



University of  
Zurich<sup>UZH</sup>

# BIO392 Bioinformatics of Genome Variations

Genomes: Core of "Personalized Health" & "Precision Medicine"

Michael Baudis **UZH SIB**  
Computational Oncogenomics

# BIO392: Course Schedule

<https://compbiozurich.org/courses/UZH-BIO392/>

	Tue Sep 20	Wed Sep 21	Thu Sep 22	Fri Sep 23	Tue Sep 27	Wed Sep 28	Thu Sep 29	Fri Sep 30	Tue Oct 4	Wed Oct 5	Thu Oct 6	Fri Oct 7	Tue Oct 11	Wed Oct 12
09:00 - 10:00			Izaskun: Terminal / Unix / Files	Izaskun: File formats for human genetic variation / file handling		Michael lecture introduction to some resources, CNVs, Progenetix	Hangjia: Blast	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
10:00 - 11:00		Github exercise: create user specific directories & upload/edit test files using Markdown (Ziying)	Izaskun: Terminal / Unix / Files	Izaskun: File formats for human genetic variation / file handling		Task: Browse/explore genome resources and provide some notes (1-2 pages total) in a doc posted on Github (.md)	Hangjia: blast	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
11:00 - 12:00		Ziying: github desktop and terminal		Izaskun: File formats for human genetic variation / file handling			Hangjia: Blast exercise	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
13:00 - 14:00	* Room information * Administrative - discuss times/days - exam	Ziying: Introduction to different interfaces eg atom, jupyter, pycharm (lecture), include R things	Izaskun: SIB online introduction to Unix	Izaskun: short project (1000 genomes), reading, literature	Recap W1; Q&A	Hangjia: Progenetix as tool for CNV frequencies etc.	Hangjia: Clinvar and Clingen	Max: STRs	Survival lecture, explanations of terms used etc, cancer classifications	Rahel: Survival analysis	Rahel & Feifei: Survival analysis	Discussion of Survival results (groups of groups)	Exam revision, Q&A	
14:00 - 15:00	Tina Siegenthaler: technical introduction (room, computer, accounts)		Izaskun: SIB online introduction to Unix	Izaskun: short project (1000 genomes), reading, literature	Literature (genome analysis techniques ...)			Max: STRs		Rahel	Rahel & Feifei: Survival analysis			
15:00 - 16:30	* explore course site * create Github accounts and forward to bio392@compbiozurich.org * Michael: short introductory lecture about genome variation			Izaskun: short project (1000 genomes), reading, literature	Genome technologies - brief notes about usage scenarios, pro & con			Max: STRs		Rahel	Rahel & Feifei: Survival analysis			

1992



Heidelberg

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

2001



Stanford

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

2003



Gainesville

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

2006



Aachen

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

2007



Zürich

Professor of bioinformatics @ DMLS (2015) | systematic assembly of oncogenomic data | databases and software tools | patterns in cancer genomes | *Progenetix* & *arrayMap* resources | GA4GH | SPHN | ELIXIR

# Our Research

## Theoretical Cytogenetics & Oncogenomics

- CNV resource
    - Data - e.g. [progenetix.org](http://progenetix.org)
    - Tools - CNV remapping, visualization, API access to resources ...
  - patterns and correlations of genomic variations in cancer
  - annotation mapping
  - API, protocols and standards contributions
- Beacon

[info.baudisgroup.org](http://info.baudisgroup.org)

### baudisgroup @ UZH & SIB

[Baudisgroup Home](#)

[Some Projects](#)

[Support or Contact](#)

[Address](#)

[Latest News & Publications](#)

[Group](#)

[Publications](#)

[Presentations](#)

[Projects and Open Positions](#)

[Progenetix ↗](#)

[CompbioZurich ↗](#)

[SchemaBlocks {S}\[B\] ↗](#)

[Beacon Project ↗](#)

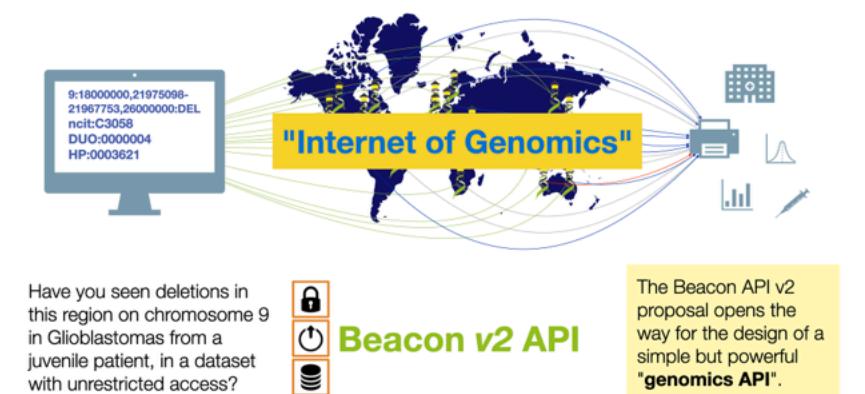
[Michael Baudis @ UZH ↗](#)

## Welcome to the *baudisgroup* Pages

The *baudisgroup* website represents projects and information by the **Computational Oncogenomics Group** of the [University of Zurich \(UZH\)](#) and the [Swiss Institute of Bioinformatics \(SIB\)](#). For visitors more interested in Particle Astrophysics, we strongly recommend the website of another, although related, [Professor Baudis](#).

### The Computational Oncogenomics

Group's research focus lies in the exploration of structural genome variations in cancer. Our work centres around our [Progenetix](#) resource of curated molecular-cytogenetic and sequencing data. Specific projects explore computational methods, genomics of selected tumour entities and genomic variant patterns across malignancies. As members of the [Global Alliance for Genomics and Health](#), the group is developing standards in biocuration and data sharing for genomic variants and phenotypic data, for instance in driving development of the [ELIXIR Beacon](#) project. Other research is related to genome data epistemology, e.g. geographic and diagnostic sampling biases in cancer studies.



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

**Beacon v2 API**

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

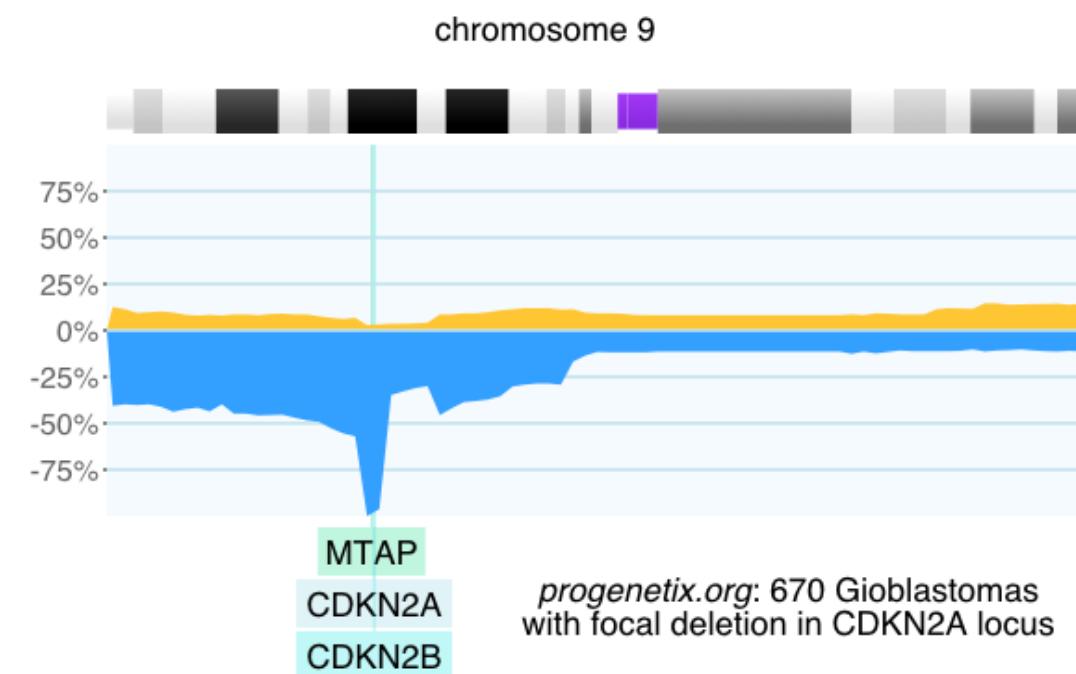
### Some Projects

Ongoing software and service projects can be visited at our Github organizations ([progenetix](#) and [baudisgroup](#)) and when looking at individual contributions to e.g. GA4GH and ELIXIR.

- [bycon](#) at Github in [Progenetix](#) - Python based implementation of a GA4GH Beacon
- [pgxRpi](#) at Github in [Progenetix](#) - An API wrapper package in R for loading & displaying data from Progenetix
- [segment-liftover](#) at Github [baudisgroup](#)
  - publication
- [SNP2pop](#) at Github [baudisgroup](#)
  - publication at [ScientificReports](#)
- [ICDOntologies](#) at Github in [Progenetix](#) - mapping disease concepts

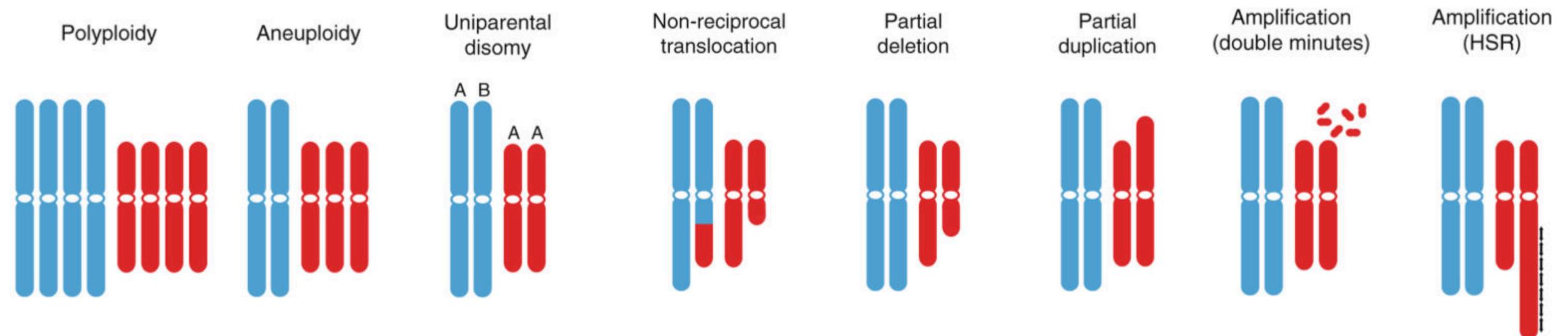
# Theoretical Cytogenetics and Oncogenomics

## Research | Methods | Standards

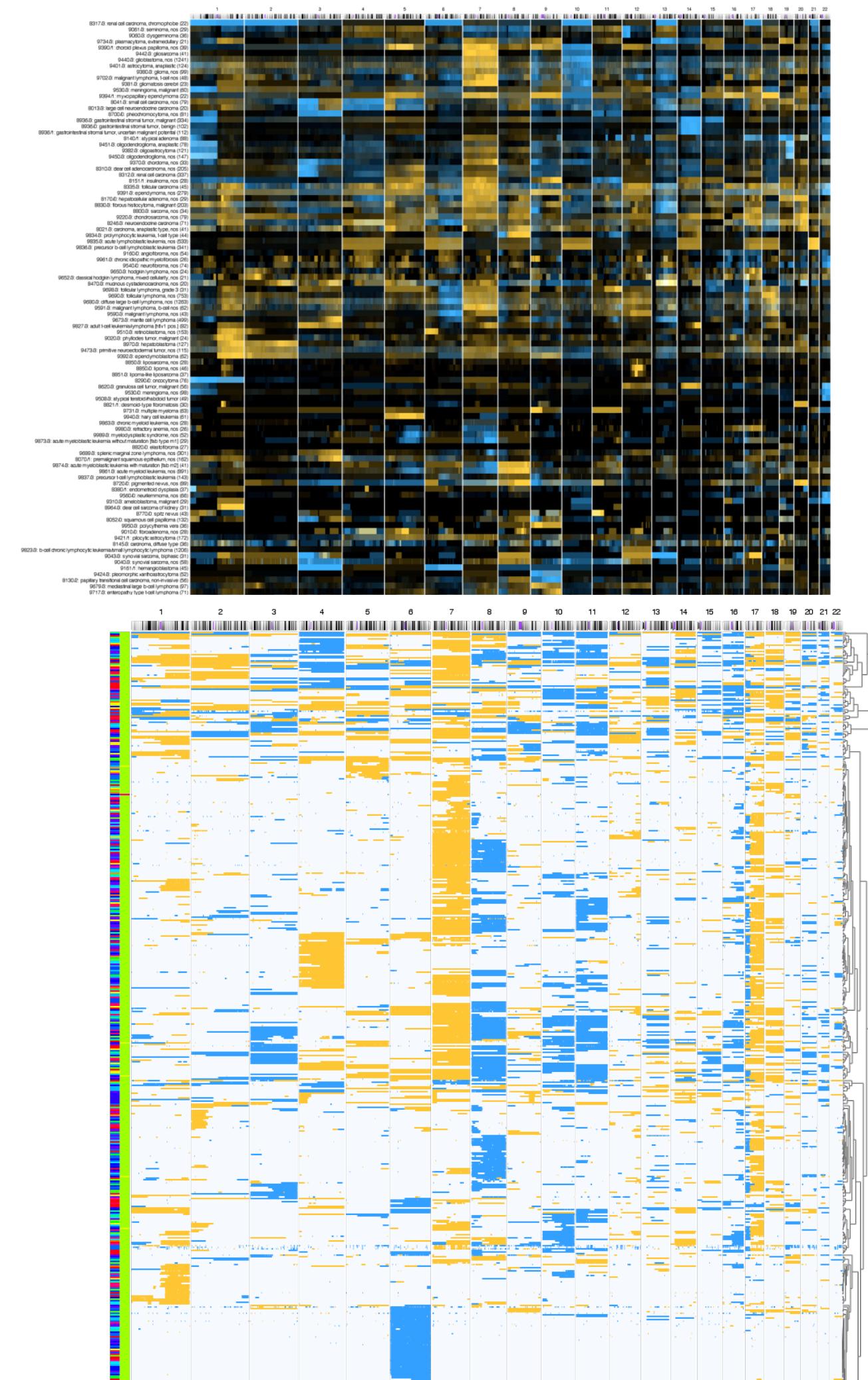


# Genomic Imbalances in Cancer - Copy Number Variations (CNV)

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



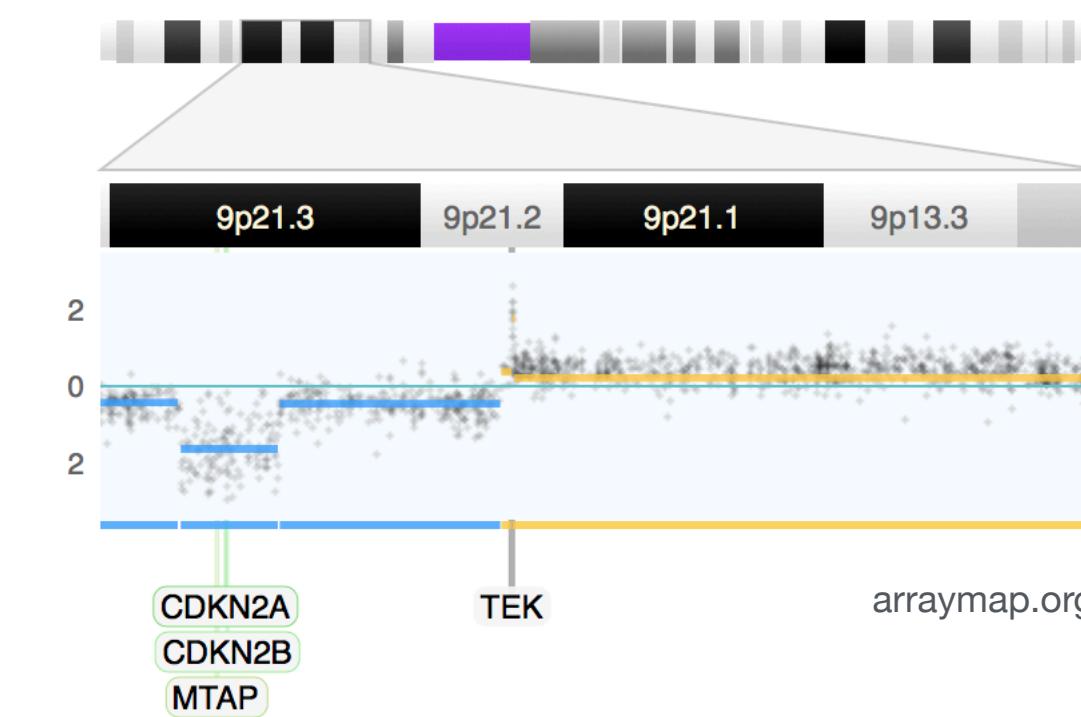
Grade et al., 2015 Recent Results Cancer Res



