

Bioinformatics and R: Visualising Genomic Data



Yari Ciani, PhD

Laboratory of Computational and Functional Oncology

Department for Cellular, Computational and Integrative Biology - CIBIO

University of Trento

yari.ciani@unitn.it

The Topic

Precision Oncology and Biomarker discovery

Diagnostic, prognostic, treatment response markers for tumor stratification and precision oncology
Hypothesis driven and agnostic studies

The Team



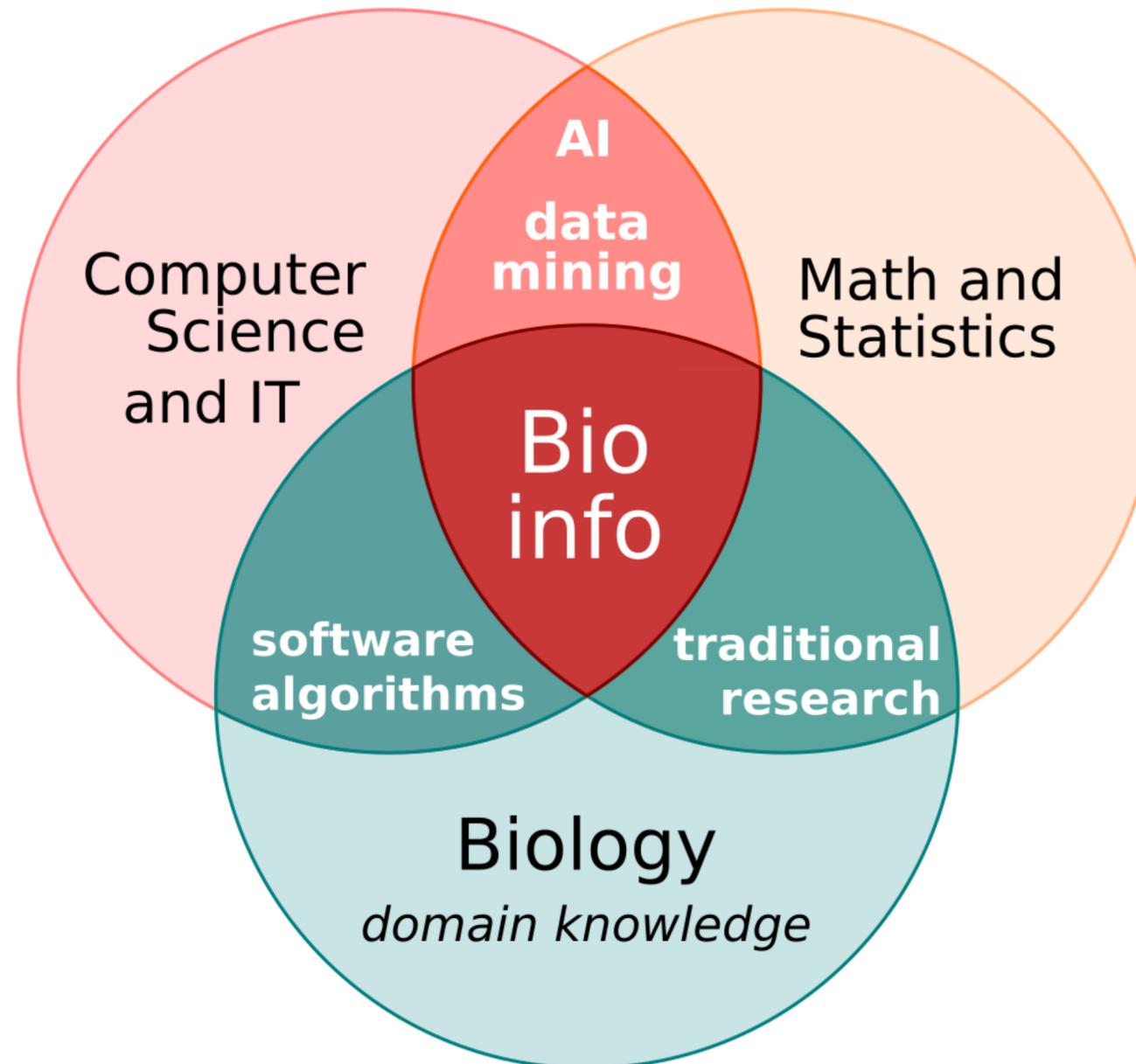
The Funding



What is Bioinformatics?

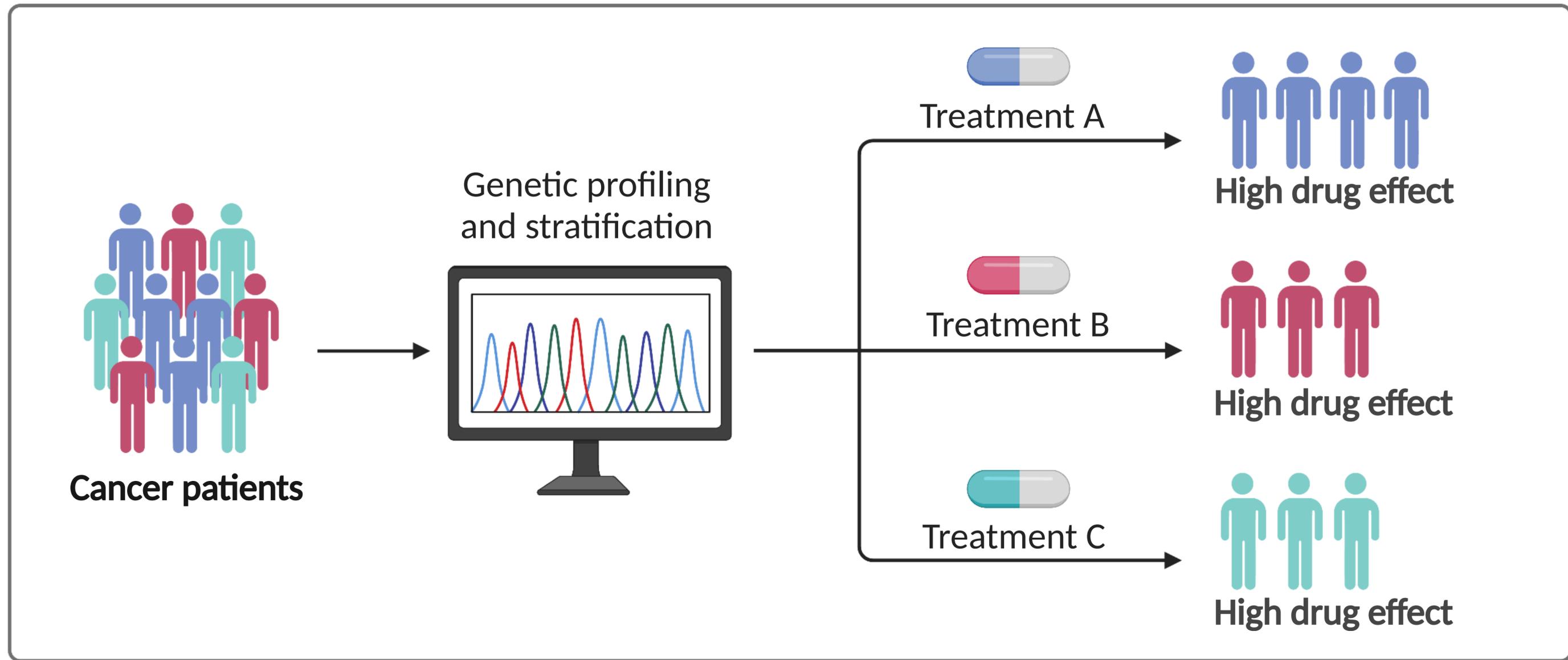
Bioinformatics: applies principles of information sciences and technologies to make the vast, diverse, and complex life sciences data understandable and useful.

e.g., Text mining,
Processing raw data,
Artificial Intelligence,
Software development,
but also **experiment design**
and **results interpretation**

