyke^{1*}, Anthony A. Philippakis², Jordi Rambla De Argila^{3,4}, Dina N. Paltoo⁵, Erin S. Luetkemeier⁵, Bartha M. Knoppers¹, Anthony J. Brookes⁶, J. Dylan Spalding⁷, Mark Thompson⁸, Marco Roos⁸, Kym M. Boycott⁹, Michael Brudno^{10,11}, Matthew Hurles¹², Heidi L. Rehm^{2,13}, Andreas Matern¹⁴, Marc Fiume¹⁵, Stephen T. Sherry¹⁶



Consent Codes		
Name	Abbreviation	Description
Primary Categories (I^{IV})		
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.
Secondary Categories (II^{IV}) (can be one or more extra conditions, in addition to I ^{IV} 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).
Requirements		
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.

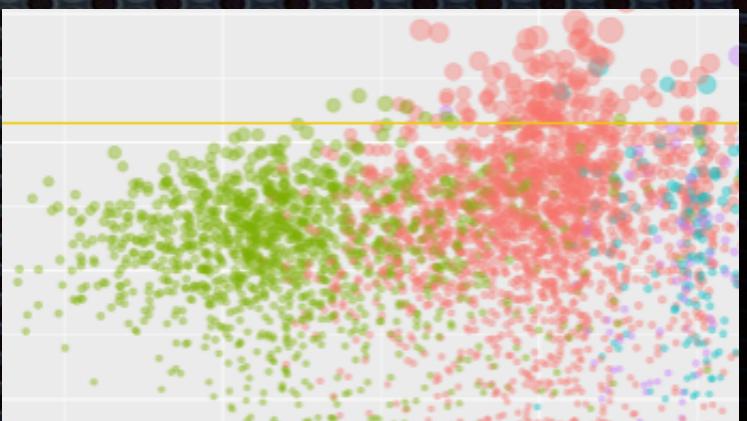
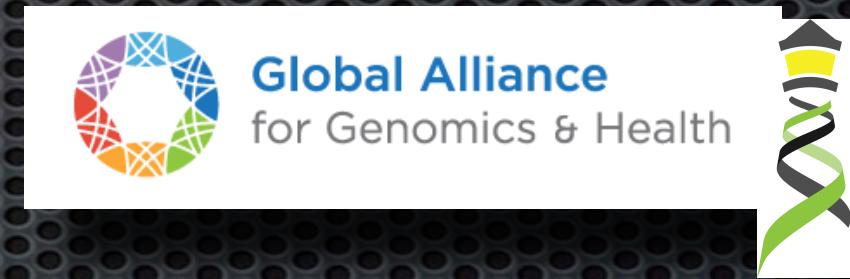
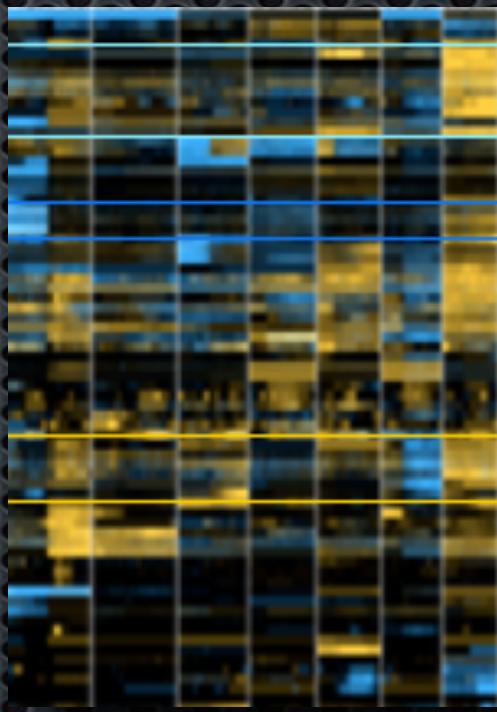
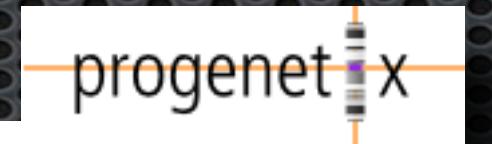
SOM Dyke, et al. Consent Codes: Upholding Standard Data Use Conditions. *PLoS Genetics* 12(1): e1005772. <http://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1005772>