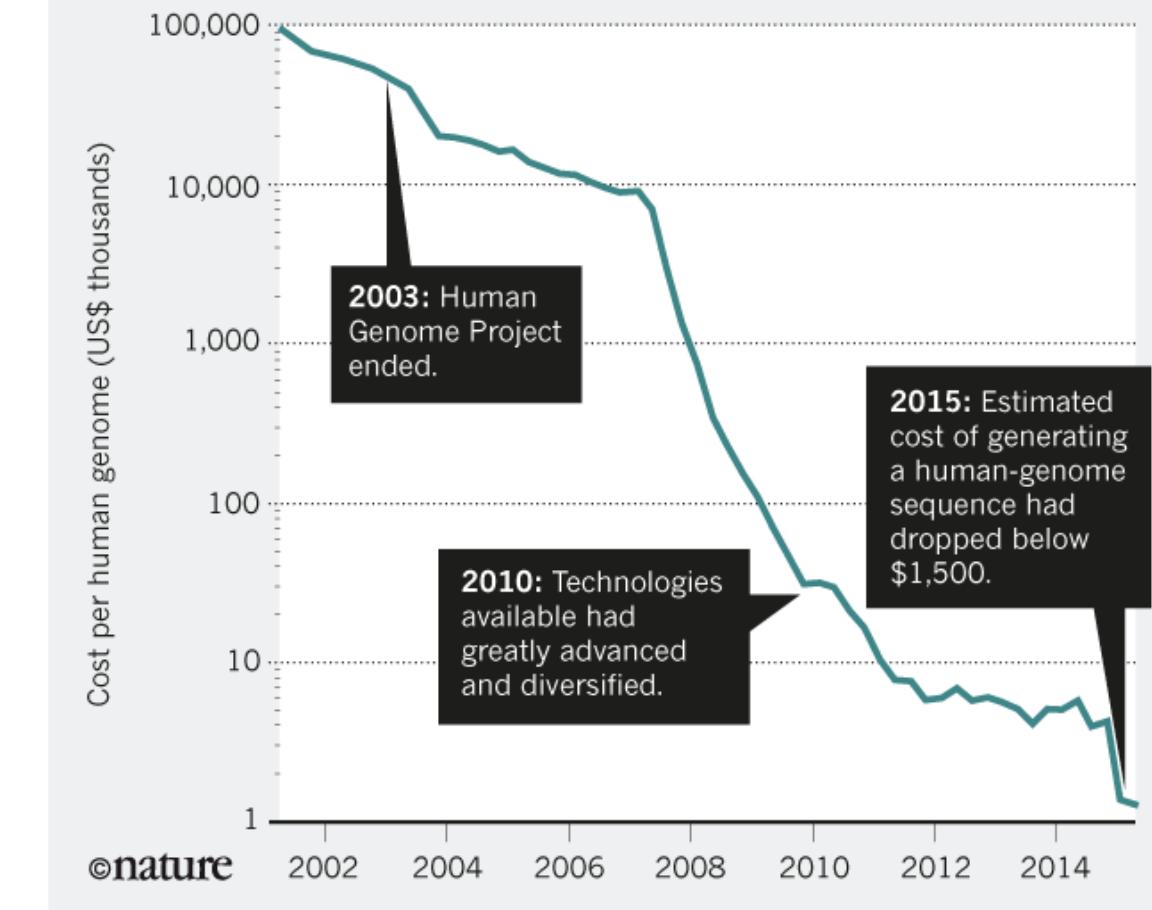


## Genome screening at the core of “Personalised Health”

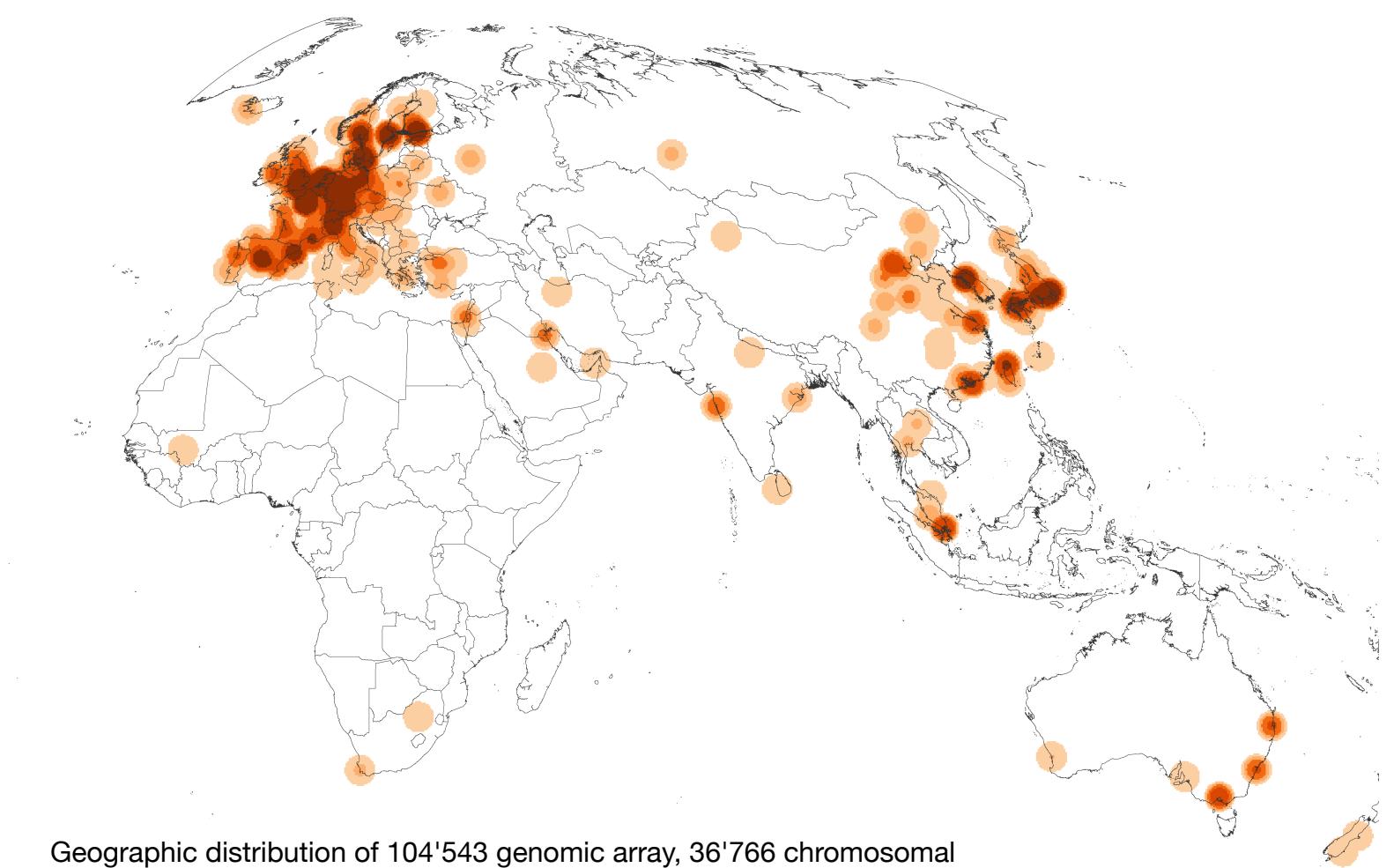
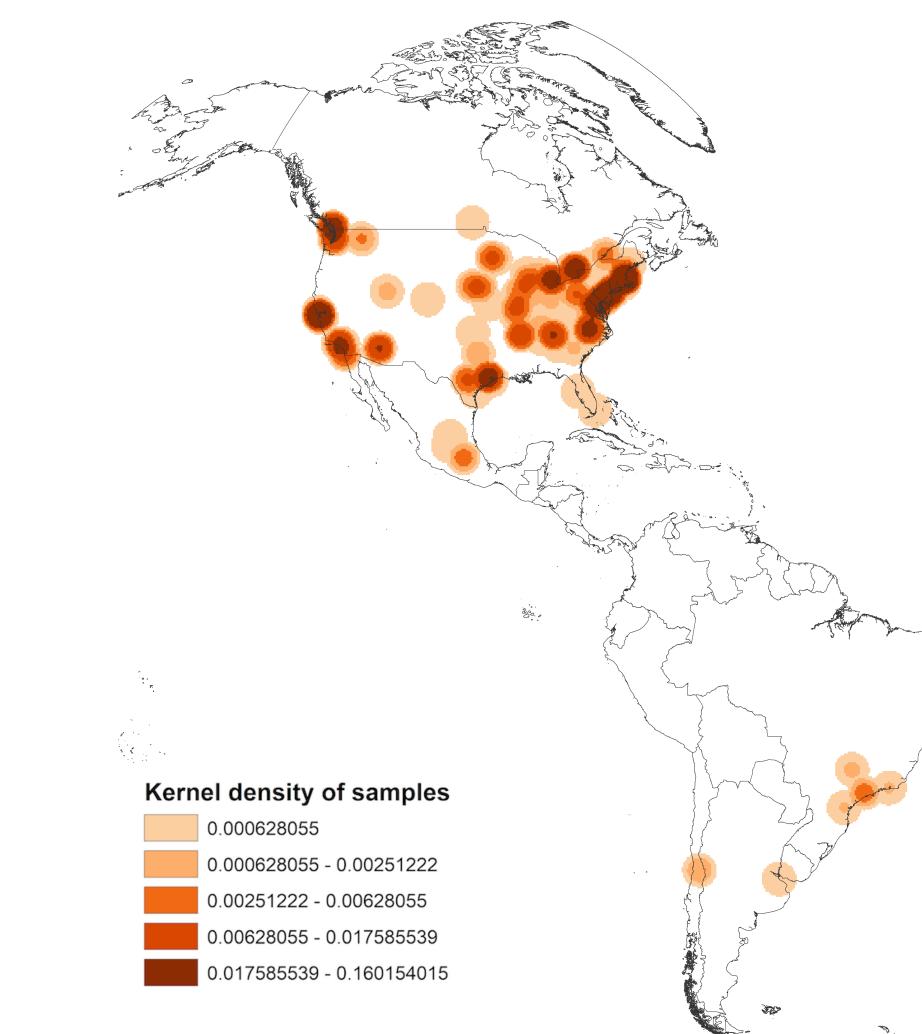
- ▶ **Genome analyses** (including transcriptome, metagenomics) are core technologies for Personalised Health™ applications
- ▶ The unexpectedly large amount of **sequence variants** in human genomes - germline and somatic/cancer - requires huge analysis efforts and creation of **reference repositories**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



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



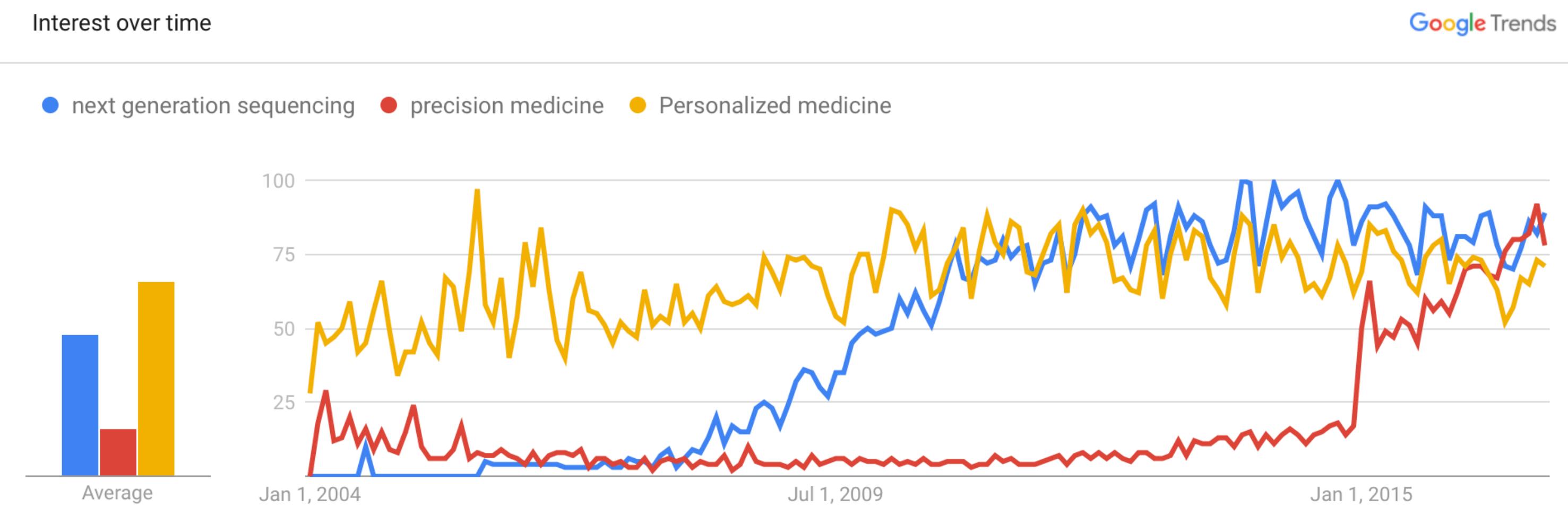
The future of DNA sequencing. Eric D. Green, Edward M. Rubin & Maynard V. Olson. Nature; 11 October 2017 (News & Views)



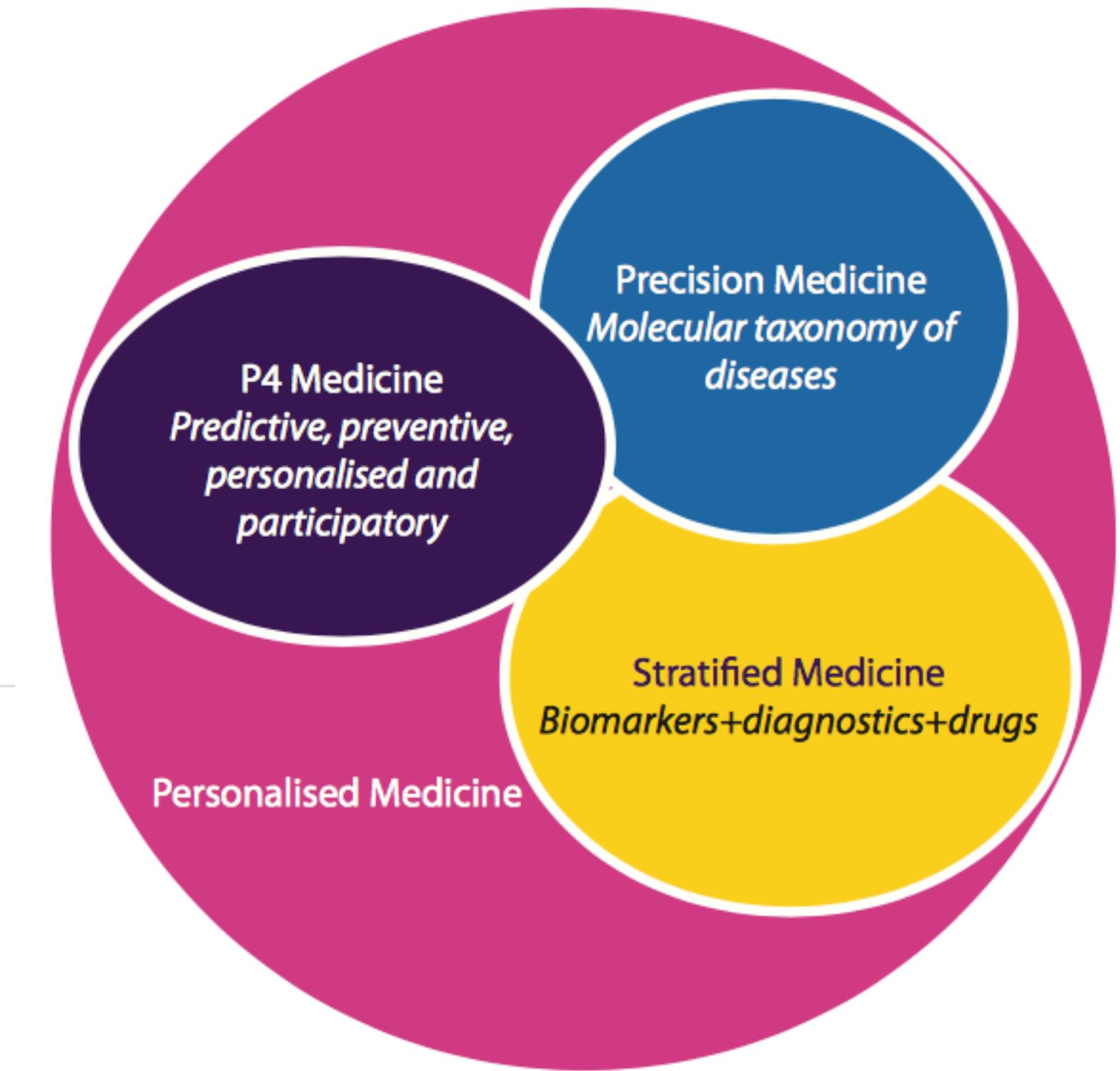
Geographic distribution of 104'543 genomic array, 36'766 chromosomal CGH and 15'409 whole genome/exome based cancer genome datasets

# 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?



Worldwide. 2004 - present.