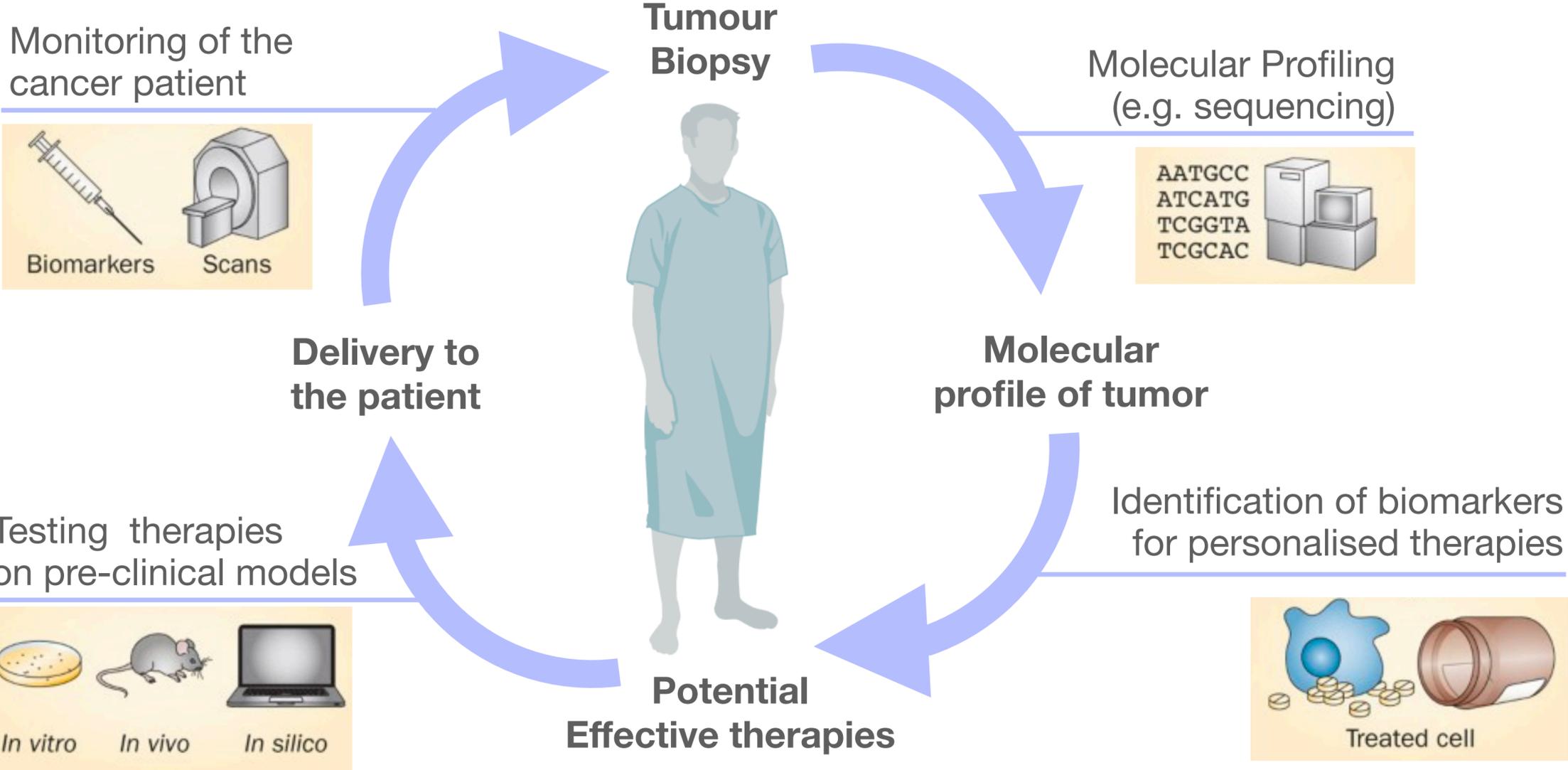


Precision Oncology

Precision cancer therapy



The Precision Drug Discovery Cycle

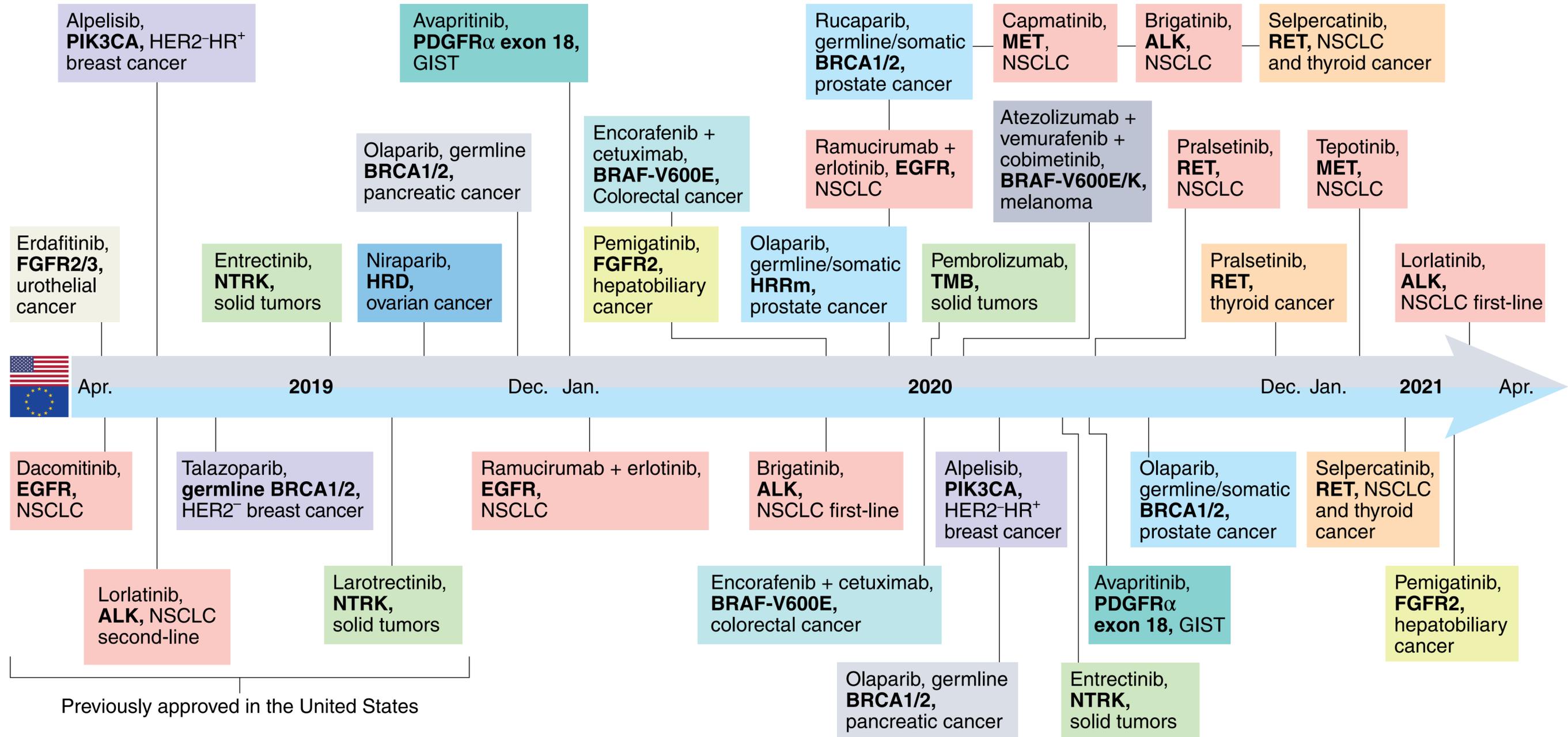


The cost of developing a single FDA/EMA-approved drug:
1B \$
and
10 - 15 years

FDA: Food and Drug Administration
EMA: European Medicines Agency

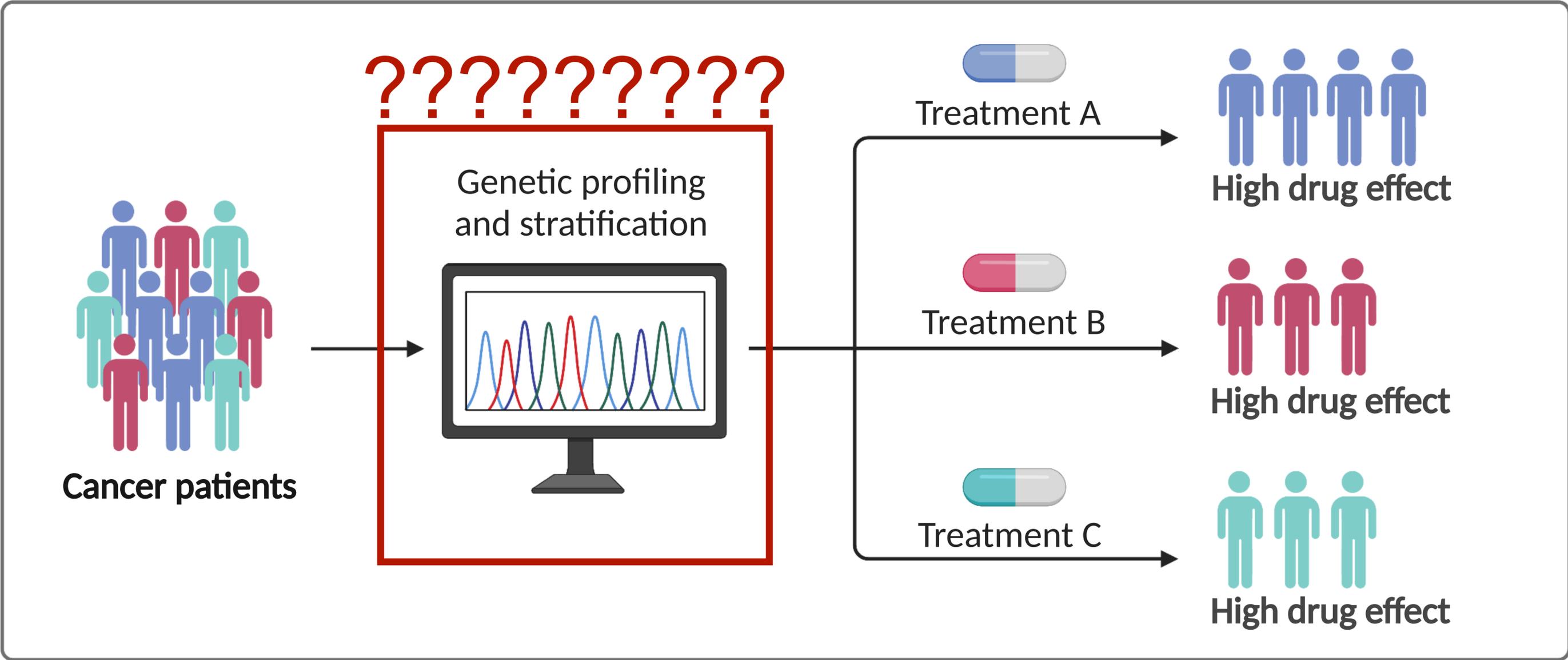
Adapted by T Cantore from Shrager, Jeff, and Jay M. Tenenbaum. "Rapid learning for precision oncology." *Nature reviews Clinical oncology* 11.2 (2014): 109-118.