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

Cancer Genome Data & Beyond

Contributing our share

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



Reference Resources for Cancer Genome Profiling

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



arrayMap

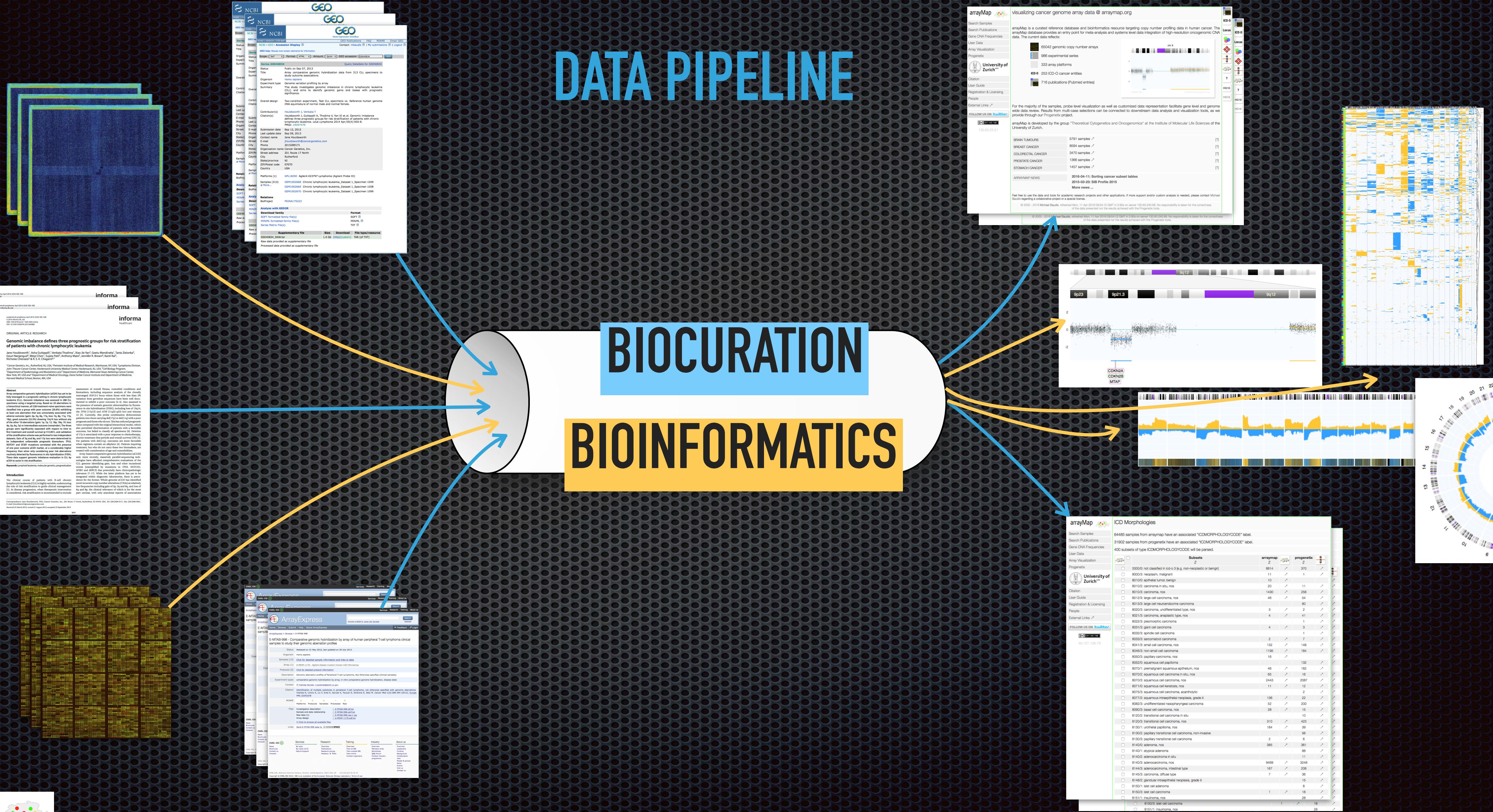


techniques	cCGH, aCGH, WES, WGS	aCGH (+?)
scope	sample (e.g. combination of several experiments)	experiment
content	>31000 samples	>60000 arrays
raw data presentation	no (link to sources if available)	yes (raw, log2, segmentation if available)
per sample re-analysis	no; supervised result (mostly as provided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

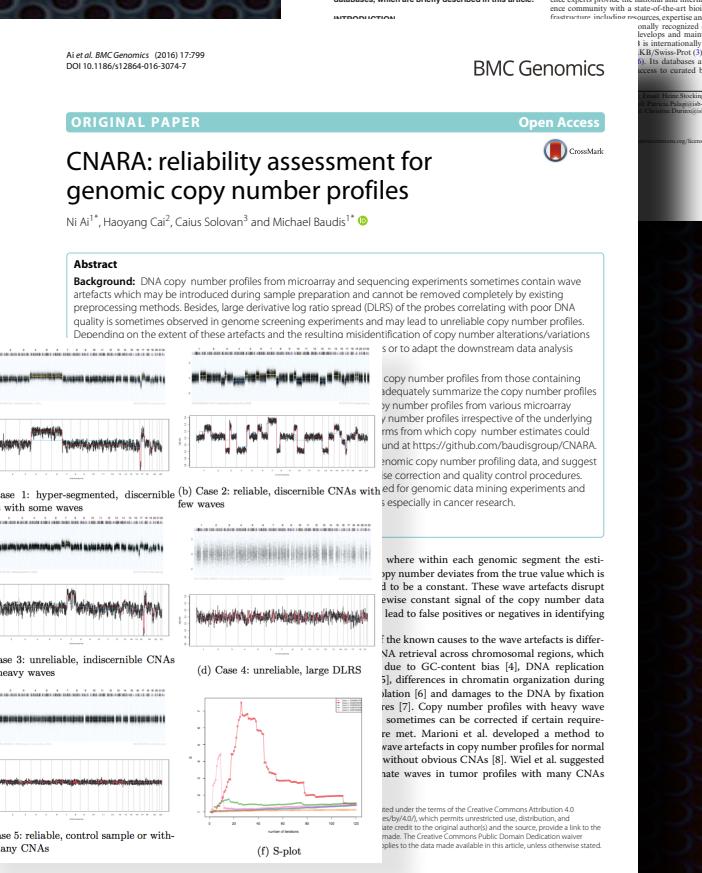
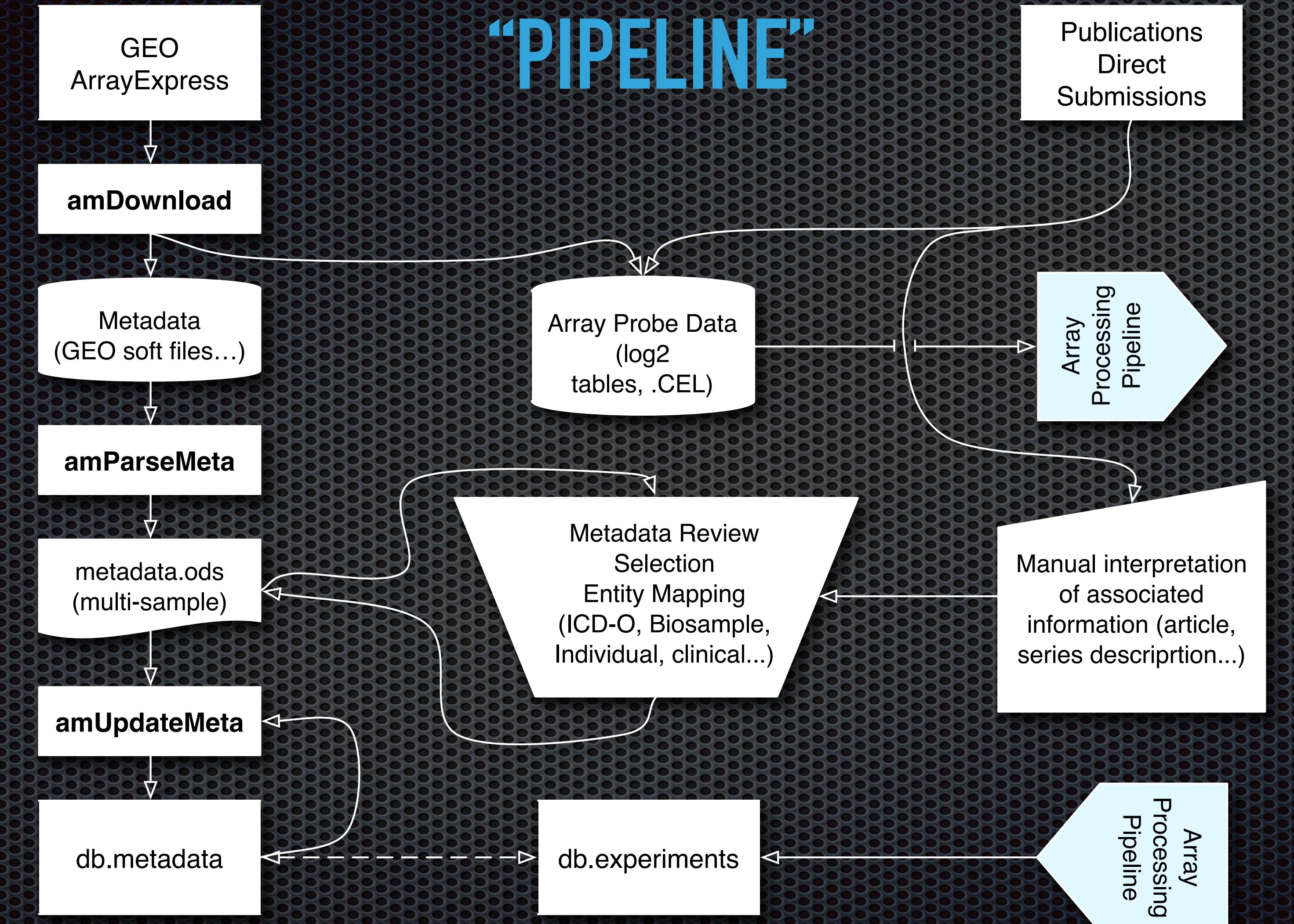
DATA PIPELINE