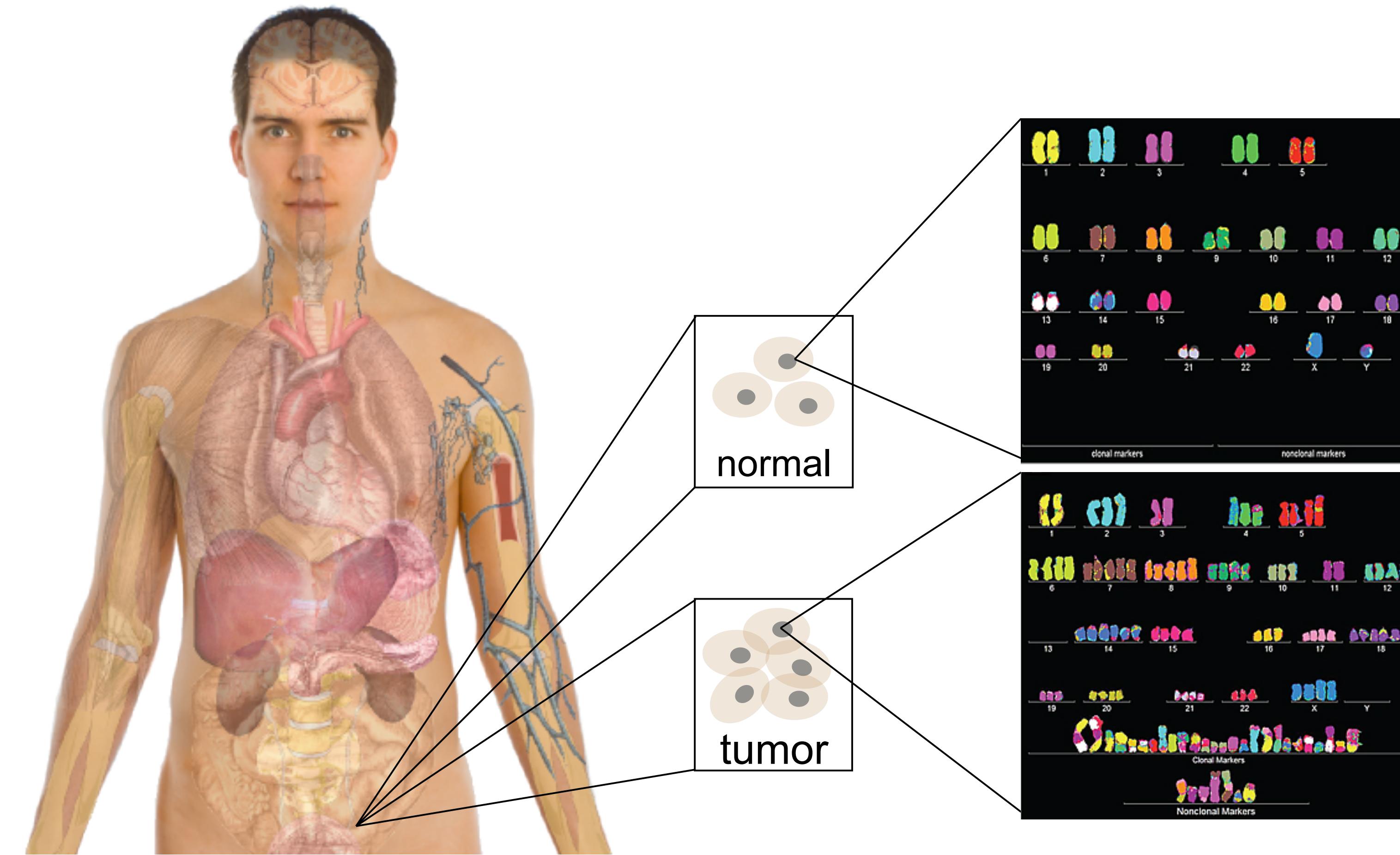


Source: PHG Foundation

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

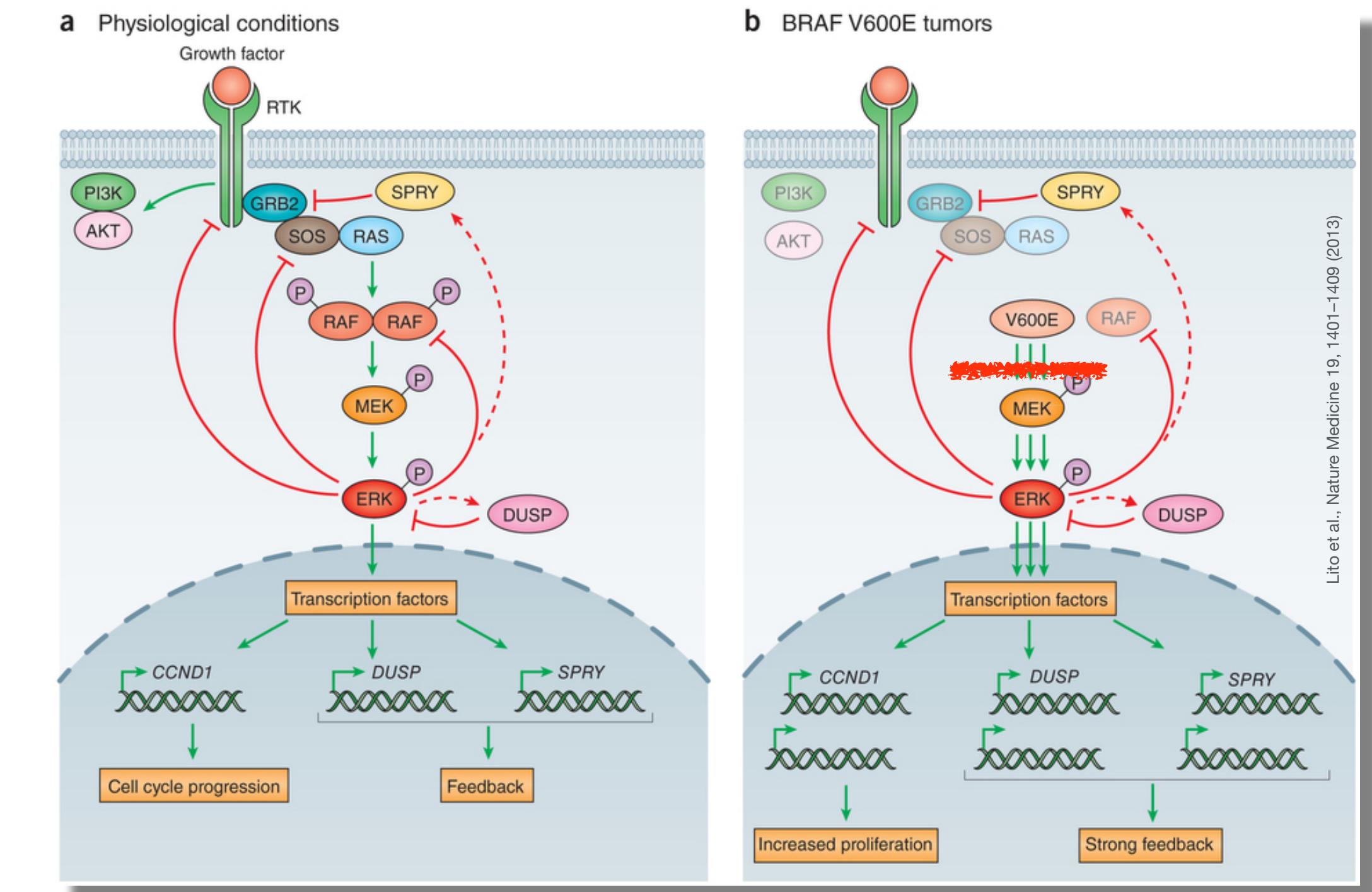
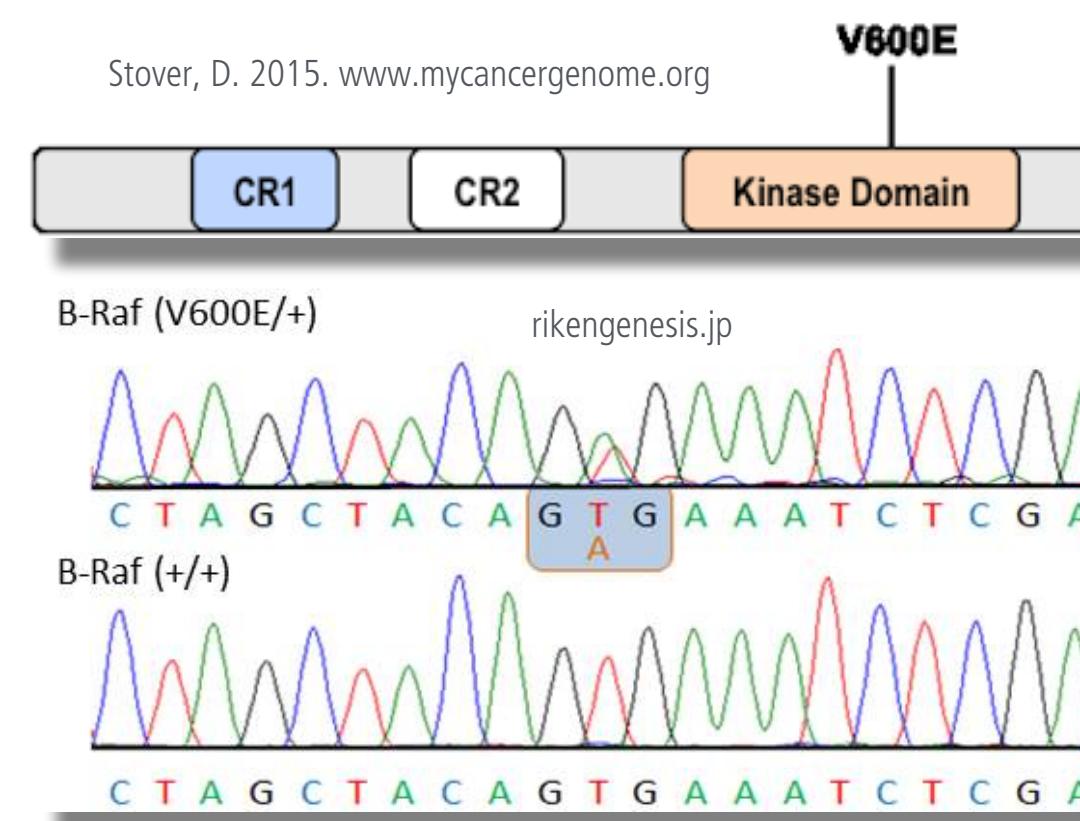
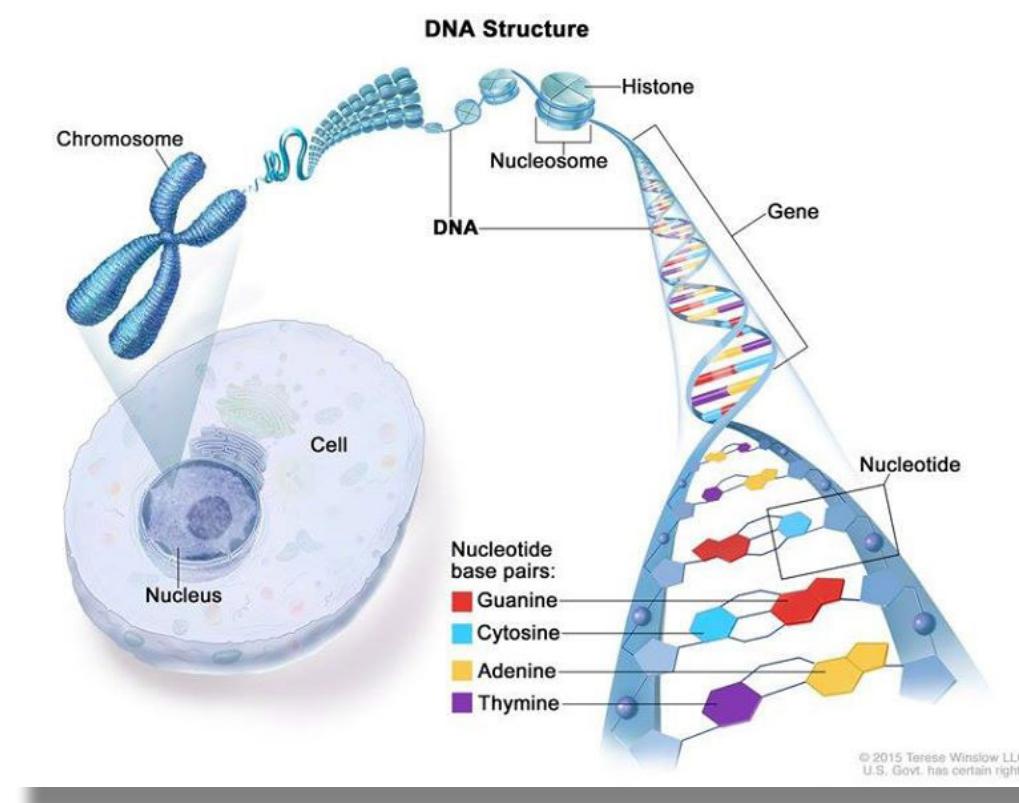
Personal genomes soon will become a commonplace part of medical research & eventually treatment (esp. for cancer). They will provide a primary connection for biological science to the general public.