...but it is worth



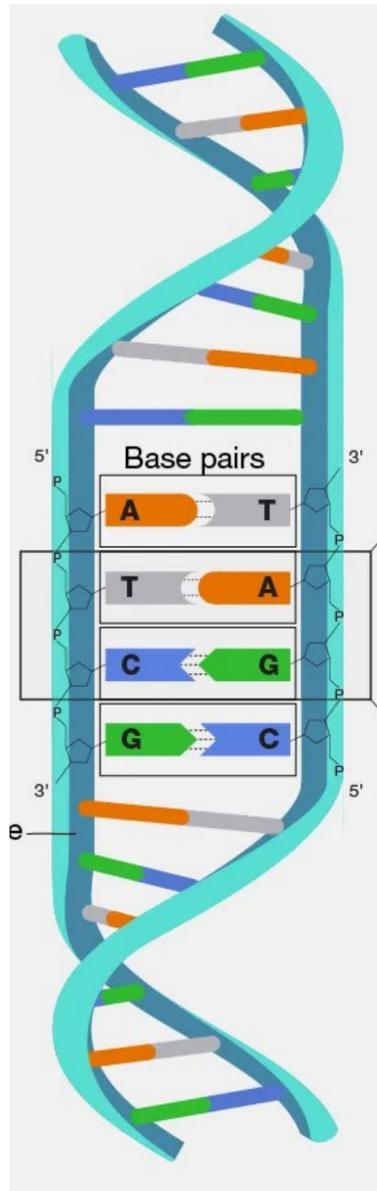
Precision Oncology

Precision cancer therapy

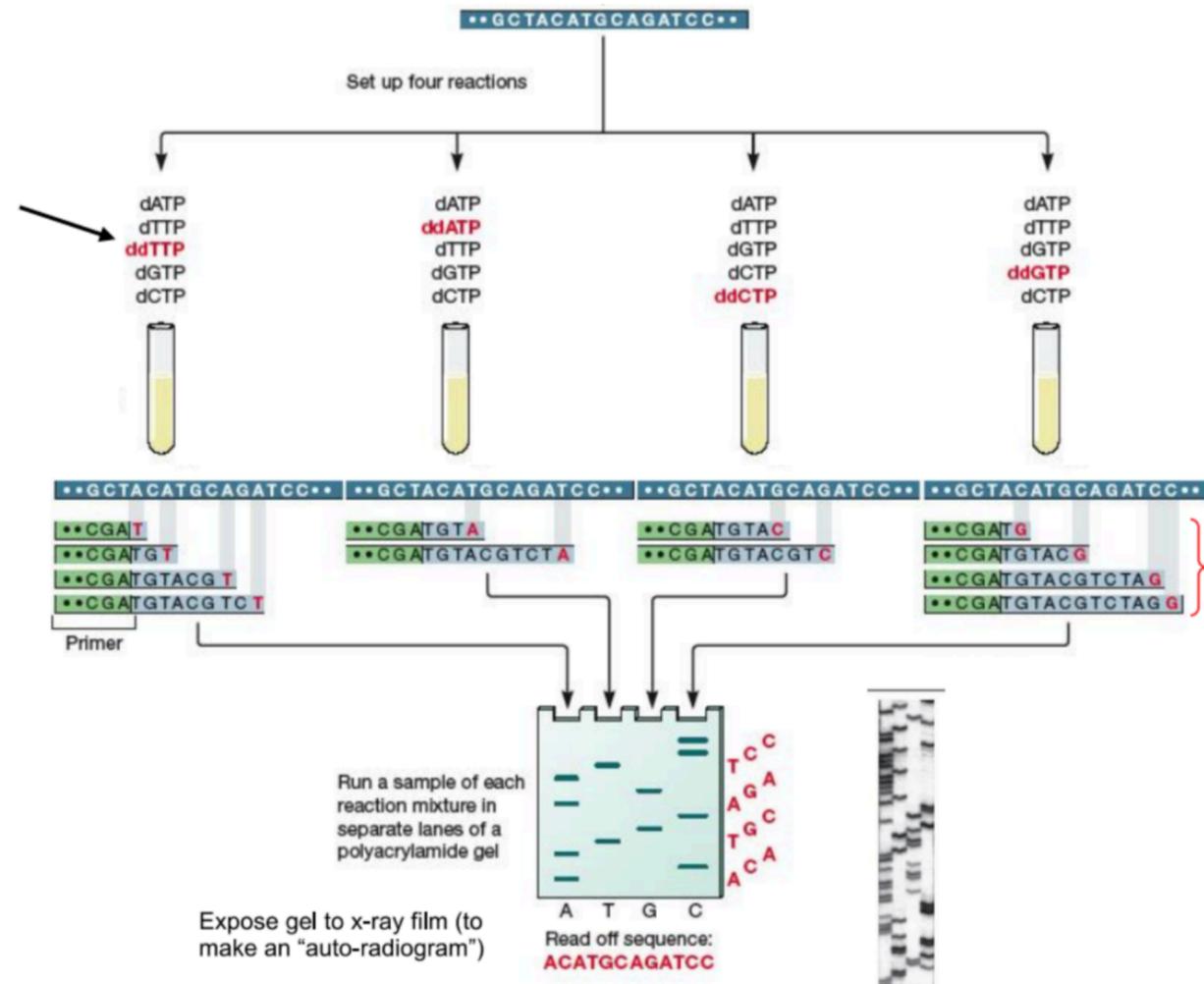


The origin of DNA sequencing

DNA Sequencing is figuring out the order of DNA nucleotides, or bases (A T G C), in a genome that make up an organism's DNA.