BIOCURATION

BIOINFORMATICS



ARRAYMAP DATA



arrayMap

The largest resource for copy number variation data in cancer

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)

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.

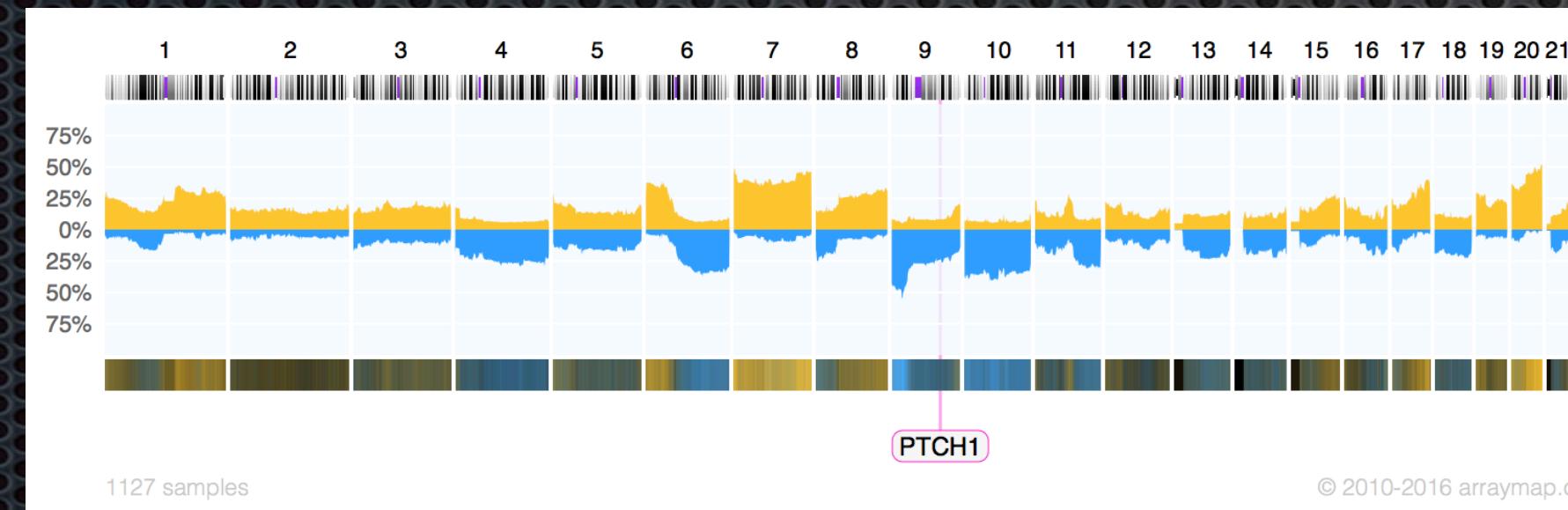
© 2000 - 2016 Michael Baudis, refreshed Mon, 19 Sep 2016 10:20:09 GMT in 6.87s on server 130.60.240.68. No responsibility is taken for the correctness of the data presented nor the results achieved with the Progenetix tools.



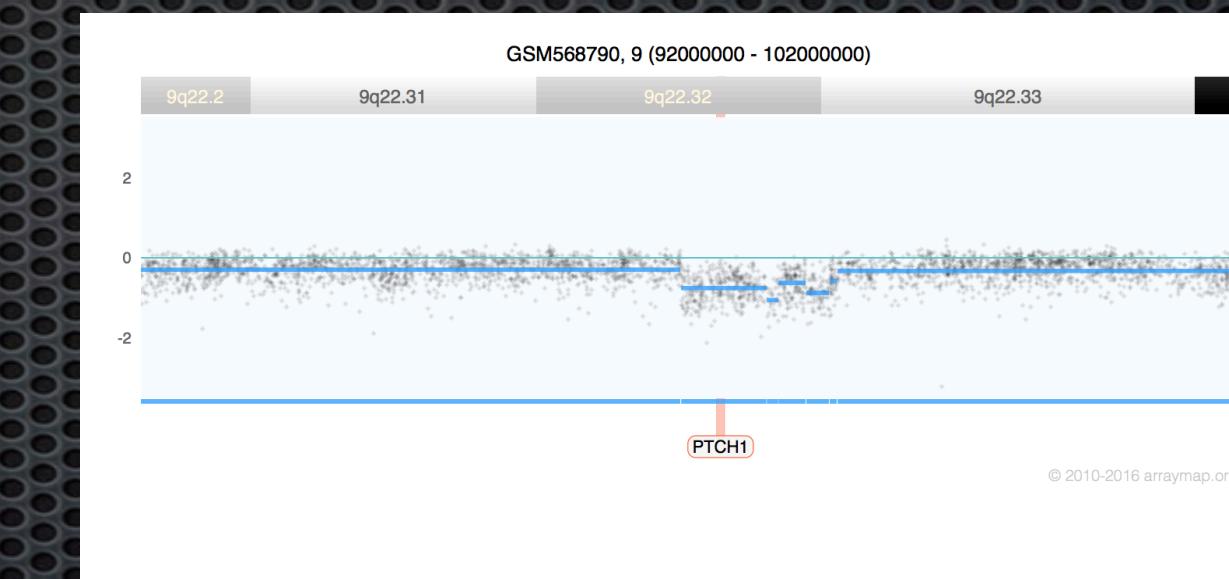
Rare Events & Hidden Therapeutic Options?