# 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.

# 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<sup>1,2,3,4</sup>, Konrad J. Karczewski<sup>1,2\*</sup>, Eric V. Minikel<sup>1,2,5\*</sup>, Kaitlin E. Samocha<sup>1,2,5,6\*</sup>, Eric Banks<sup>2</sup>, Timothy Fennell<sup>2</sup>, Anne H. O'Donnell-Luria<sup>1,2,7</sup>, James S. Ware<sup>2,8,9,10,11</sup>, Andrew J. Hill<sup>1,2,12</sup>, Beryl B. Cummings<sup>1,2,5</sup>, Taru Tukiainen<sup>1,2</sup>, Daniel P. Birnbaum<sup>2</sup>, Jack A. Kosmicki<sup>1,2,6,13</sup>, Laramie E. Duncan<sup>1,2,6</sup>, Karol Estrada<sup>1,2</sup>

Rapid whole genome sequencing and precision neonatology

CrossMark

Joshua E. Petrikirin, MD<sup>a,\*</sup>, Laurel K. Willig, MD, FAAP<sup>b</sup>, Laurie D. Smith, MD, PhD<sup>c</sup>, and Stephen F. Kingsmore, MB, BAO, ChB, Dsc, FRCPath<sup>d,e</sup>

## Genomic Classification of Cutaneous Melanoma

The Cancer Genome Atlas Network<sup>1,\*,\*\*</sup>

<sup>1</sup>Cancer 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>

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<sup>1,2</sup> and James R. Lupski<sup>1,3,4,5</sup>

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

## Consequences of genomic diversity in *Mycobacterium tuberculosis*

Mireia Coscolla<sup>a,b</sup>, Sébastien Gagneux<sup>a,b,\*</sup>

<sup>a</sup> Department of Medical Parasitology and Infection Biology, Swiss Tropical and Public Health Institute, Socinstrasse 57, 4002 Basel, Switzerland

<sup>b</sup> University of Basel, Petersplatz 1, Basel 4003, Switzerland

## Common gene variants, mortality and extreme longevity in humans

B.T. Heijmans<sup>a,b</sup>, R.G.J. Westendorp<sup>b</sup>, P.E. Slagboom<sup>a,\*</sup>

RESEARCH ARTICLE

Open Access

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

Alfred Beletz<sup>1,5\*</sup>, Philip Zimmermann<sup>2</sup>, Michael Baudis<sup>3</sup>, Nicole Brun<sup>4</sup>, Peter Bühlmann<sup>4</sup>, Oliver Laule<sup>2</sup>, Hu-Duc Luu<sup>1</sup>, Wilhelm Gruissem<sup>2</sup>, Peter Schraml<sup>1,\*</sup> and Holger Moch<sup>1</sup>

NEURODEVELOPMENT

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

Mustafa Sahin\* and Srivatsa 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

Rameen Beroukhim<sup>1,3,4,5,\*</sup>, Craig H. Mermel<sup>1,3,\*</sup>, Dale Porter<sup>8</sup>, Guo Wei<sup>1</sup>, Soumya Raychaudhuri<sup>1,4</sup>, Jerry Donovan<sup>8</sup>, Jordi Barretina<sup>1,3</sup>, Jesse S. Boehm<sup>1</sup>, Jennifer Dobson<sup>1,3</sup>, Mitsuyoshi Urashima<sup>9</sup>, Kevin T. McHenry<sup>8</sup>, Reid M. Pinchback<sup>1</sup>, Azra H. Ligon<sup>4</sup>, Yoon-Jae Cho<sup>6</sup>, Leila Haery<sup>1,3</sup>, Heidi Greulich<sup>1,3,4,5</sup>, Michael Reich<sup>1</sup>, Wendy Winckler<sup>1</sup>, Michael S. Lawrence<sup>1</sup>, Barbara A. Weir<sup>1,3</sup>, Kumiko E. Tanaka<sup>1,3</sup>, Derek Y. Chiang<sup>1,3,13</sup>, Adam J. Bass<sup>1,3,4</sup>, Alice Loo<sup>8</sup>, Carter Hoffman<sup>1,3</sup>, John Prensner<sup>1,3</sup>, Ted Liefeld<sup>1</sup>, Qing Gao<sup>1</sup>, Derek Yecies<sup>3</sup>, Sabina Signoretti<sup>3,4</sup>, Elizabeth Maher<sup>10</sup>, Frederic J. Kaye<sup>11</sup>, Hidefumi Sasaki<sup>12</sup>, Joel E. Tepper<sup>13</sup>, Jonathan A. Fletcher<sup>4</sup>, Josep Tabernero<sup>14</sup>, José Baselga<sup>14</sup>, Ming-Sound Tsao<sup>15</sup>, Francesca Demichelis<sup>16</sup>, Mark A. Rubin<sup>16</sup>, Pasi A. Janne<sup>3,4</sup>, Mark J. Daly<sup>1,17</sup>, Carmelo Nucera<sup>7</sup>, Ross L. Levine<sup>18</sup>, Benjamin L. Ebert<sup>1,4,5</sup>, Stacey Gabriel<sup>1</sup>, Anil K. Rustgi<sup>19</sup>, Cristina R. Antonescu<sup>18</sup>, Marc Ladanyi<sup>18</sup>, Anthony Letai<sup>3</sup>, Levi A. Garraway<sup>1,3</sup>, Massimo Loda<sup>3,4</sup>, David G. Beer<sup>20</sup>, Lawrence D. True<sup>21</sup>, Aikou Okamoto<sup>22</sup>, Scott L. Pomeroy<sup>6</sup>, Samuel Singer<sup>18</sup>, Todd R. Golub<sup>1,3,23</sup>, Eric S. Lander<sup>1,2,5</sup>, Gad Getz<sup>1</sup>, William R. Sellers<sup>8</sup> & Matthew Meyerson<sup>1,3,5</sup>

1

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<sup>1,2,3</sup>, Branko Aleksić, MD, PhD<sup>3</sup> and Norio Ozaki, MD, PhD<sup>3,\*</sup>

RESEARCH ARTICLE

Open Access

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

Haoyang Cai<sup>1,2</sup>, Nitin Kumar<sup>1,2</sup>, Homayoun C Bagheri<sup>3</sup>, Christian von Mering<sup>1,2</sup>, Mark D Robinson<sup>1,2\*</sup>

# Genome analyses at the core of Personalized Health™

## There'll be Sequencing Everywhere...

- 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

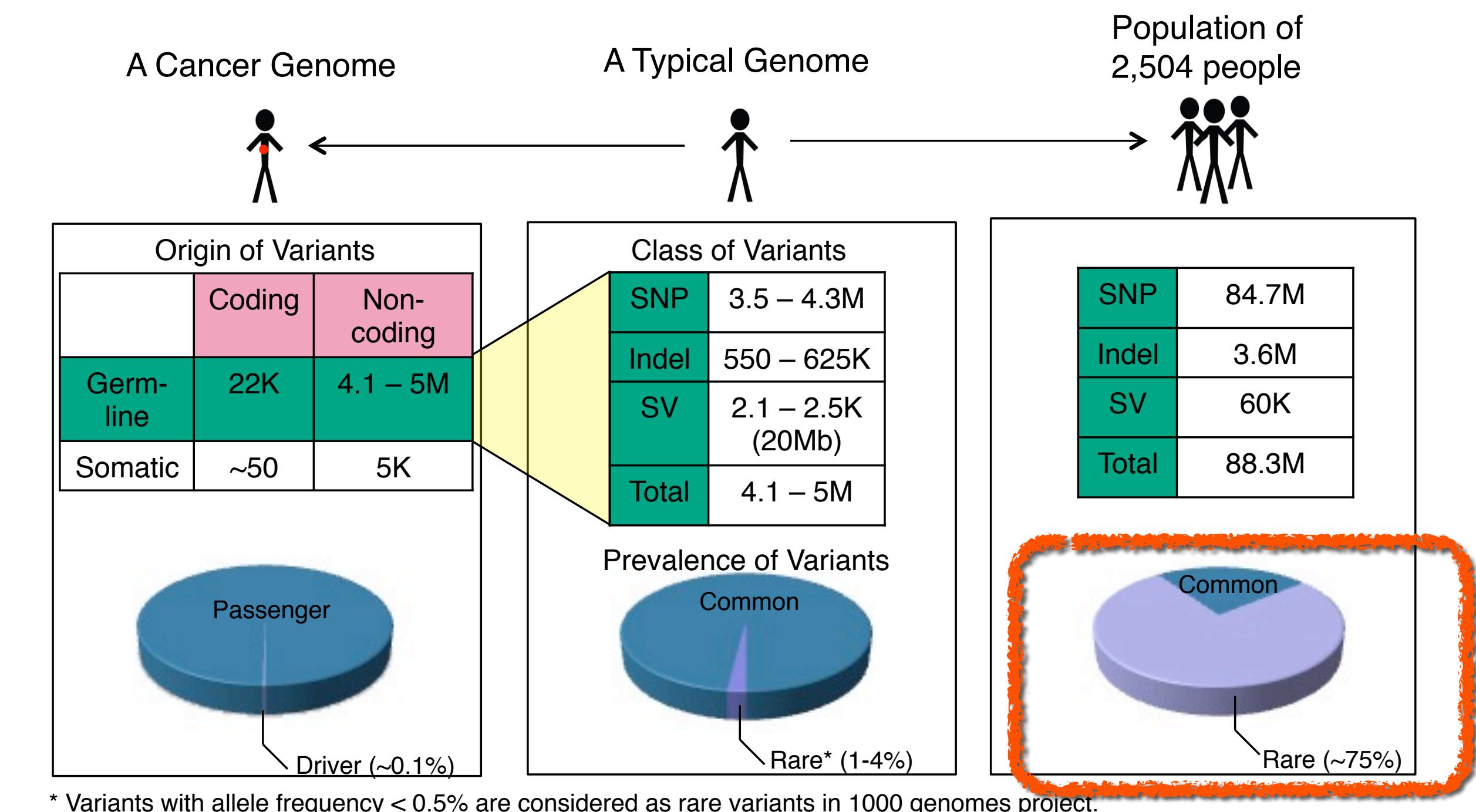
The trouble with human genome variation



# 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

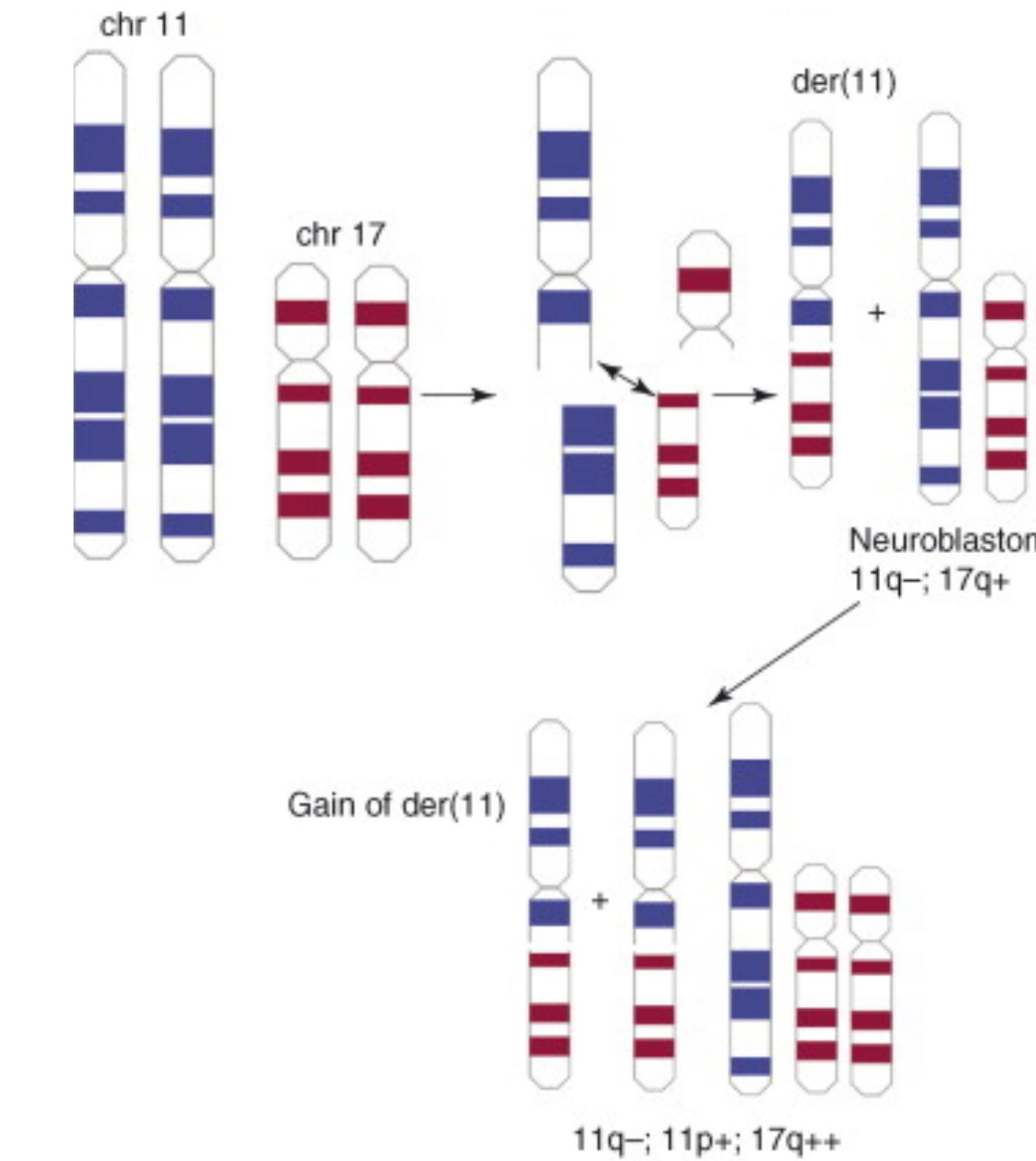
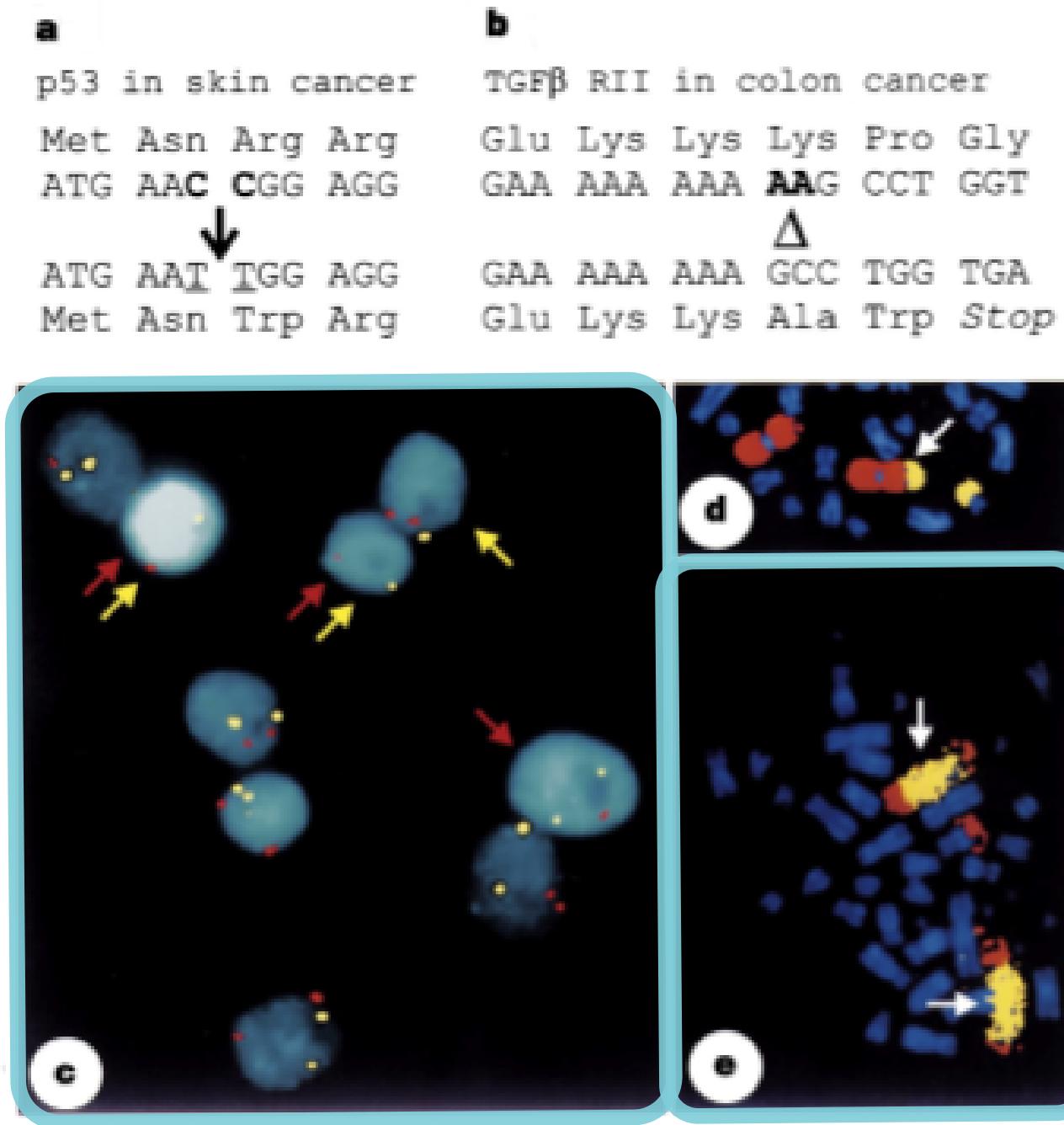