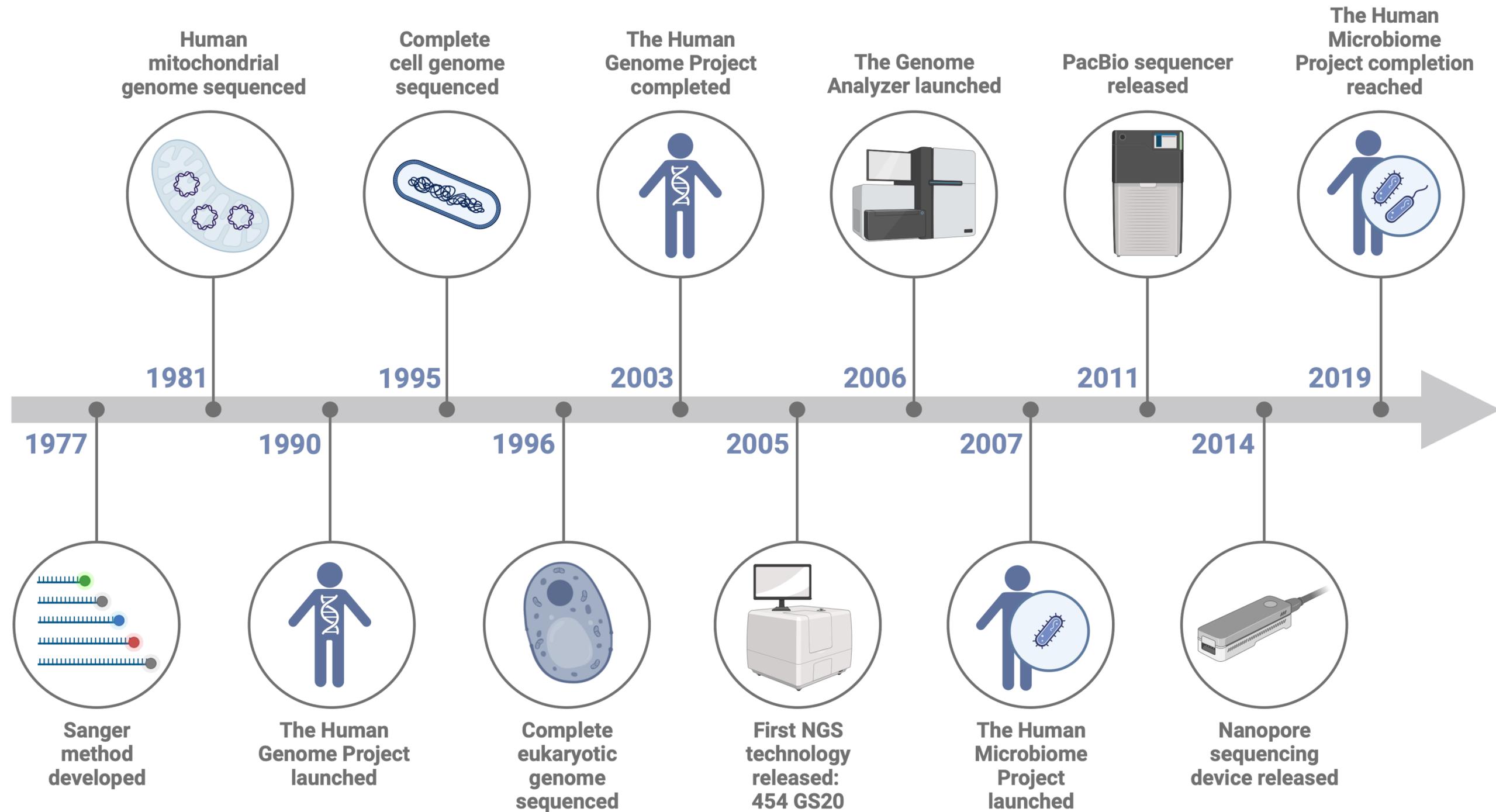
ddNTP are at low concentrations to permit elongation of fragments



The original Sanger sequencing method (1977)

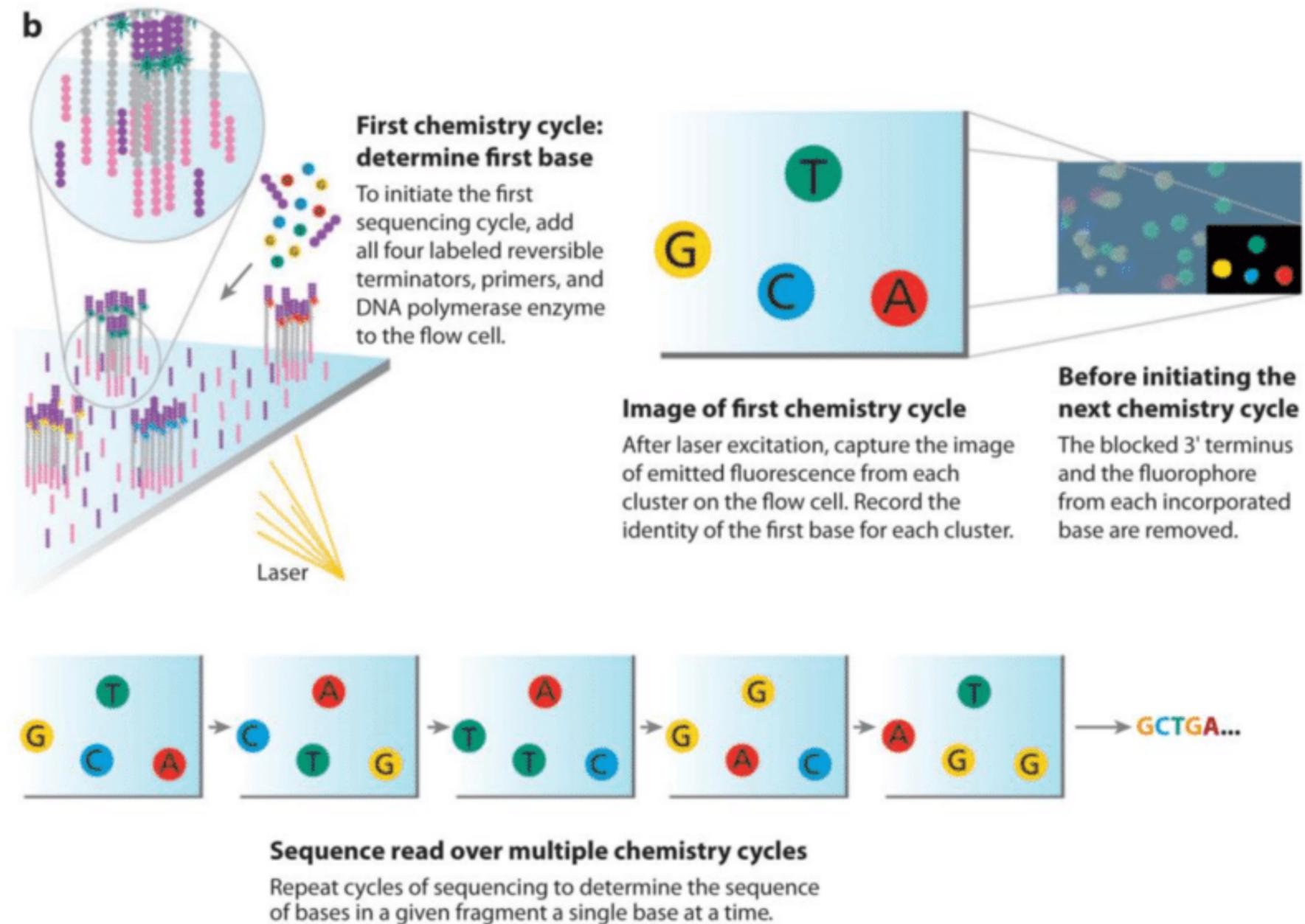
A nested series of DNA fragments ending in the base specified by the terminator-ddNTP

The NGS revolution



NGS approaches in a nutshell

Illumina - Sequencing by synthesis



How to get from NGS data to biological interpretation?



Up to 8B read pairs (2x150bp)

BASECALLING

ACGTCGATCGATCGATCGATCG
TCGATCGCGCGAGATGGCTGAA
CGAGCTAGCTAGCTGGCTAGAGCT
CAGCGAGCTAGCTAGCATCGAT
CGATGCTAGCTAGCTAGCTAGC

- Sequencing produces high-resolution TIFF images
- 100 tiles per lane, 8 lanes per flow cell, 100 cycles
- 4 images (A,G,C,T) per tile per cycle = 320,000 images
- Each *TIFF* image ~ 7Mb = 2,240,000 Mb of data (**2.24TB**)

How to get from NGS data to biological interpretation?

Up to 8B read pairs (2x150bp)

SEQUENCING READ

ACGTCGATCG**G**TTCGATCGATCG

Single nucleotide variant (SNV)

...CGATCGATCGGATCGAC**G**TTCGATCGATCGATCGATCGCGATCGATCGATCGG...