Example: PTCH1 deletions in malignant melanomas

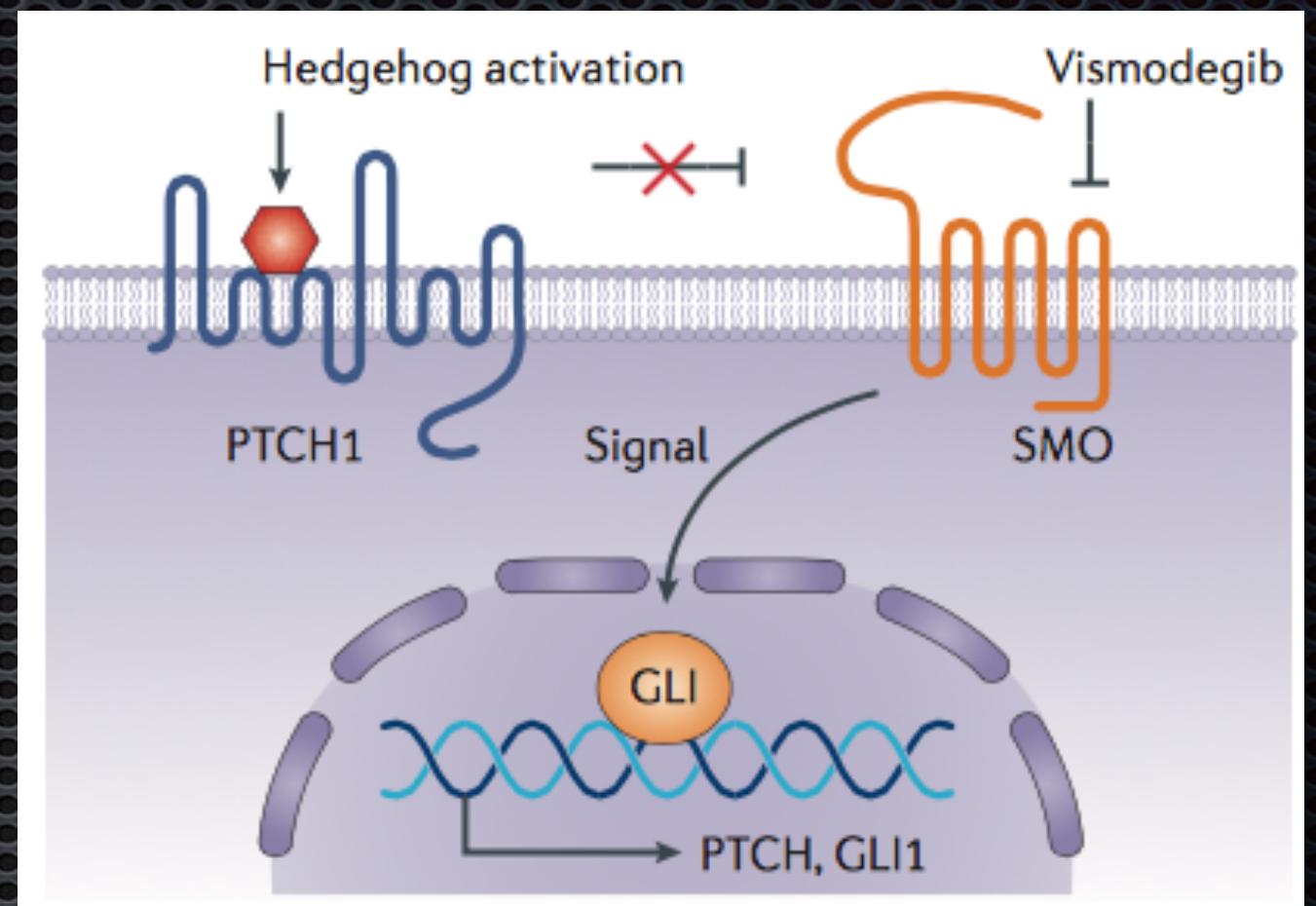
- PTCH1 is a actionable tumor suppressor gene, which has been demonstrated in e.g. basalomas and medulloblastomas
- analysis of 1127 samples from 26 different publications could identify **focal** deletions in 4 samples
- a current project addresses the focal involvement of all mapped genes, in >50'000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



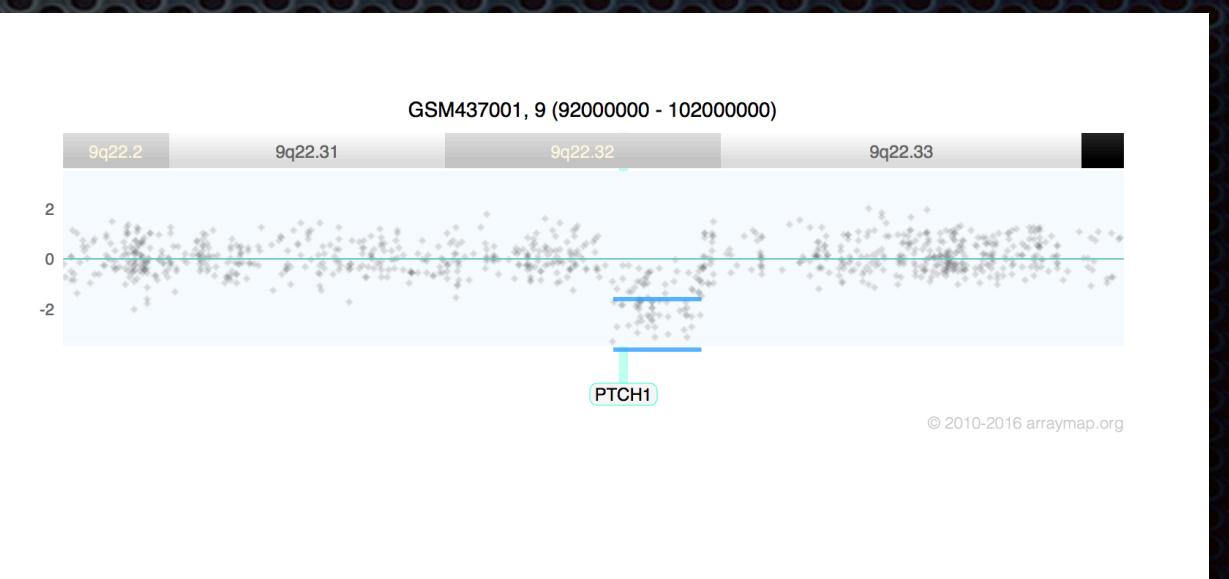
Summary of somatic copy number aberrations from the analysis of 1127 genome profiles of malignant melanomas, collected in our arraymap.org cancer genome resource. While PTCH1 does not represent a deletion hotspot, the genomic locus is part of larger deletions in ~25% of melanoma samples.