The 1000 Genomes Project Consortium, Nature. 2015. 526:68-74  
Khurana E. et al. Nat. Rev. Genet. 2016. 17:93-108

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

# Mutations & genomic rearrangements in cancer

Lengauer et al. Genetic instabilities in human cancers. Nature (1998) vol. 396 (6712) pp. 643-9

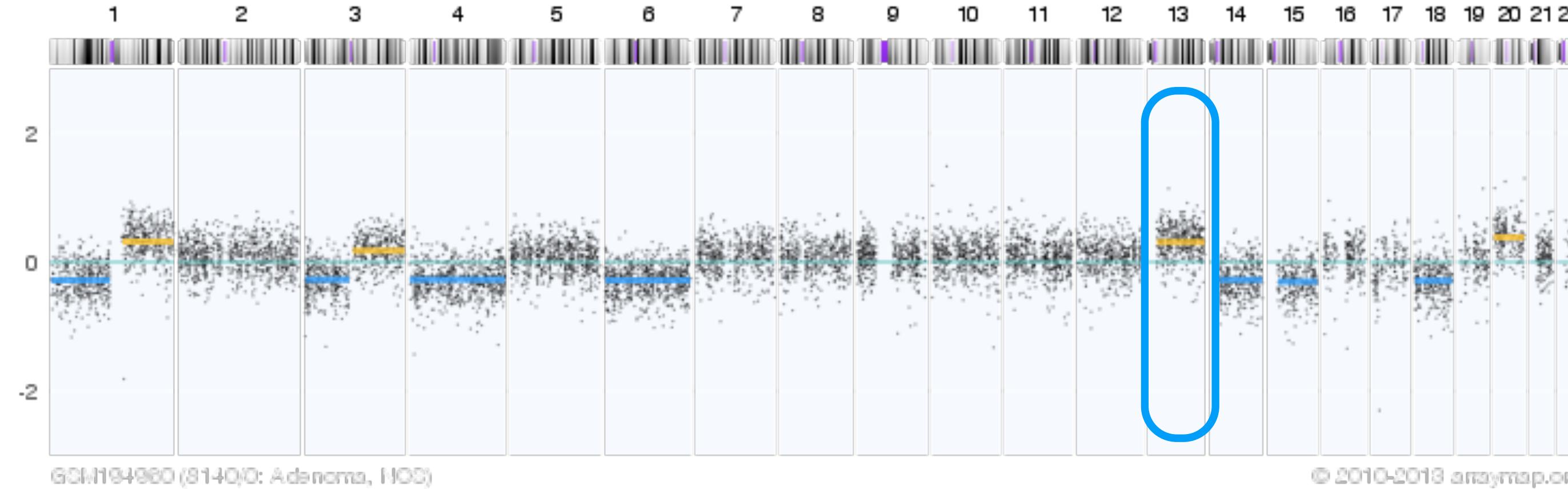


- a. small mutation (di-pyrimidine exchange at p53 in Xeroderma pigmentosum patient)
- b. two-base deletion in *TGFB* in a colorectal cancer patient with mismatch repair deficiency
- c. chromosomal losses (FISH; red=3, yellow=12) in CRC
- d. t(1;17) in neuroblastoma, whole-chromosomal painting
- e. *MYCN*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

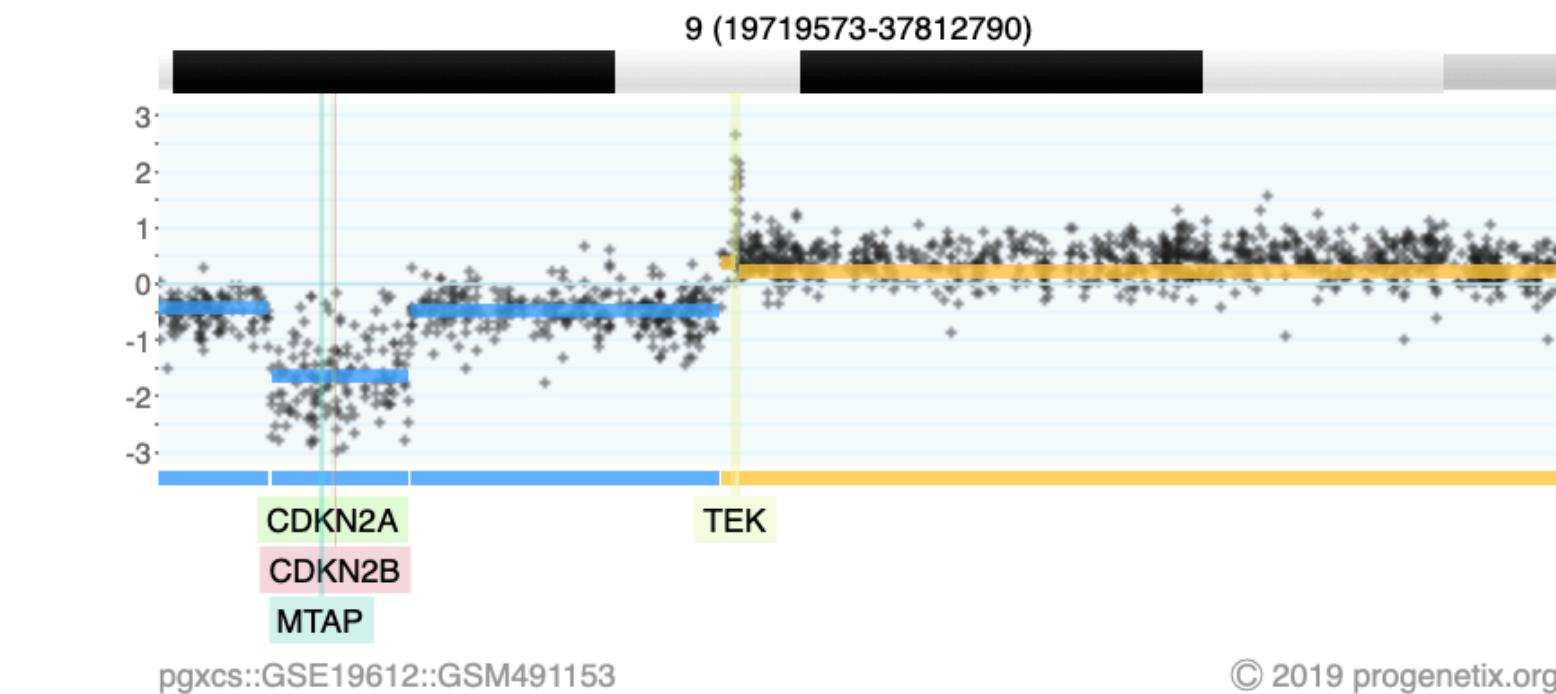
Generation of copy number imbalances in cancer through imbalanced cytogenetic rearrangements - partial deletion of 11q, gain of 11pterq21 and 2 addl. copies of 17q

RL Stallings: Are chromosomal imbalances important in cancer? Volume 23, Issue 6, p278-283, 2007