CHROMOSOME SEQUENCE
(reference genome)

Human Genome = 3.3B bases

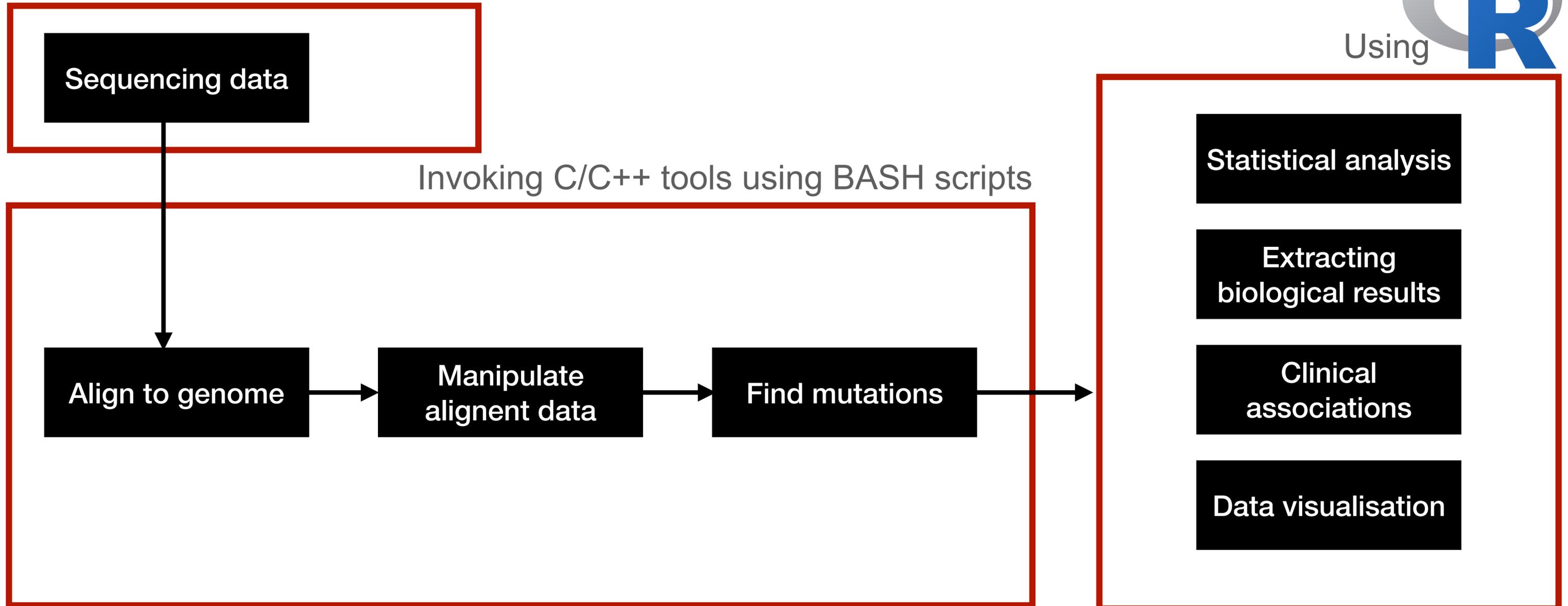
Approximate String Matching with Bounded Edit Distance, where the goal is to find all (locally) similar substrings of a large reference that align to a given read, allowing for a limited number of edits.

R is not optimized for low-level memory access or speed.

Genomic-scale alignment is typically done in **C/C++** (e.g., BWA, minimap2) due to performance constraints.

The role of R in Bioinformatics

Proprietary software



The R Bioinformatics community



[About](#) [Learn](#) [Packages](#) [Developers](#)

[Get Started >](#)

Home > **BiocViews**

Bioconductor version 3.21 (Release)

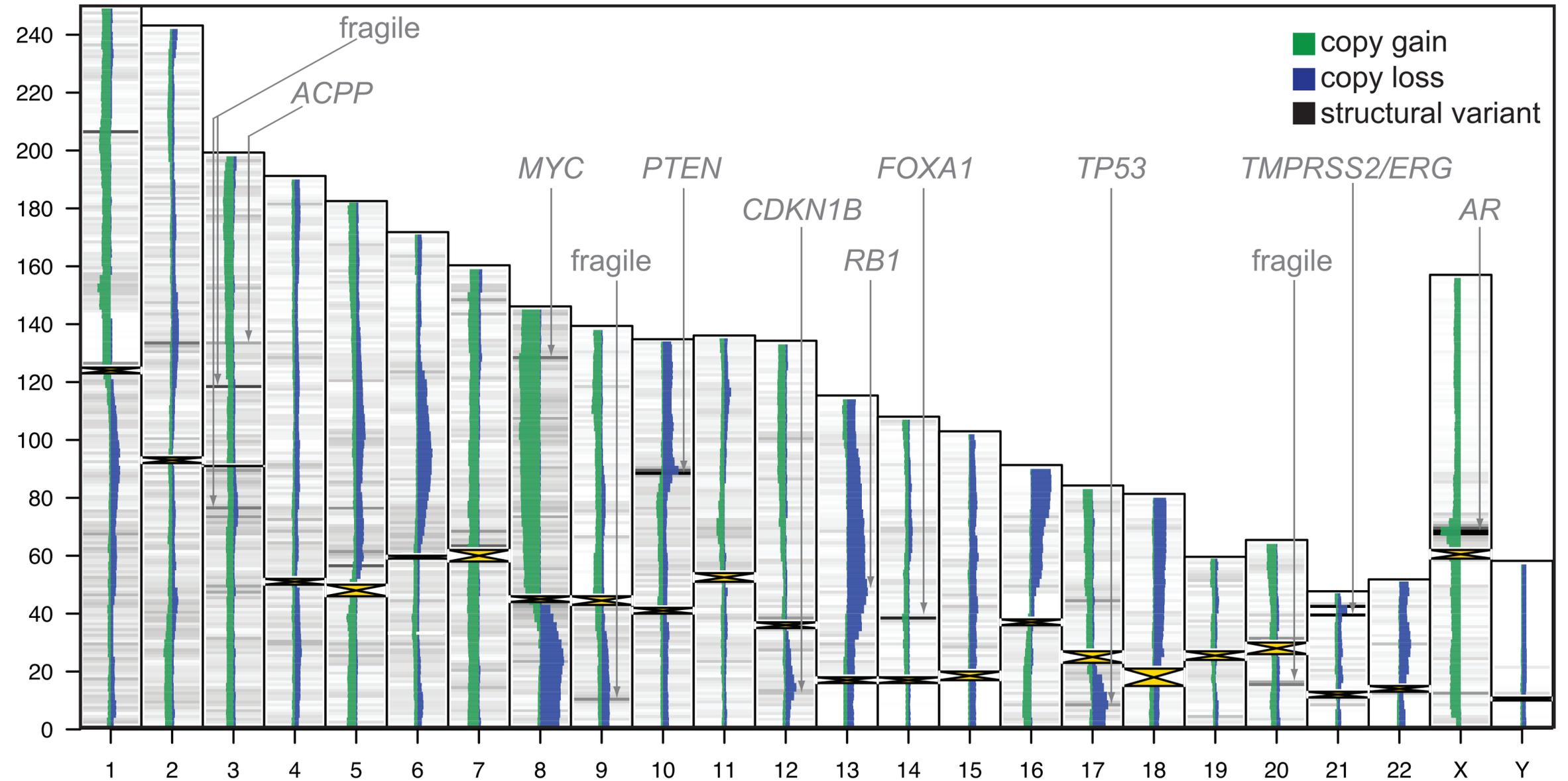
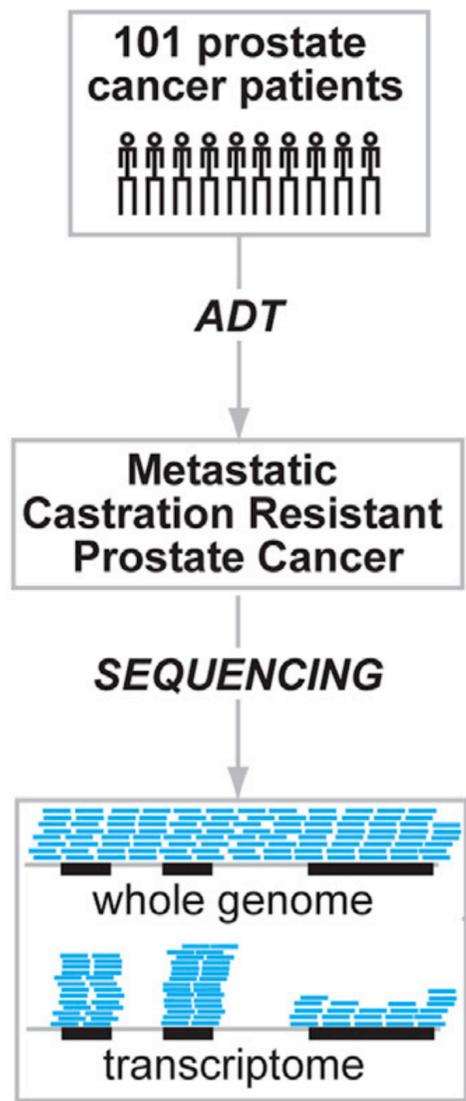
[Go to 3.22 \(Devel\) >](#)

Bioconductor uses the R statistical programming language, and is open source and open development. It has two releases each year, and an active user community.