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

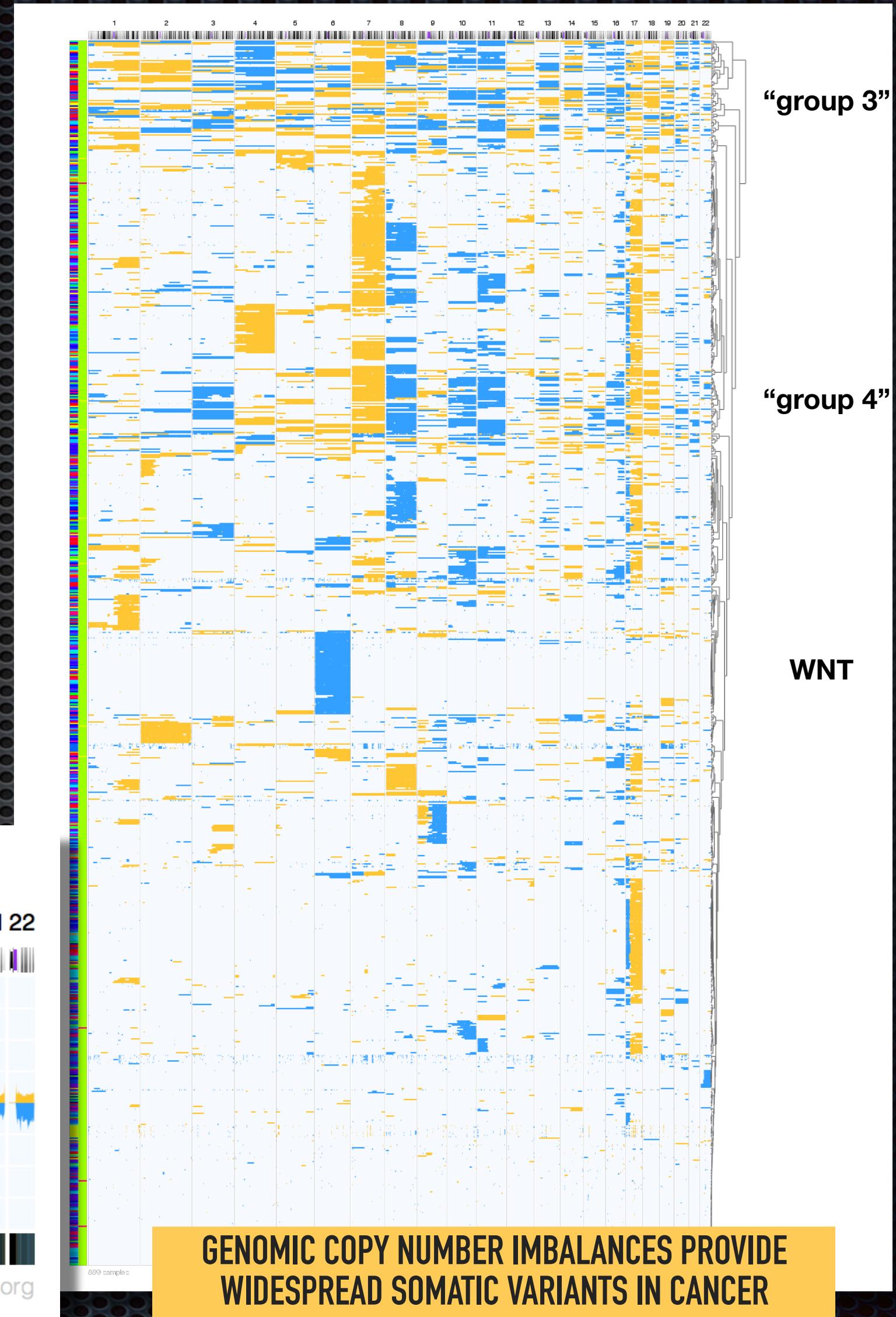
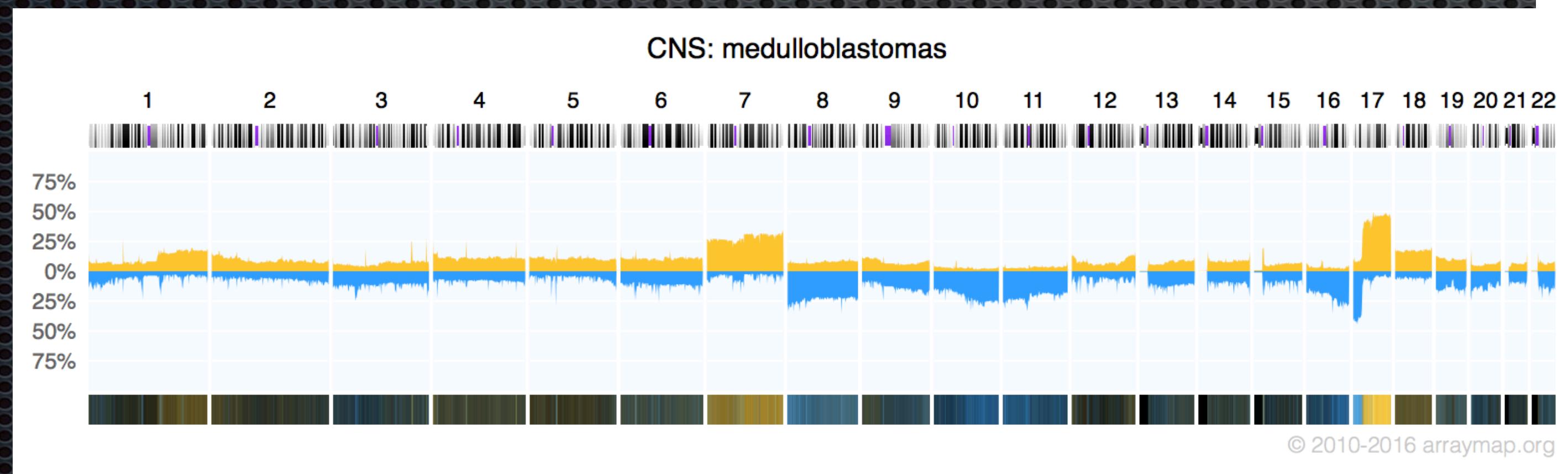


In its normal function, PTCH1 is a tumor suppressor gene in the sonic hedgehog pathway and inhibits SMO driven transcriptional activation. A loss of PTCH1 function (mutation, deletion) can be mitigated through drugs antagonistic to SMO activation.