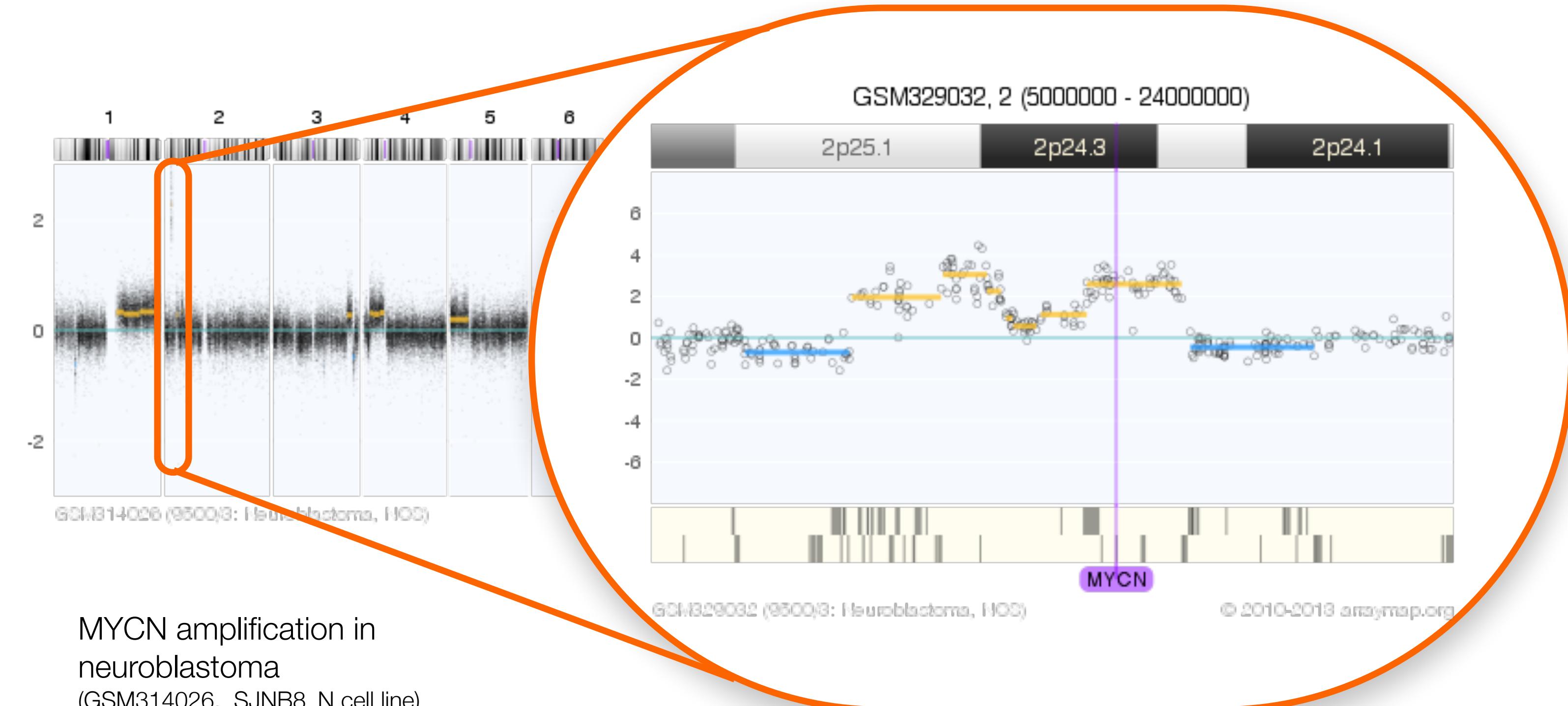
# Somatic Copy Number Variations



Gain of chromosome arm 13q in colorectal carcinoma



2-event, homozygous deletion in a Glioblastoma

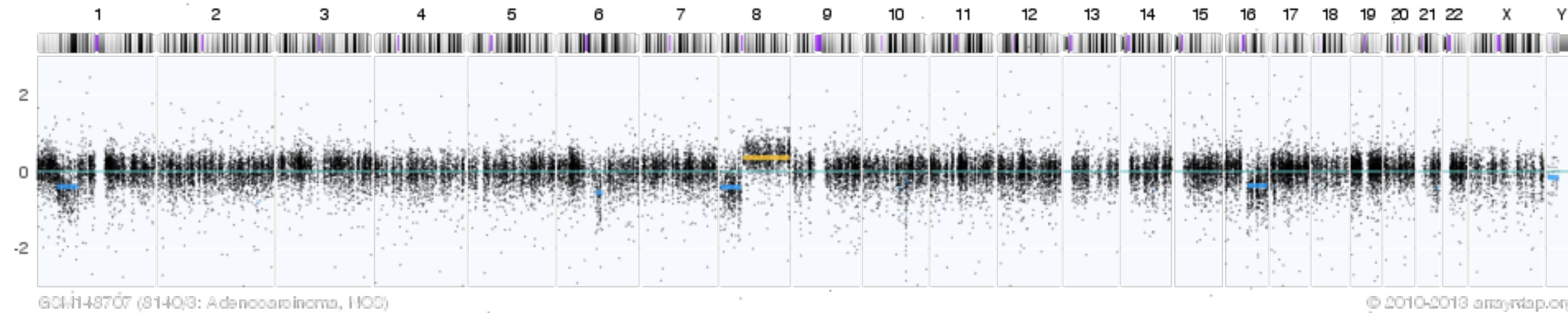
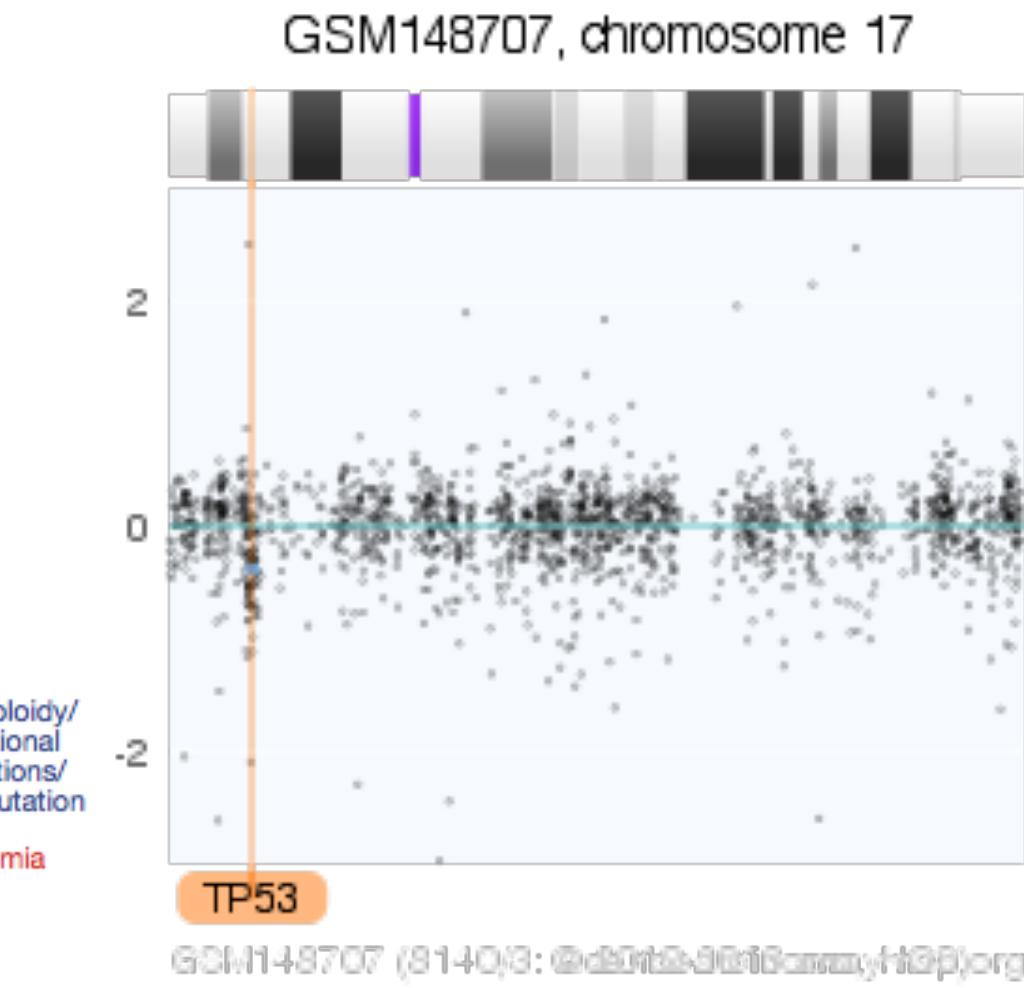
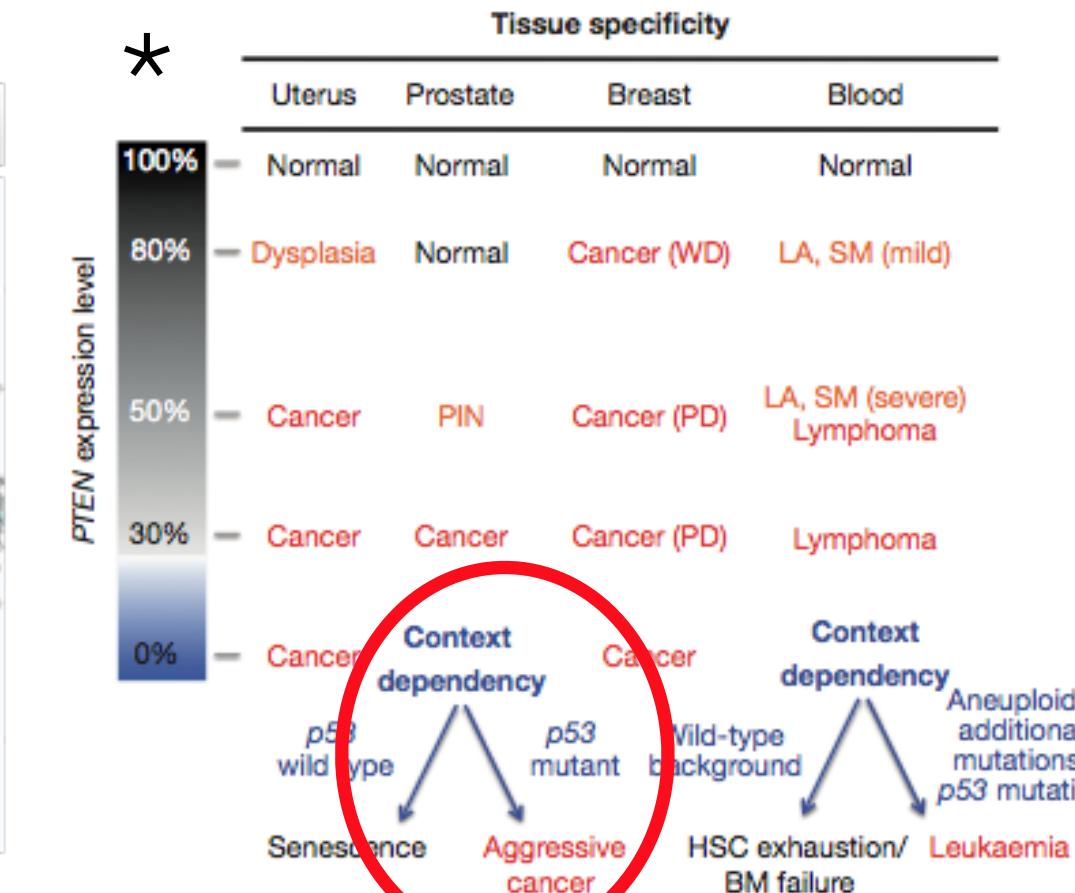
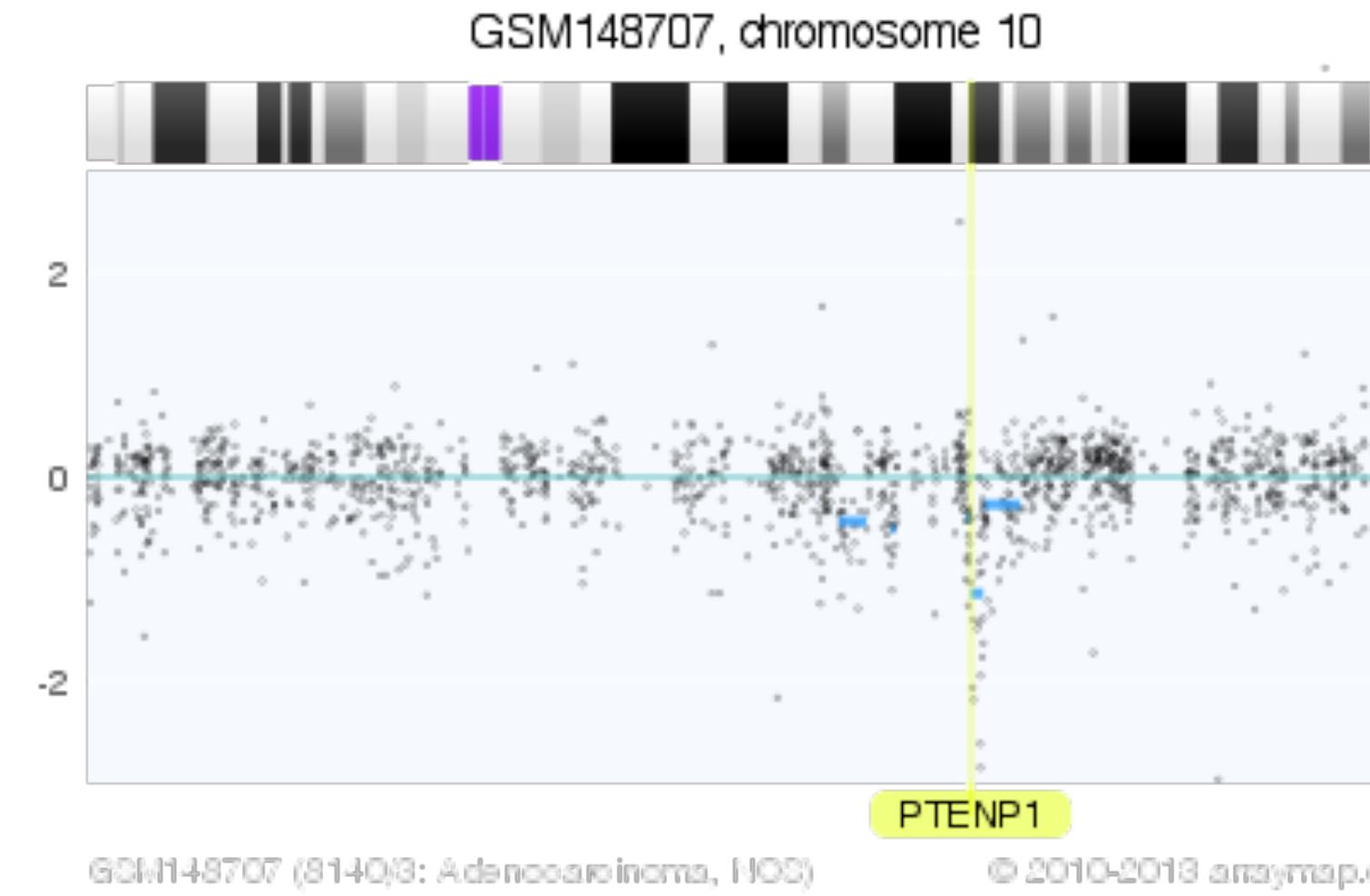


**low level/high level** copy number alterations (CNAs)

arrayMap

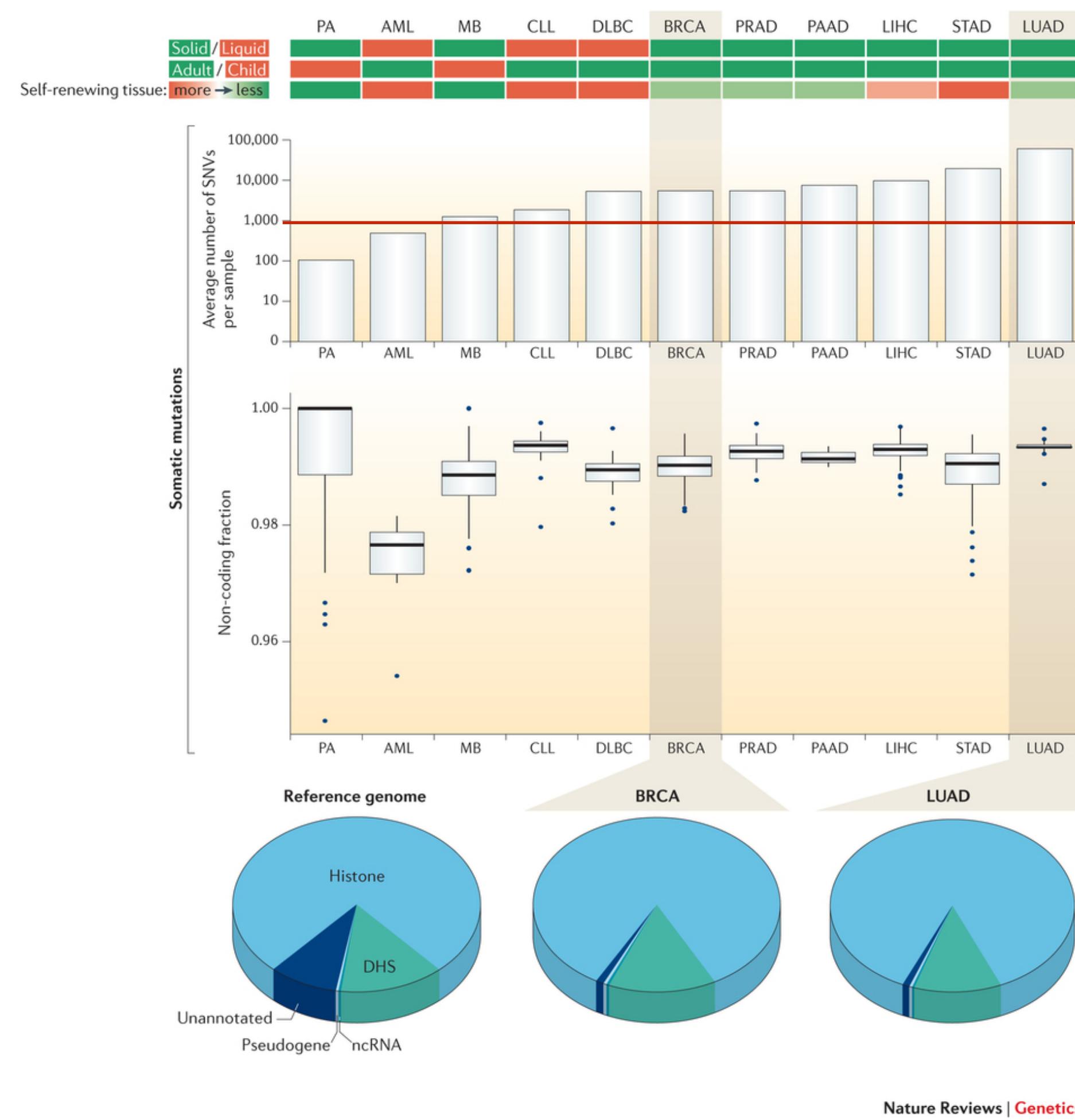


# Gene dosage phenomena beyond simple on/off effects



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

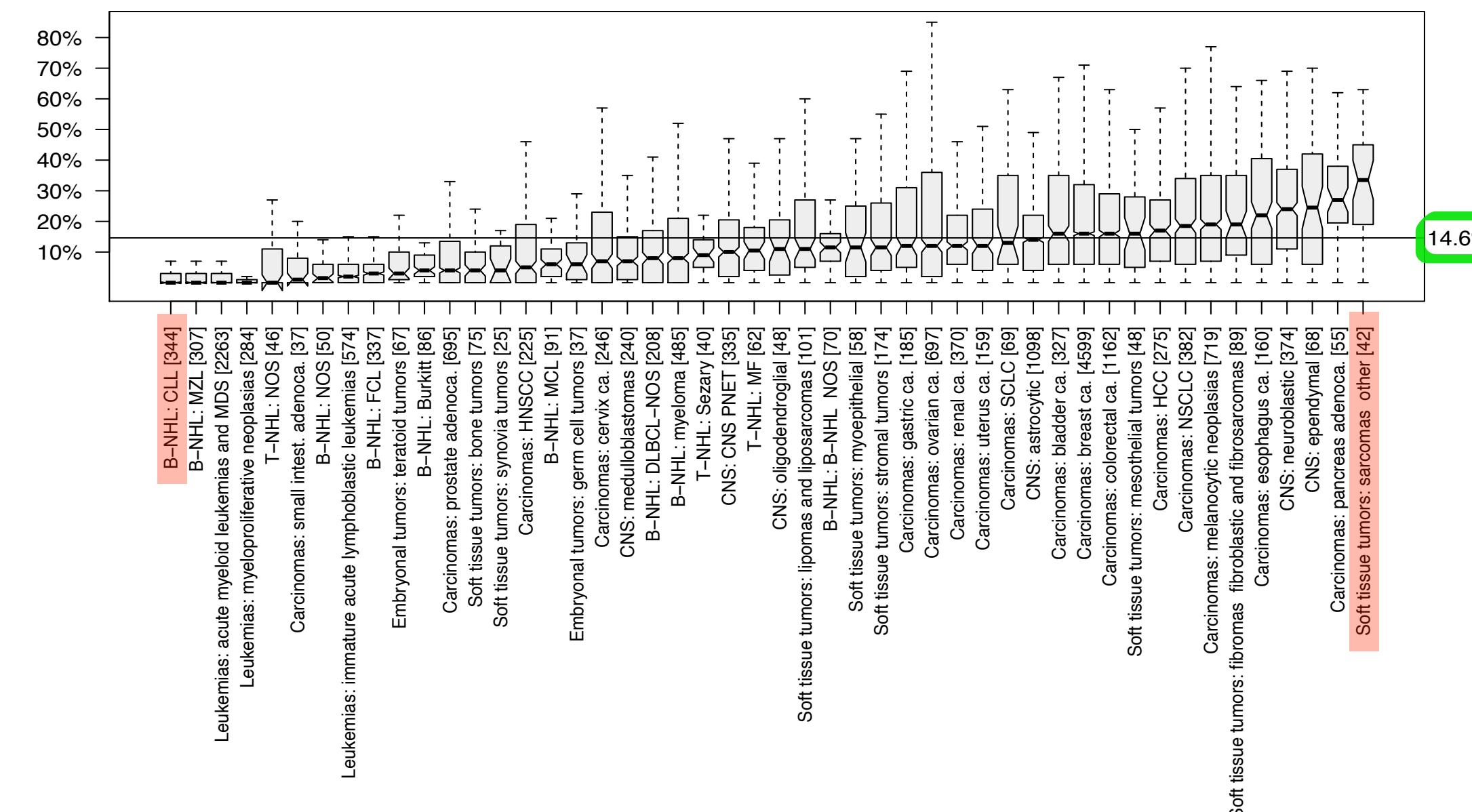
\* A. H. Berger, A. G. Knudson, and P. P. Pandolfi, "A continuum model for tumour suppression," *Nature*, vol. 476, no. 7359, pp. 163–169, Aug. 2011.