Find biocViews:

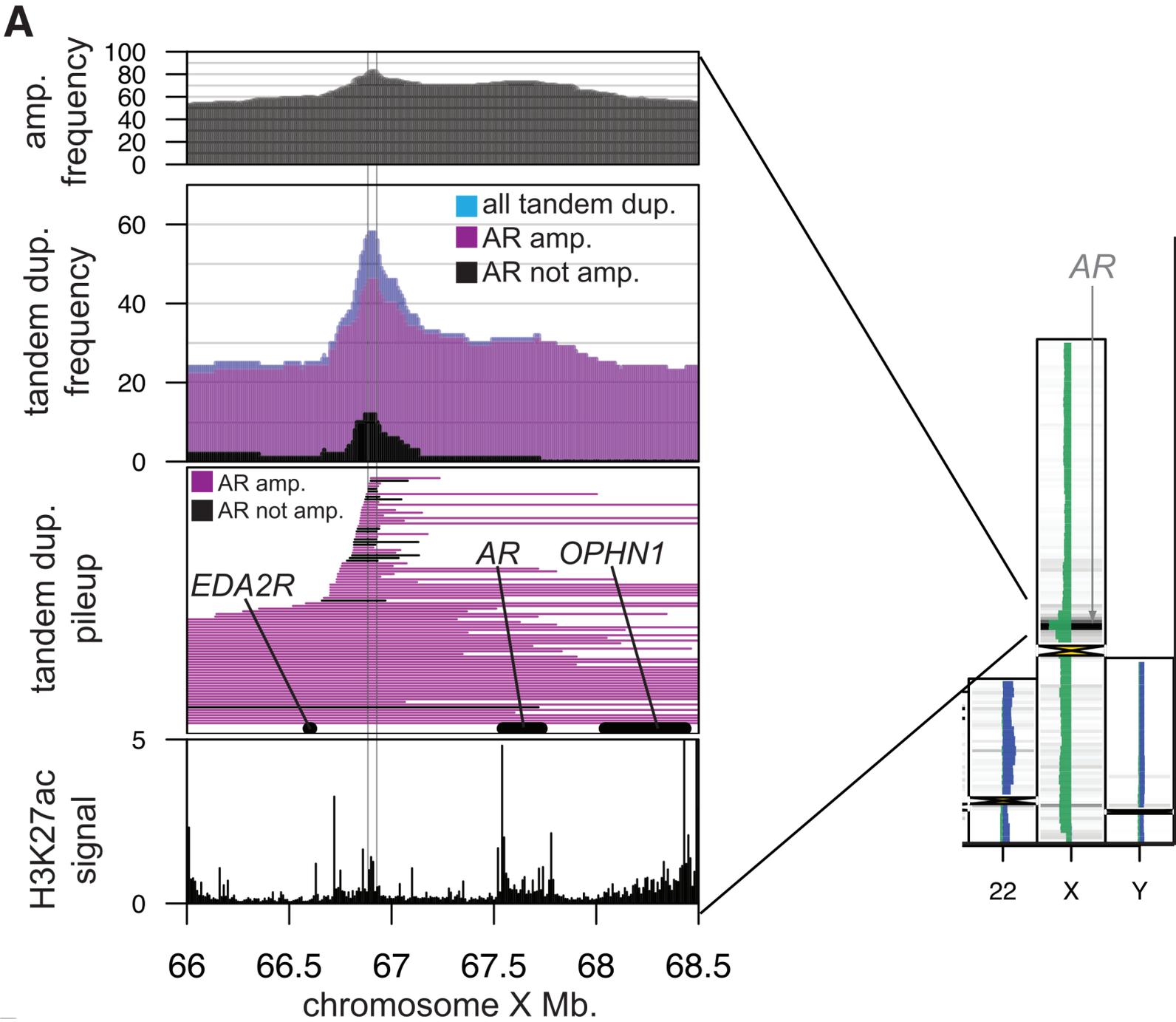
- ▶ **Software (2341)**
- ▶ AnnotationData (928)
- ▶ ExperimentData (432)
- ▶ Workflow (30)

Visualising the Genome: an example



Showing frequencies of mutations of an entire cohort across the entire genome

From Overview to Detail



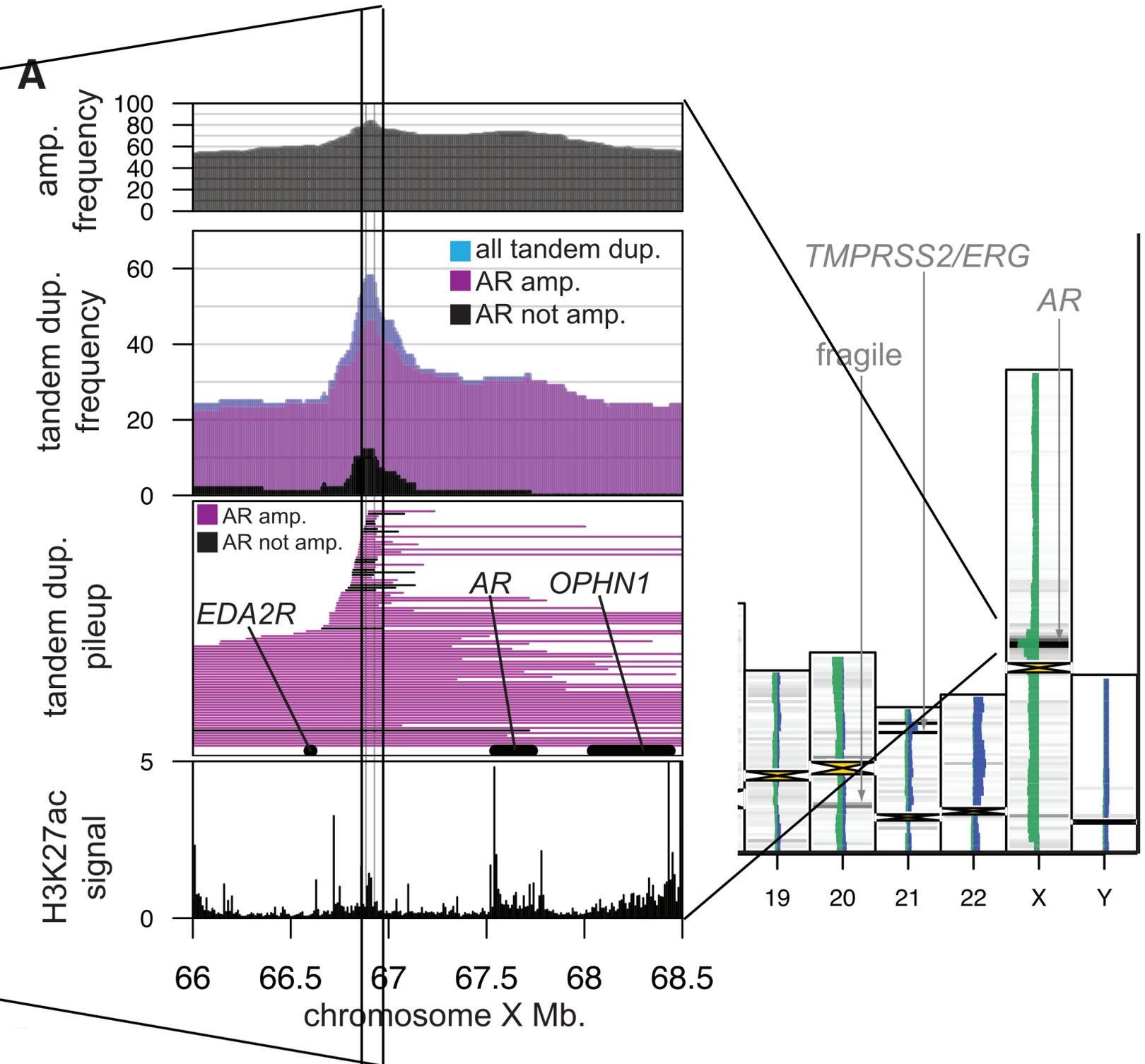
AR = androgen receptor. Gene involved in prostate cancer