Somatic Mutations In Cancer: Patterns

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

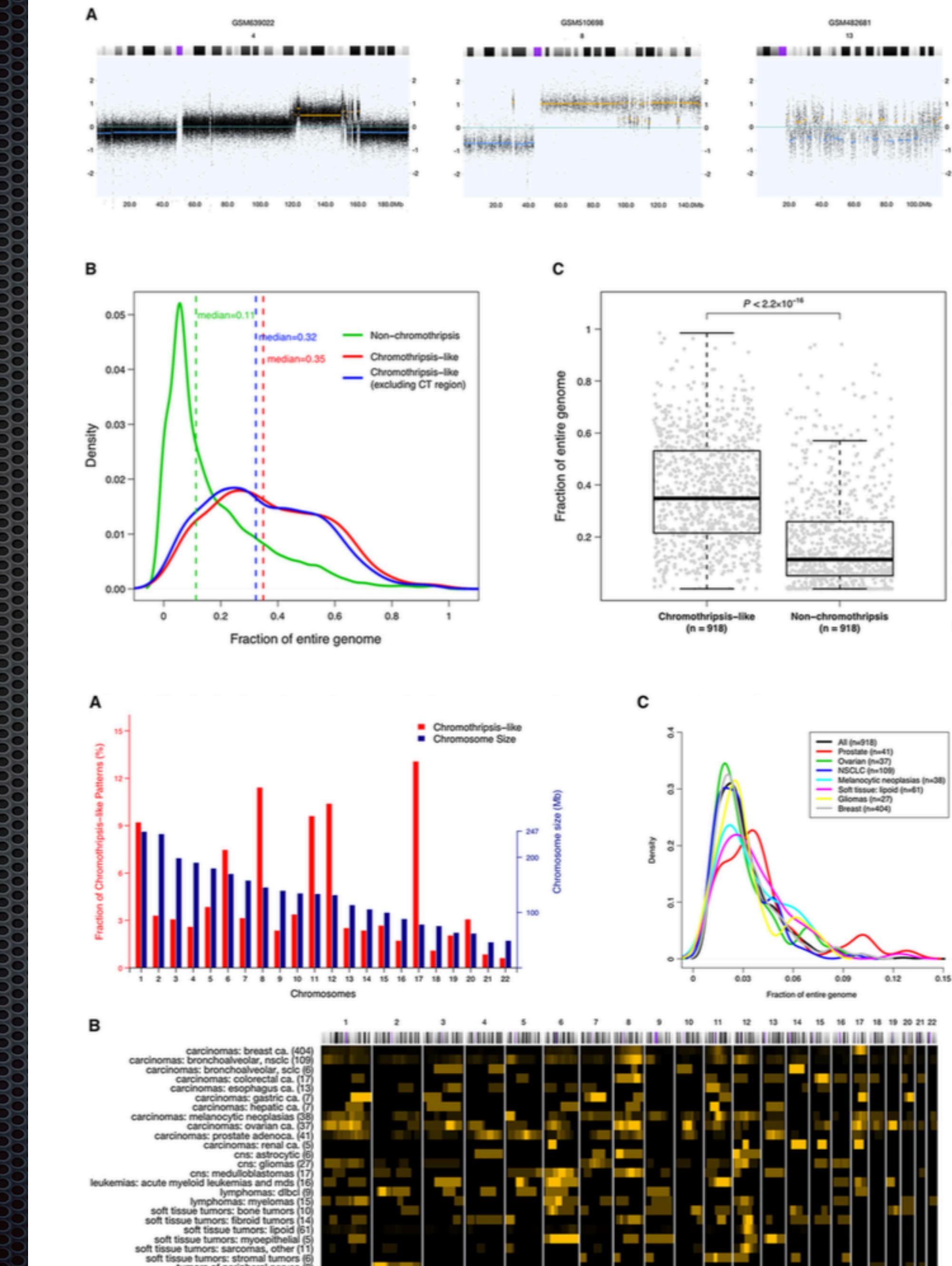


A genomic copy number histogram for malignant medulloblastomas, the most frequent type of pediatric brain tumors, displaying regions of genomic duplications and deletions. These can be decomposed into individual tumor profiles which segregate into several clusters of related mutation patterns with functional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

Somatic Mutations In Cancer: Patterns II

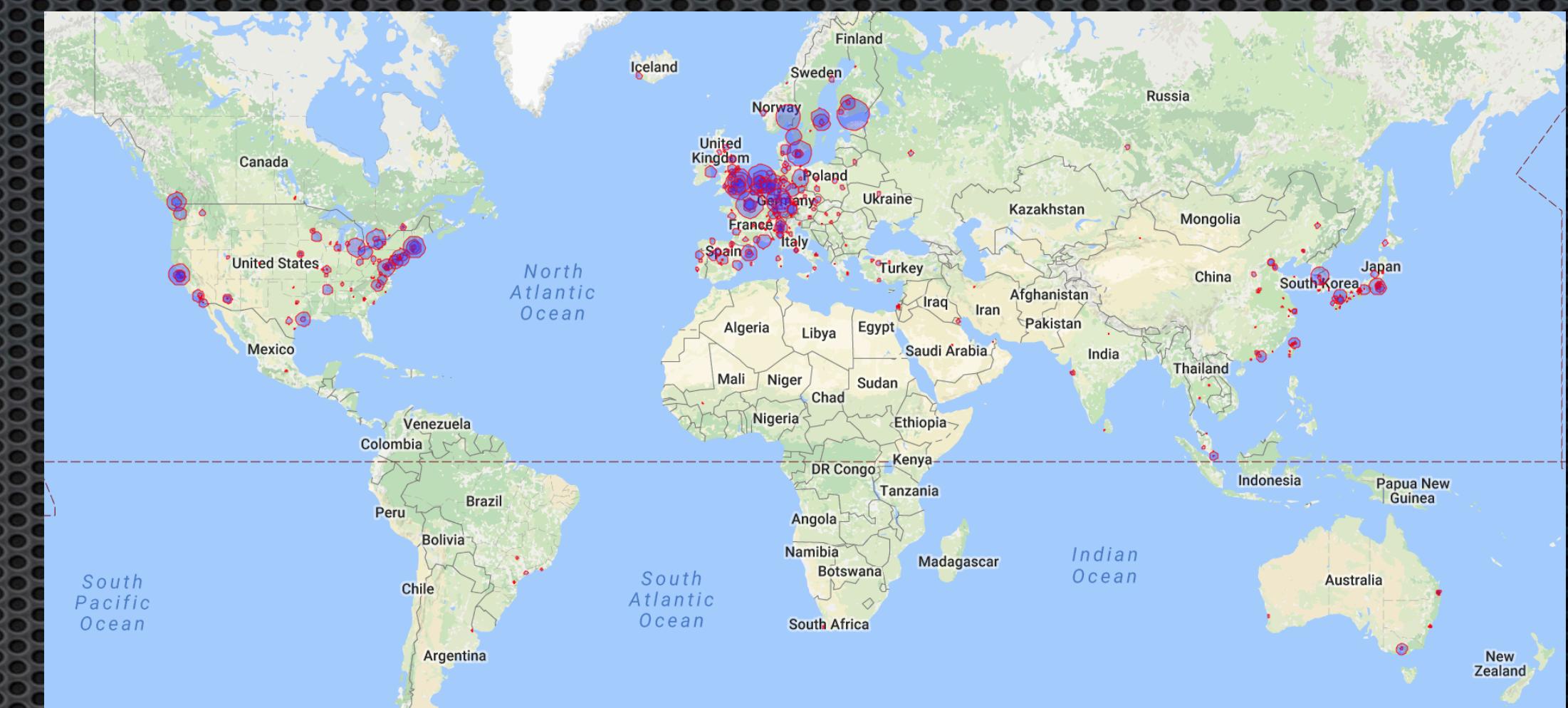
Chromothripsy-like patterns are recurring but heterogeneously distributed features in a survey of 22,347 cancer genome screens