CANCERS SHOW THOUSANDS OF SINGLE NUCLEOTIDE VARIANTS PER SAMPLE, MOSTLY IN NON-CODING REGIONS

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))

# Quantifying Somatic Mutations In Cancer

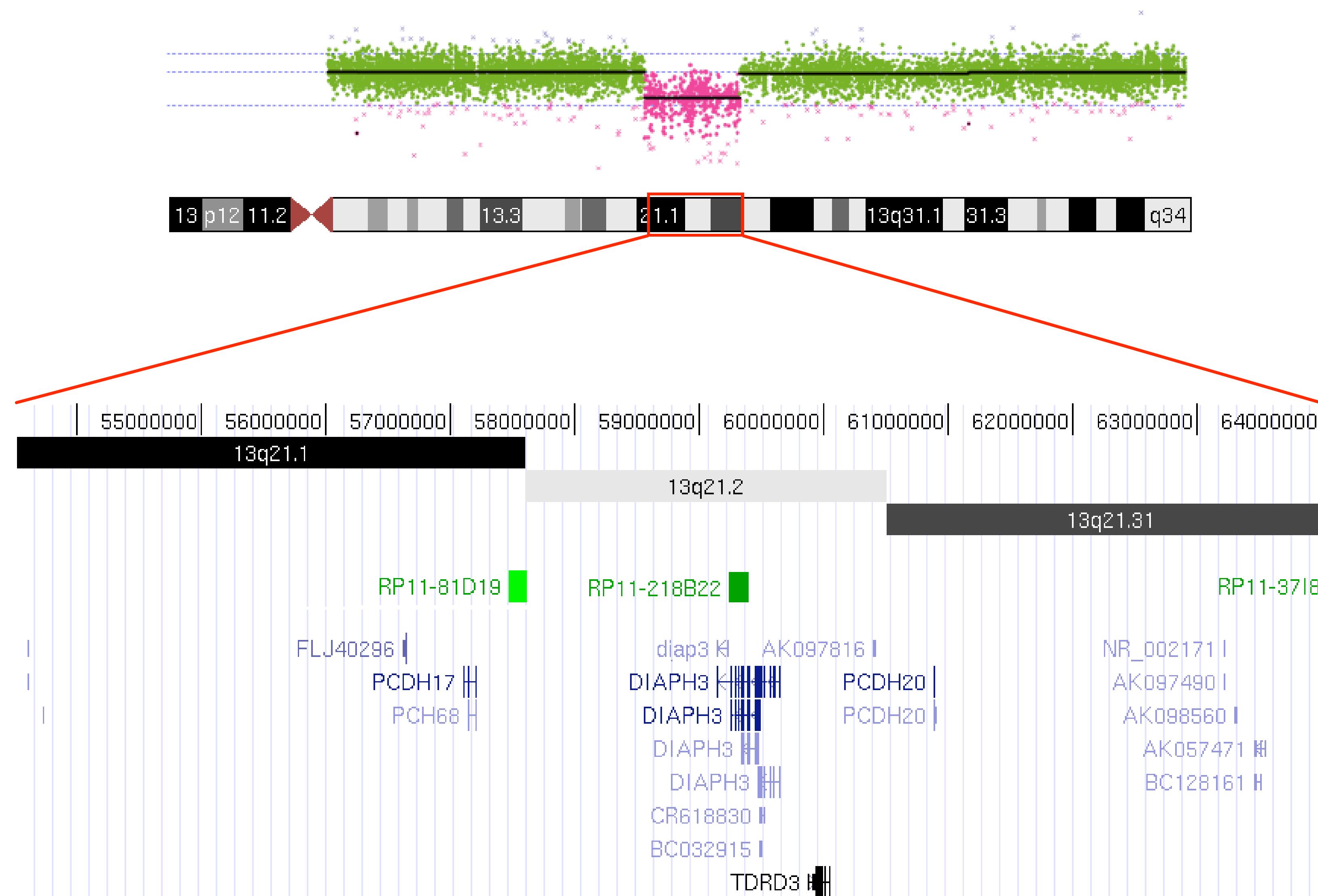


GENOMIC COPY NUMBER IMBALANCES PROVIDE WIDESPREAD SOMATIC VARIANTS IN CANCER

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

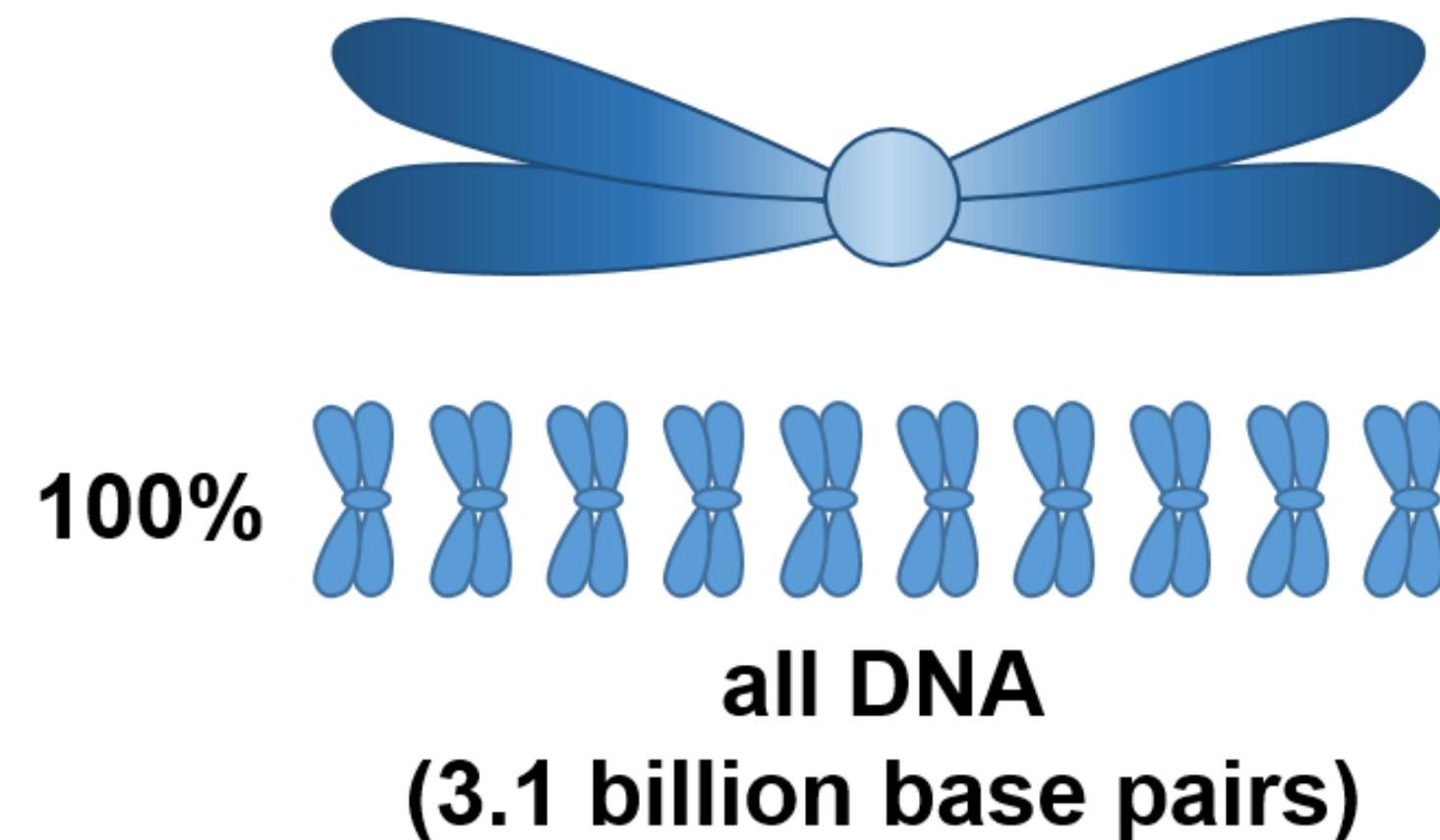
# Nobody is perfect (?)

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)

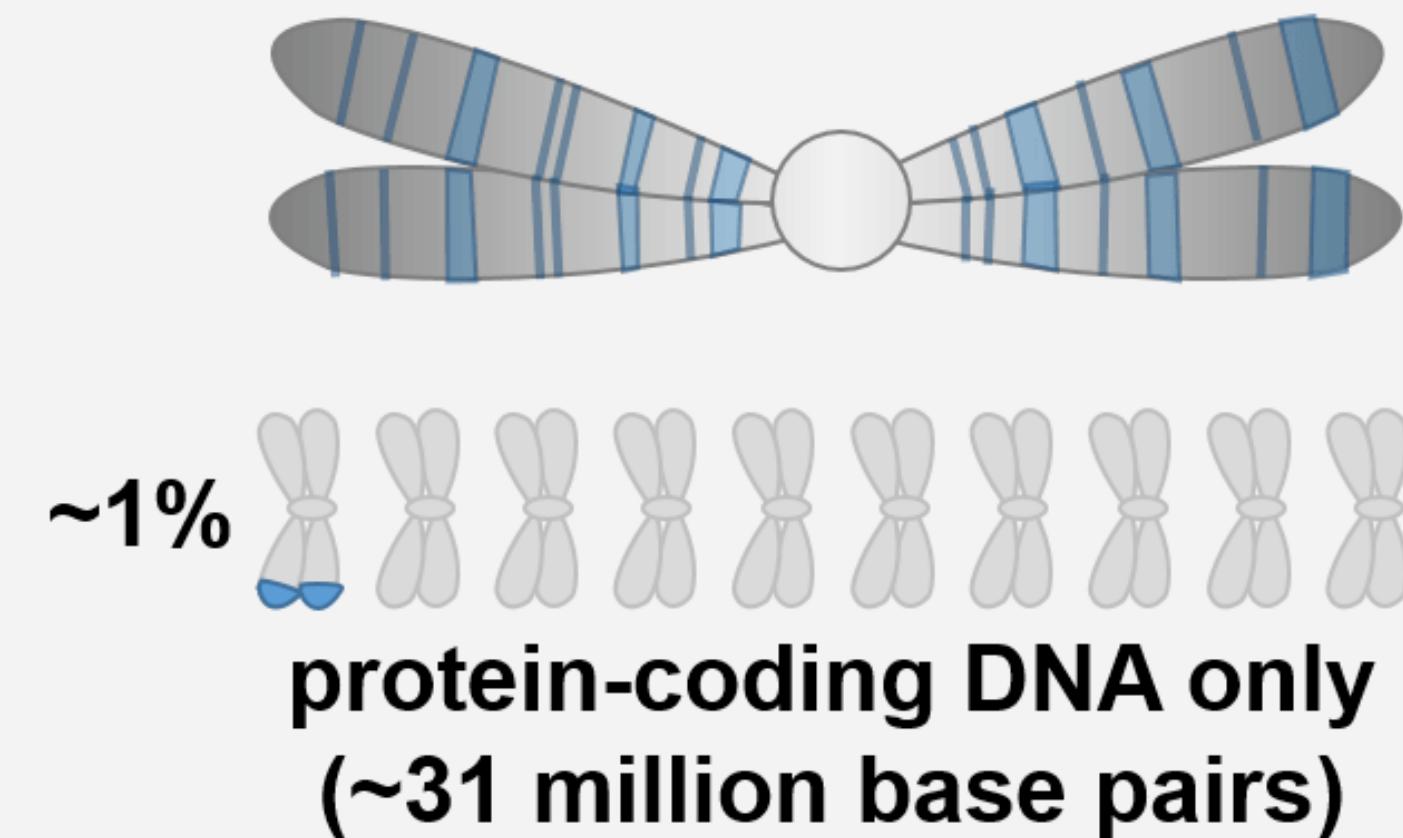


# Genome Sequencing

## 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:

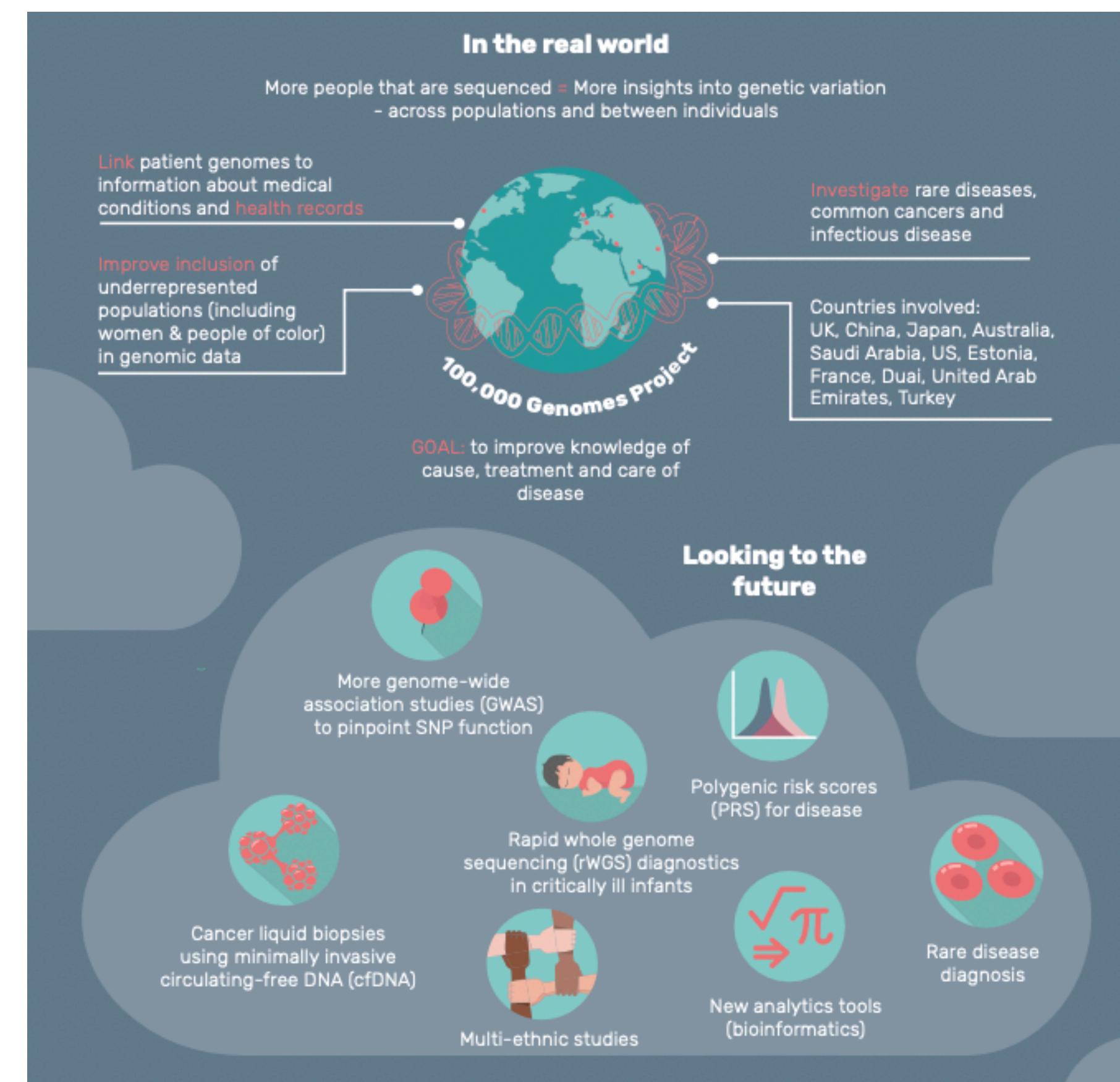
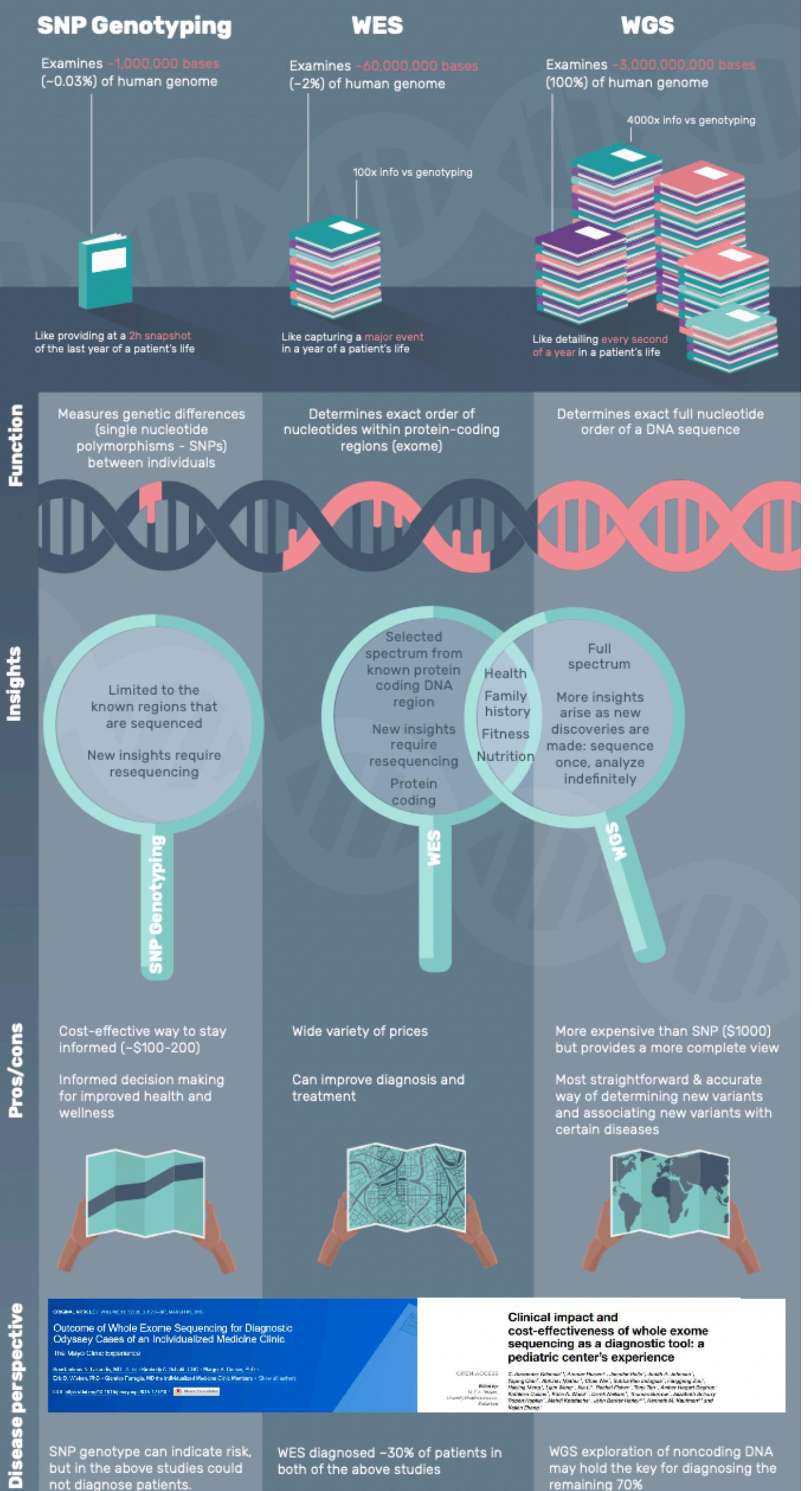
2017