From Overview to Detail



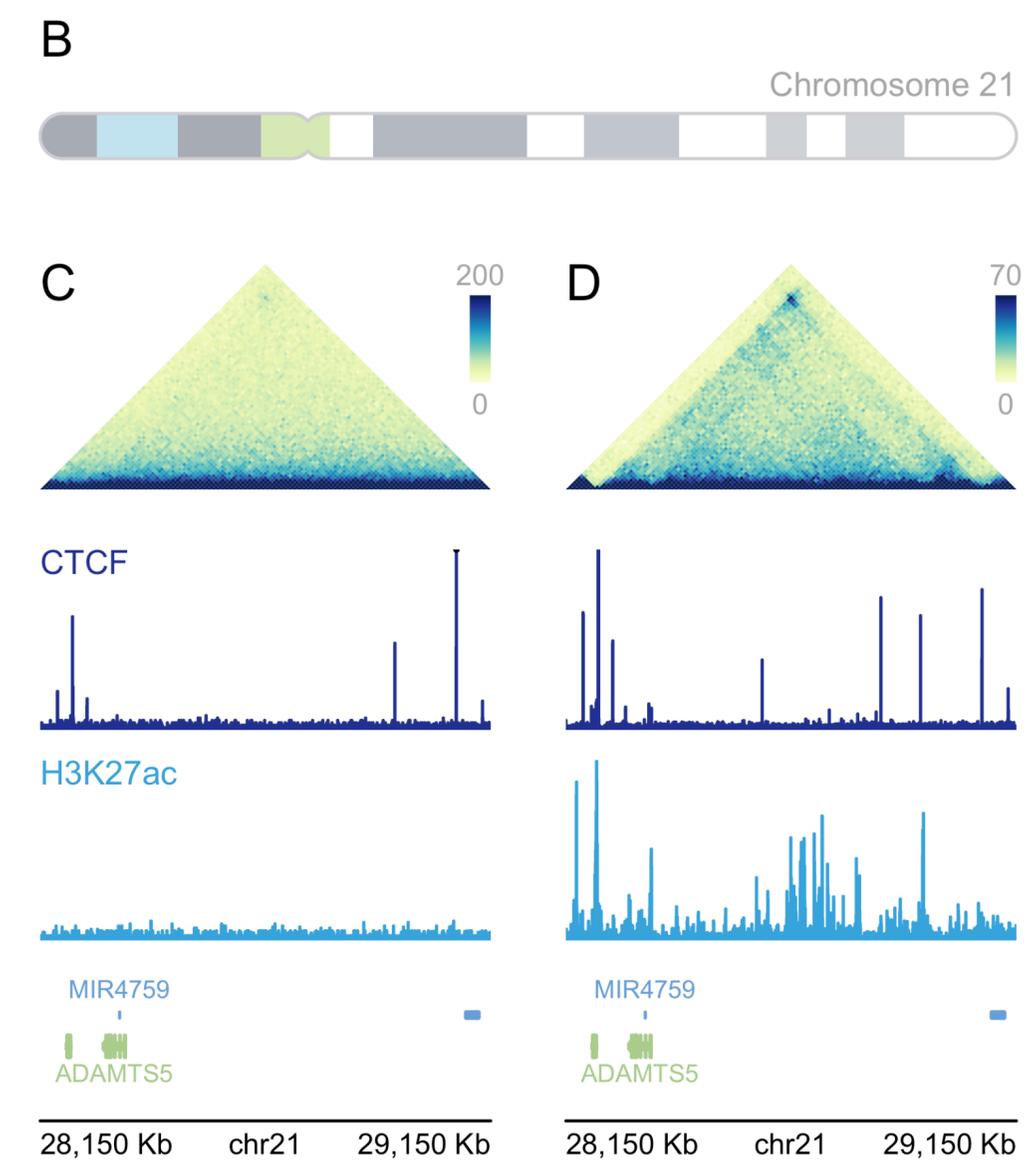
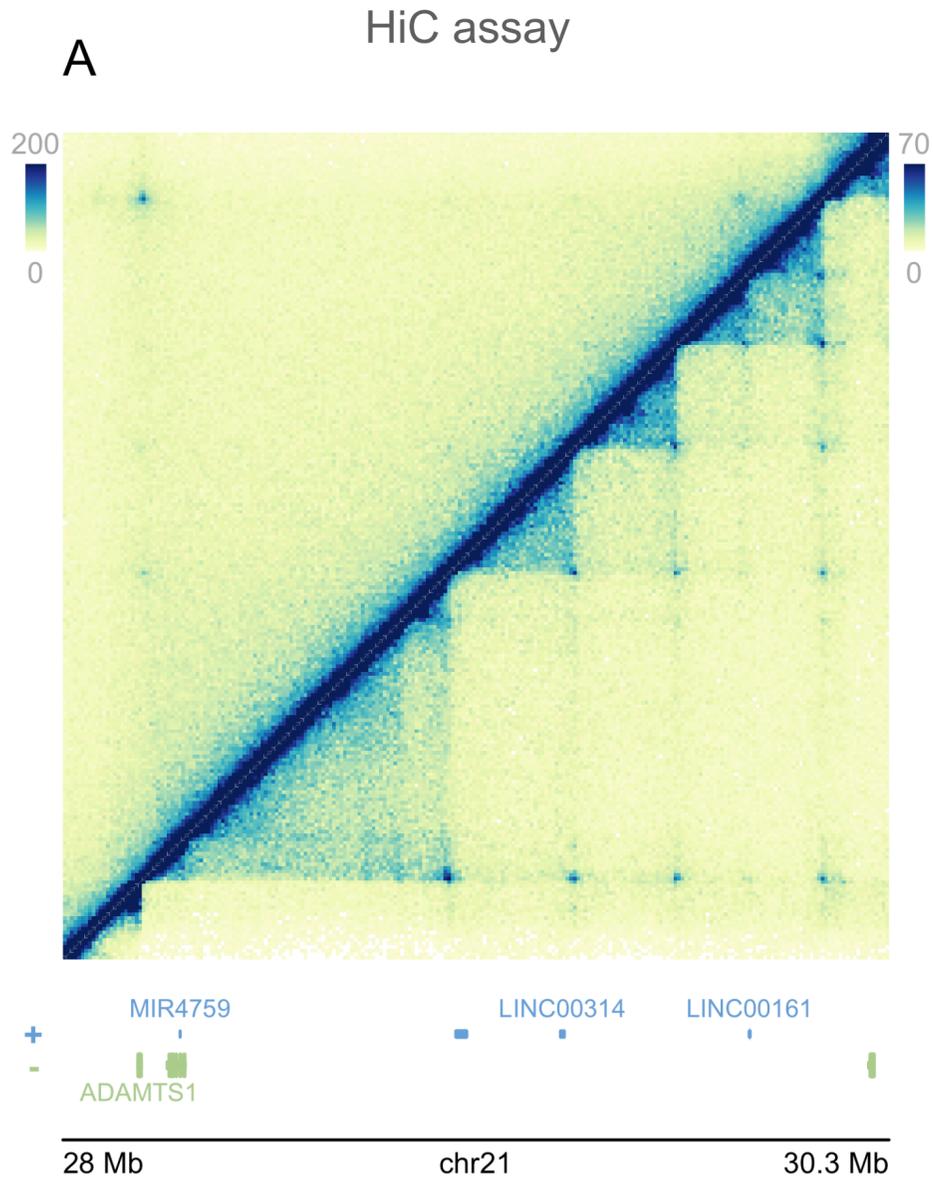
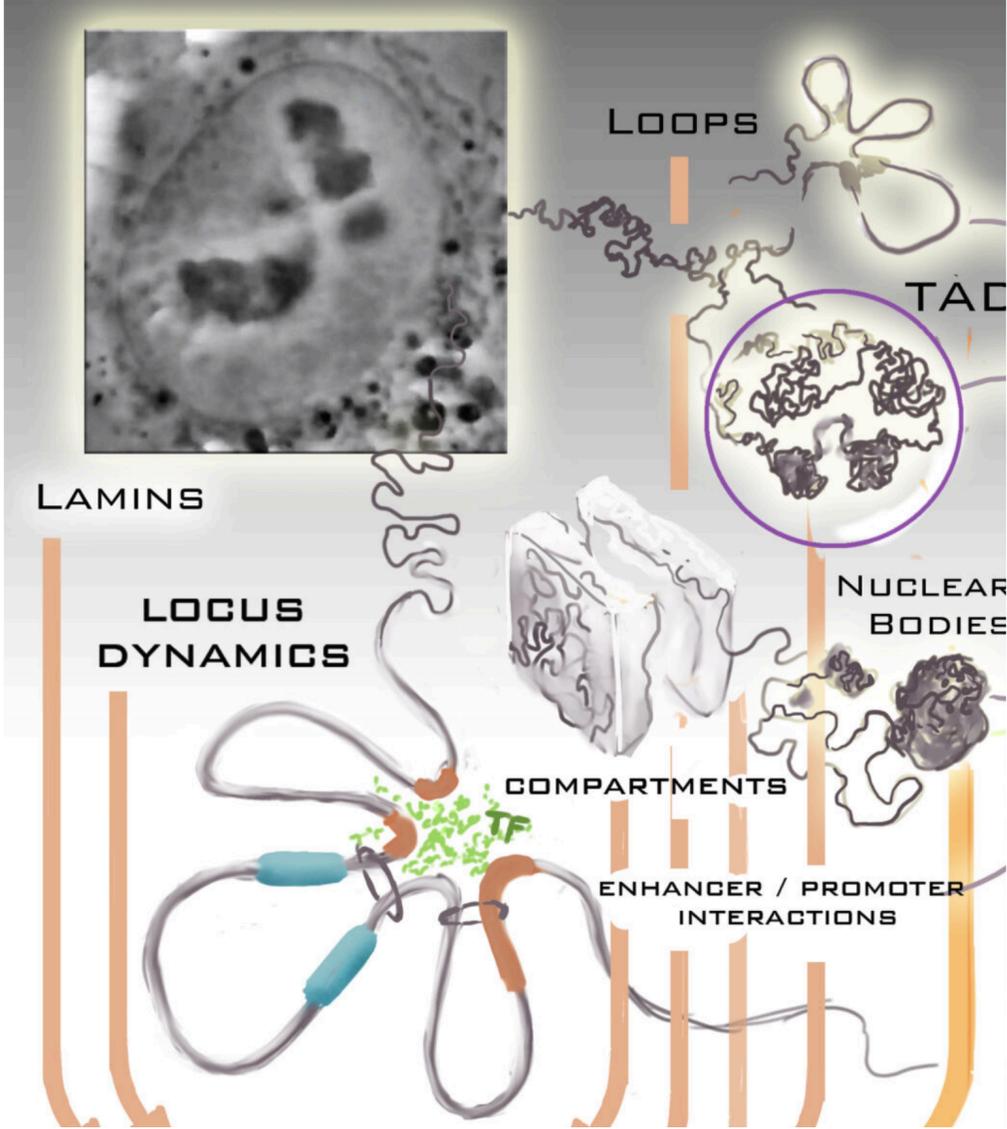
CONS:

- graphical interface
- slow
- graphical customisation
- JAVA!