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



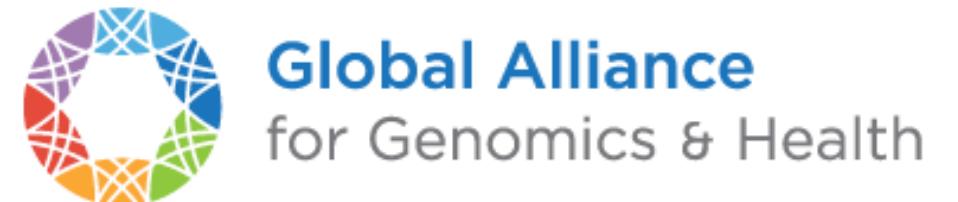
Bias in Ascertainment / Background / Environment in Cancer Genome Studies

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



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

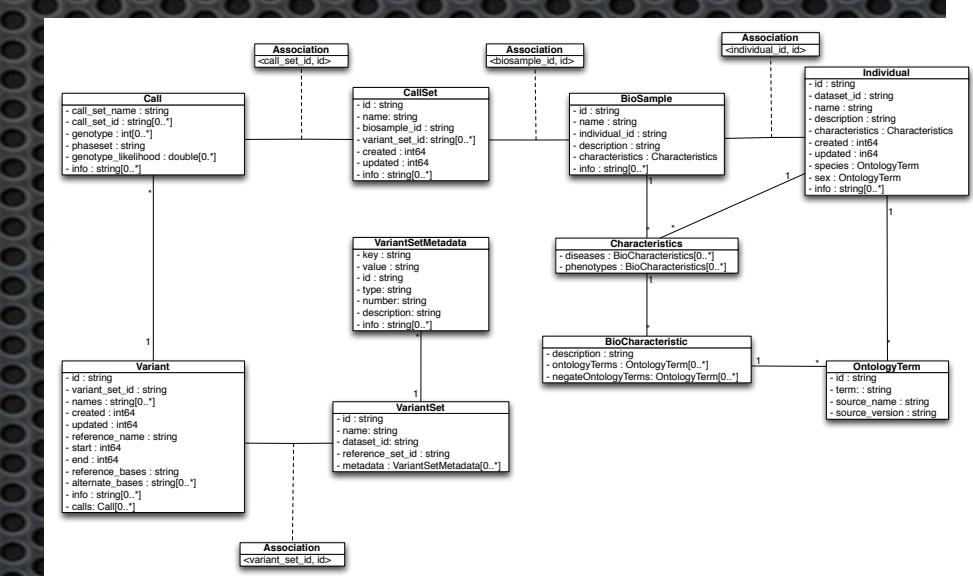




Developing the GA4GH Metadata Schema

▶ arrayMap for GA4GH

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



Driving Beacon Development

▶ Beacon+

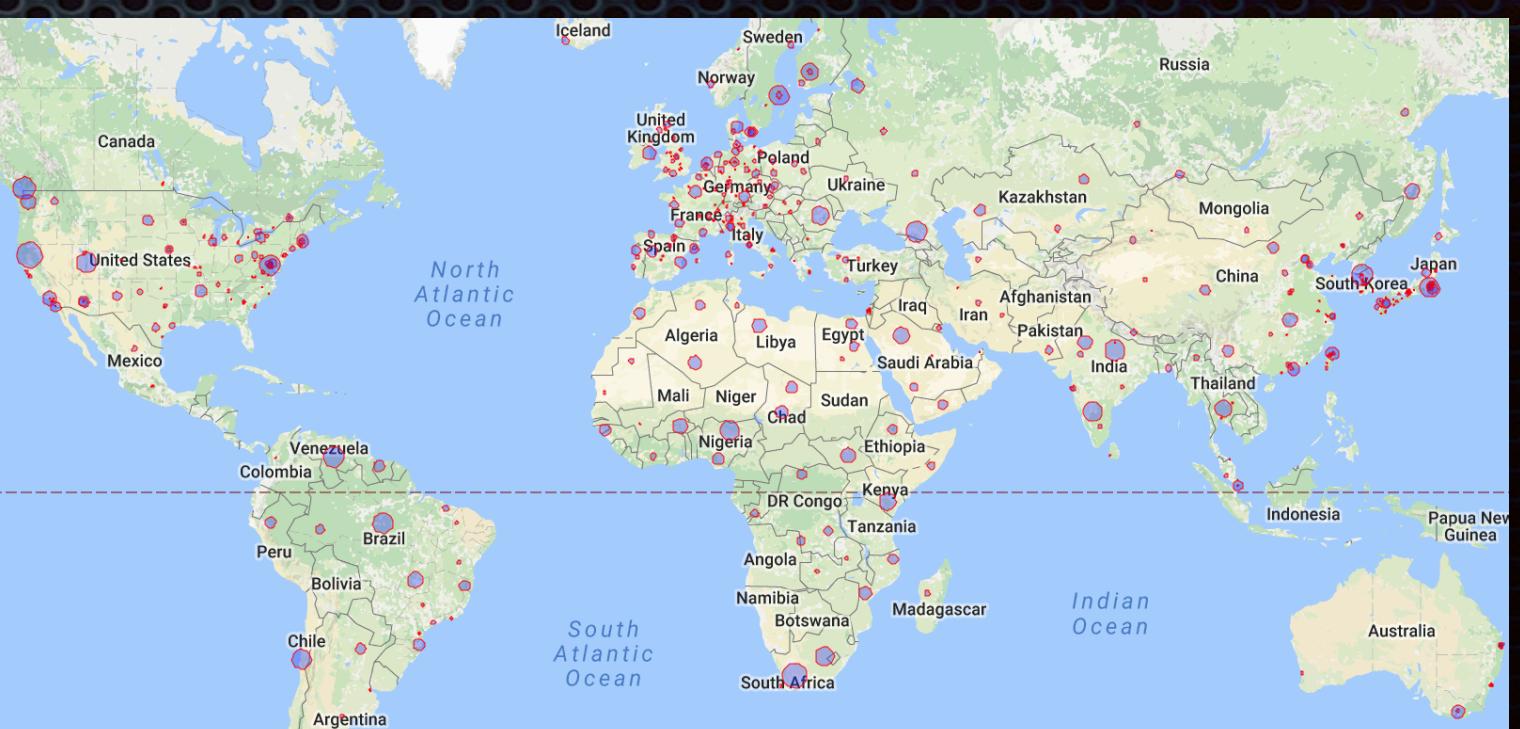
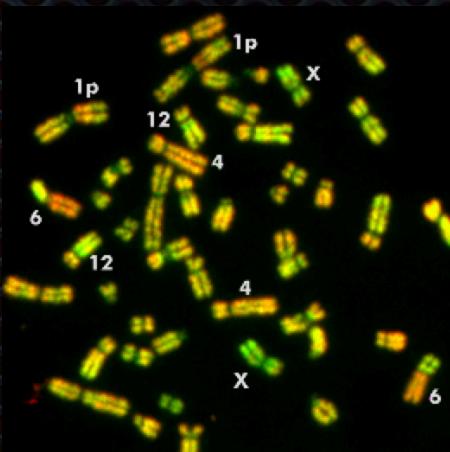
- CNV/CNA as first type of structural variants
 - disease specific queries
 - quantitative reporting