~\$530

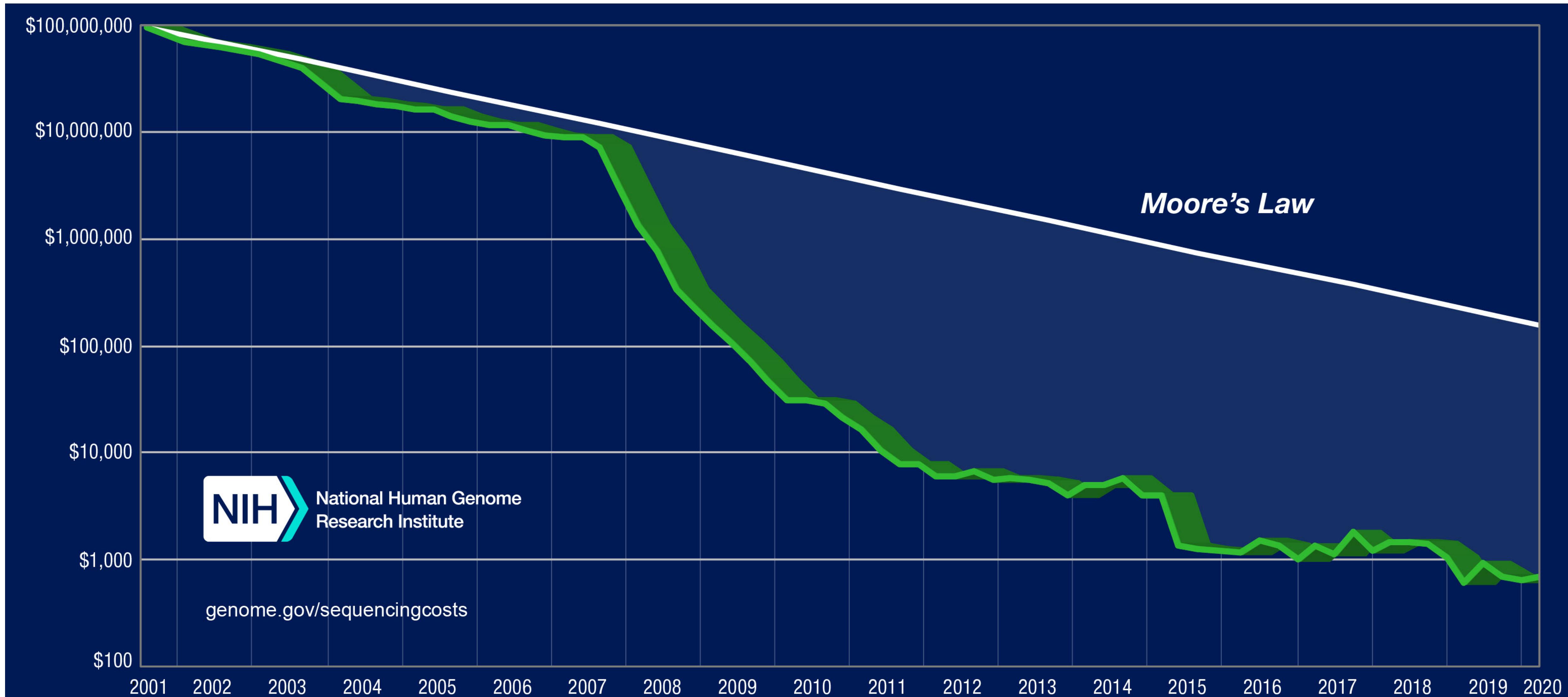
~3 days

# Genome Analysis

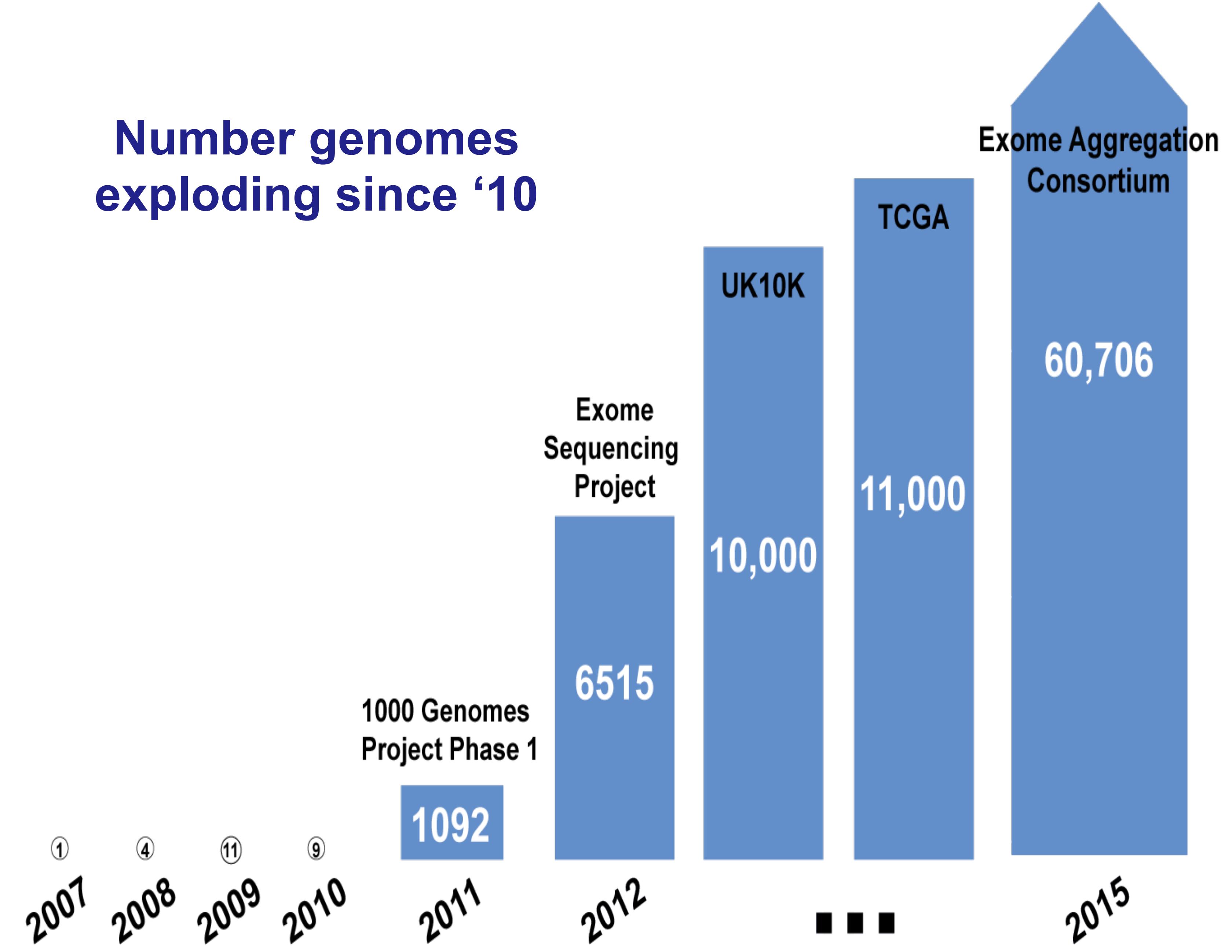
## A “progressing technologies” view



# The Cost of Sequencing a Human Genome



## Number genomes exploding since ‘10



# Genomes Everywhere

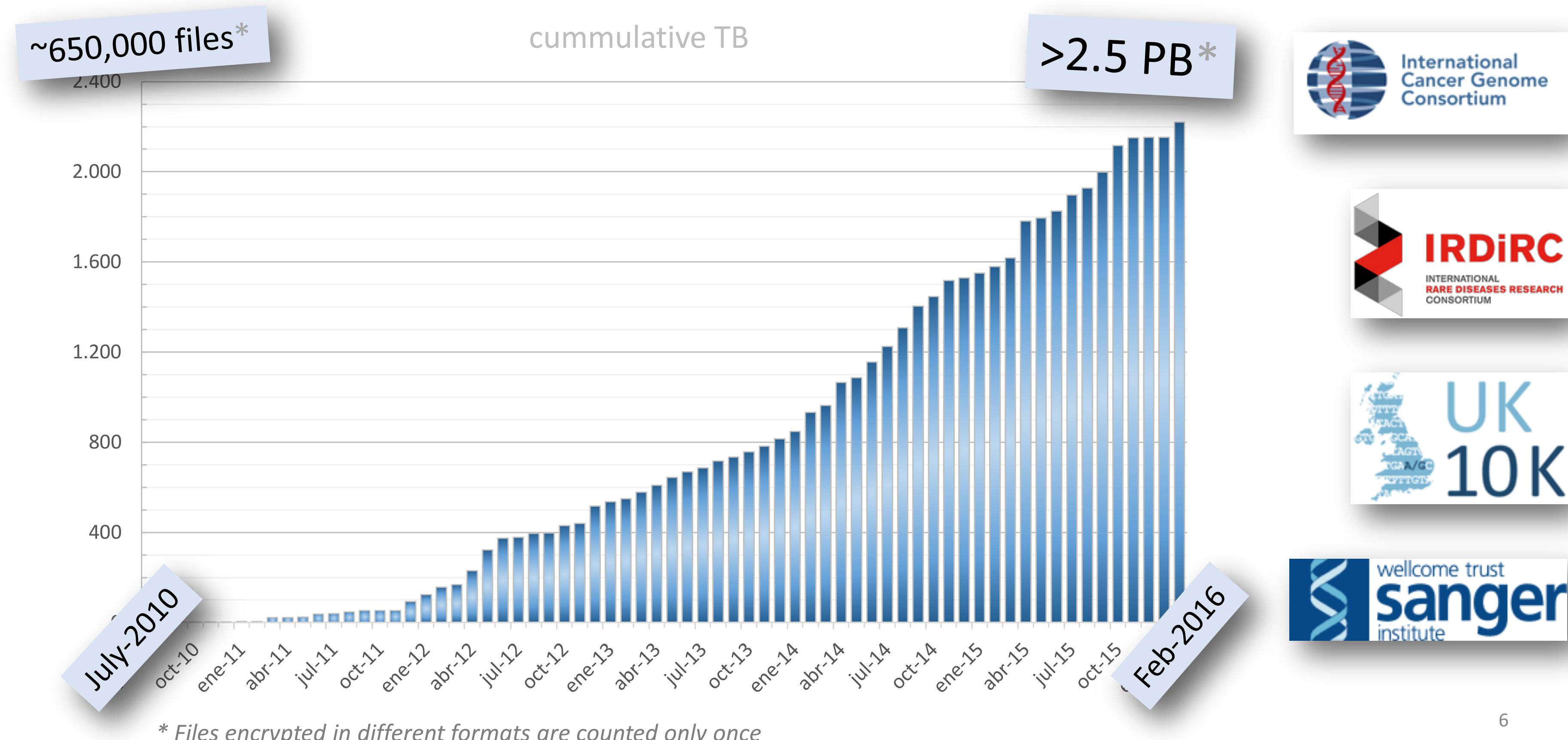
Organization / Initiative: Name	Organization / Initiative: Category	Cohort
100K Wellness Project	Research Project	107 unaffected individuals (scaling up to 100,000)
23andMe	Organization	>1 million customers (>80% consented to research)
Actionable Cancer Genome Initiative (ACGI)	Data-Sharing Project	Goal: 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

# National Medical Genome Projects and Cohorts (2018)



# Growth of Genome Data Repositories: Example EGA

The EGA contains a growing amount of data



# Total worldwide sequencing capability?!

## Deployed sequencers of major platforms

- deployed sequencers of major platforms (estimate end of 2021)
- "While 3 Exabases seems like a awful lot, it's worth noting that this still isn't enough to sequence every human born.  
Somewhere in the region of 4 children are born a second, our 100Gb/s wouldn't even let us fully sequence one of them (at 30x)" - Nava Whiteford

Platform	Estimated Instruments	Runs/Week	Total Runs/Year	Run Yield (Tb)	Total Sequencing Capacity Tb/year
ONT MinION	5501	3	858156	0.05	42907.8
ONT GridION	782	3	121992	0.25	30498
ONT PromethION 48	67	2	6968	14	97552
Illum Novaseq 6000	1485	3	231660	6	1389960
Illum NextSeq	5430	3	847080	0.36	304948.8
Illum Miseq/Mini/iSeq	12340	3	1925040	0.015	28875.6
Ion Torrent	2220	14	1616160	0.05	80808
PacBio	577	5	150020	0.03	4500.6
MGI - Mid/Low	2000	3	312000	0.72	224640
MGI -T7	10	5	2600	6	15600
Total	30412		6071676		2220290.8

Tb/year

2331523584 Samples

73.9 Samples/s

70.4 Gb/s

Estimated sequencing capability if everything would run continuously at full speed...

# Task: Reading up on Genome Technologies

- General NGS technologies
- count based vs. intensity based as principle
- bonus: dig deeper for some molecular-cytogenetic techniques:
  - SNP, aCGH arrays
  - SKY, M-FISH
  - chromosomal CGH

# BIO392: Course Schedule

<https://compbiozurich.org/courses/UZH-BIO392/>

	Tue Sep 20	Wed Sep 21	Thu Sep 22	Fri Sep 23	Tue Sep 27	Wed Sep 28	Thu Sep 29	Fri Sep 30	Tue Oct 4	Wed Oct 5	Thu Oct 6	Fri Oct 7	Tue Oct 11	Wed Oct 12
09:00 - 10:00			Izaskun: Terminal / Unix / Files	Izaskun: File formats for human genetic variation / file handling		Michael lecture introduction to some resources, CNVs, Progenetix	Hangjia: Blast	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
10:00 - 11:00		Github exercise: create user specific directories & upload/edit test files using Markdown (Ziying)	Izaskun: Terminal / Unix / Files	Izaskun: File formats for human genetic variation / file handling		Task: Browse/explore genome resources and provide some notes (1-2 pages total) in a doc posted on Github (.md)	Hangjia: blast	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
11:00 - 12:00		Ziying: github desktop and terminal		Izaskun: File formats for human genetic variation / file handling			Hangjia: Blast exercise	Max: STRs		Rahel: Sequence analysis	Rahel & Feifei: Survival analysis	Rahel & Feifei: Survival analysis		Exam
13:00 - 14:00	* Room information * Administrative - discuss times/days - exam	Ziying: Introduction to different interfaces eg atom, jupyter, pycharm (lecture), include R things	Izaskun: SIB online introduction to Unix	Izaskun: short project (1000 genomes), reading, literature	Recap W1; Q&A	Hangjia: Progenetix as tool for CNV frequencies etc.	Hangjia: Clinvar and Clingen	Max: STRs	Survival lecture, explanations of terms used etc, cancer classifications	Rahel: Survival analysis	Rahel & Feifei: Survival analysis	Discussion of Survival results (groups of groups)	Exam revision, Q&A	
14:00 - 15:00	Tina Siegenthaler: technical introduction (room, computer, accounts)		Izaskun: SIB online introduction to Unix	Izaskun: short project (1000 genomes), reading, literature	Literature (genome analysis techniques ...)			Max: STRs		Rahel	Rahel & Feifei: Survival analysis			
15:00 - 16:30	* explore course site * create Github accounts and forward to bio392@compbiozurich.org * Michael: short introductory lecture about genome variation			Izaskun: short project (1000 genomes), reading, literature	Genome technologies - brief notes about usage scenarios, pro & con			Max: STRs		Rahel	Rahel & Feifei: Survival analysis			