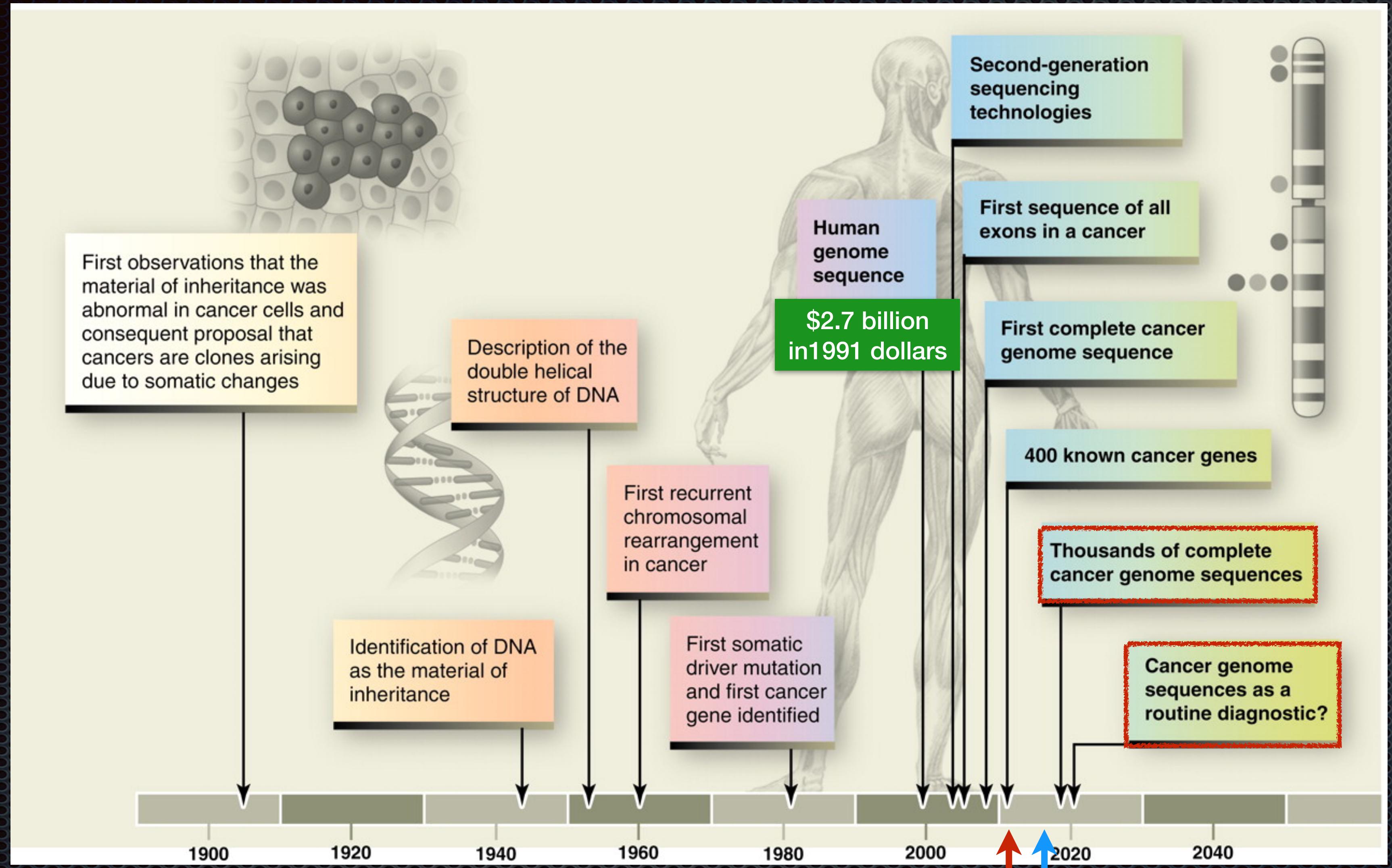


```
{  
    "_id" : ObjectId("58297ca32ca4591e5a0df054"),  
    "id" : "AM_V_1778741",  
    "variant_set_id" : "AM_VS_HG18",  
    "reference_name" : "10"  
    "start" : 579049,  
    "end" : 17236099,  
    "alternate_bases" : "DUP",  
    "reference_bases" : ".",  
    "info" : {  
        "svlen":16657050,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

Cancer Genome Data: Where Do We Need To Go

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





Th. Boveri

Nowell/Hungerford
J. Rowley

2016



Not yet practical
for full human
genomes, but real
product!
(OxfordNanopore)

Michael R. Stratton.
Exploring the Genomes
of Cancer Cells:
Progress and Promise.
Science (2011)

BAUDISGROUP @ UZH

NI AI
MICHAEL BAUDIS
(HAOYANG CAI)
PAULA CARRIO CORDO
BO GAO
(LINDA GROB)
SAUMYA GUPTA)
(ROMAN HILLJE
(NITIN KUMAR)
(ALESSIO MILANESE)

SIB

HEINZ STOCKINGER
SÉVERINE DUVAUD
VASSILIOS IOANNIDIS
DANIEL TEIXEIRA

THOMAS EGGERMANN
ROSA NOGUERA
REINER SIEBERT
CAIUS SOLOVAN



... MY COLLEAGUES AT THE INSTITUTE OF MOLECULAR
LIFE SCIENCES AND OTHER MEMBERS OF THE UZH



University of
Zurich UZH

GA4GH DWG + CWG

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

ELIXIR & CRG

JORDI RAMBLA DE ARGILA
SABELA DE LA TORRE PERNAS
SUSANNA REPO
SERENA SCOLLEN