Screenshot from IGV, a Java based alignment visualisation tool

Long Range Interactions (physical)



<https://4dnucleome.org/>

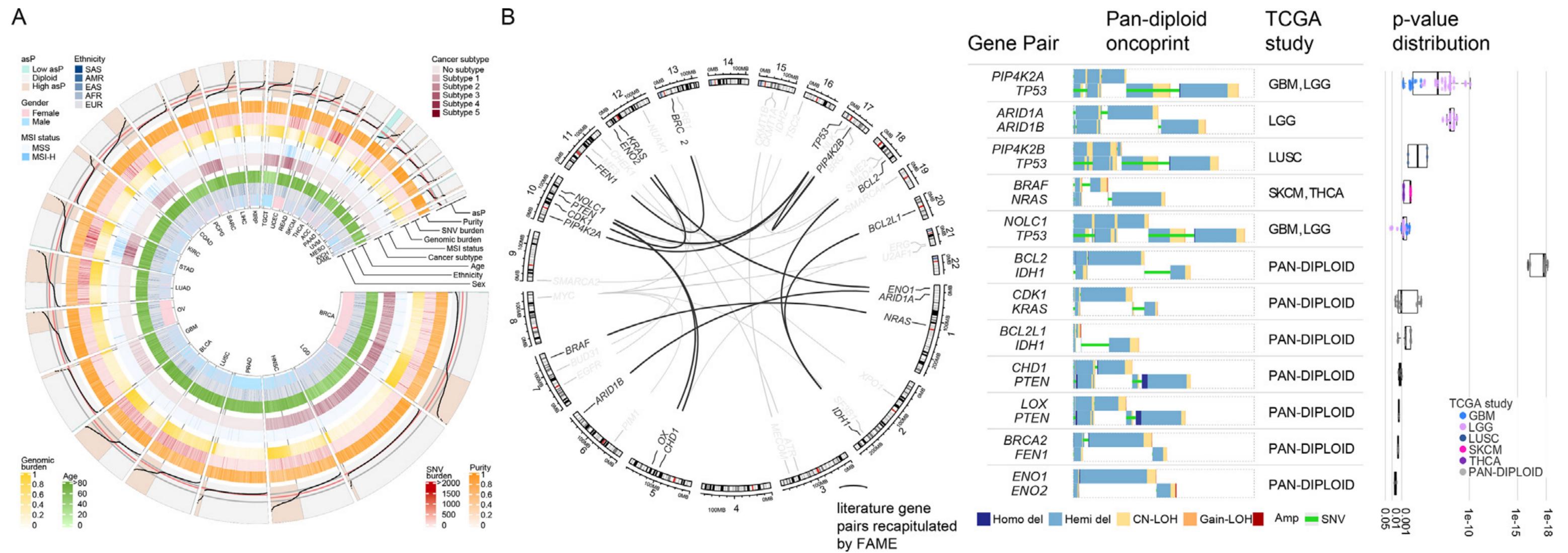
<https://phanstiellab.github.io/plotgardener/>



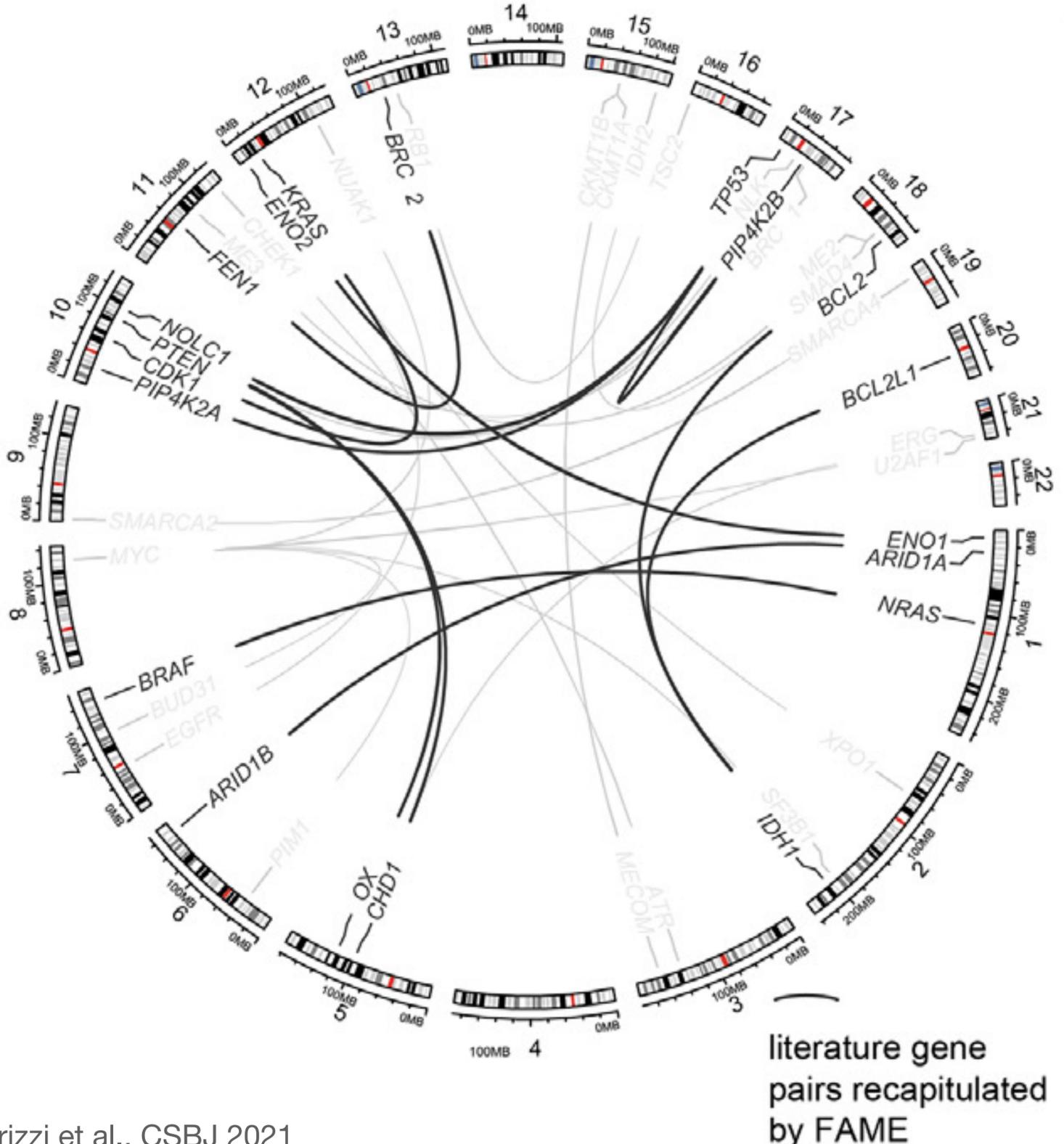
Long Range Interactions (functional)

T. Fedrizzi, Y. Ciani, F. Lorenzin et al.

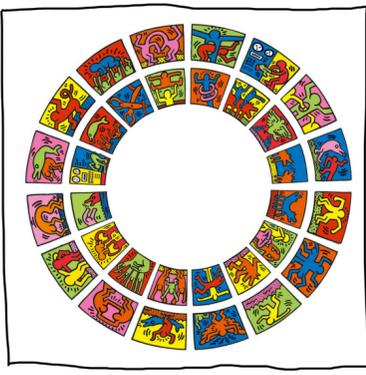
Computational and Structural Biotechnology Journal 19 (2021) 4394–4403



Long Range Interactions (functional)



circlize implements and enhances circular visualization in R Bioinformatics, 2014



Zuguang Gu

https://jokergoo.github.io/circlize_book/book/index.html

```

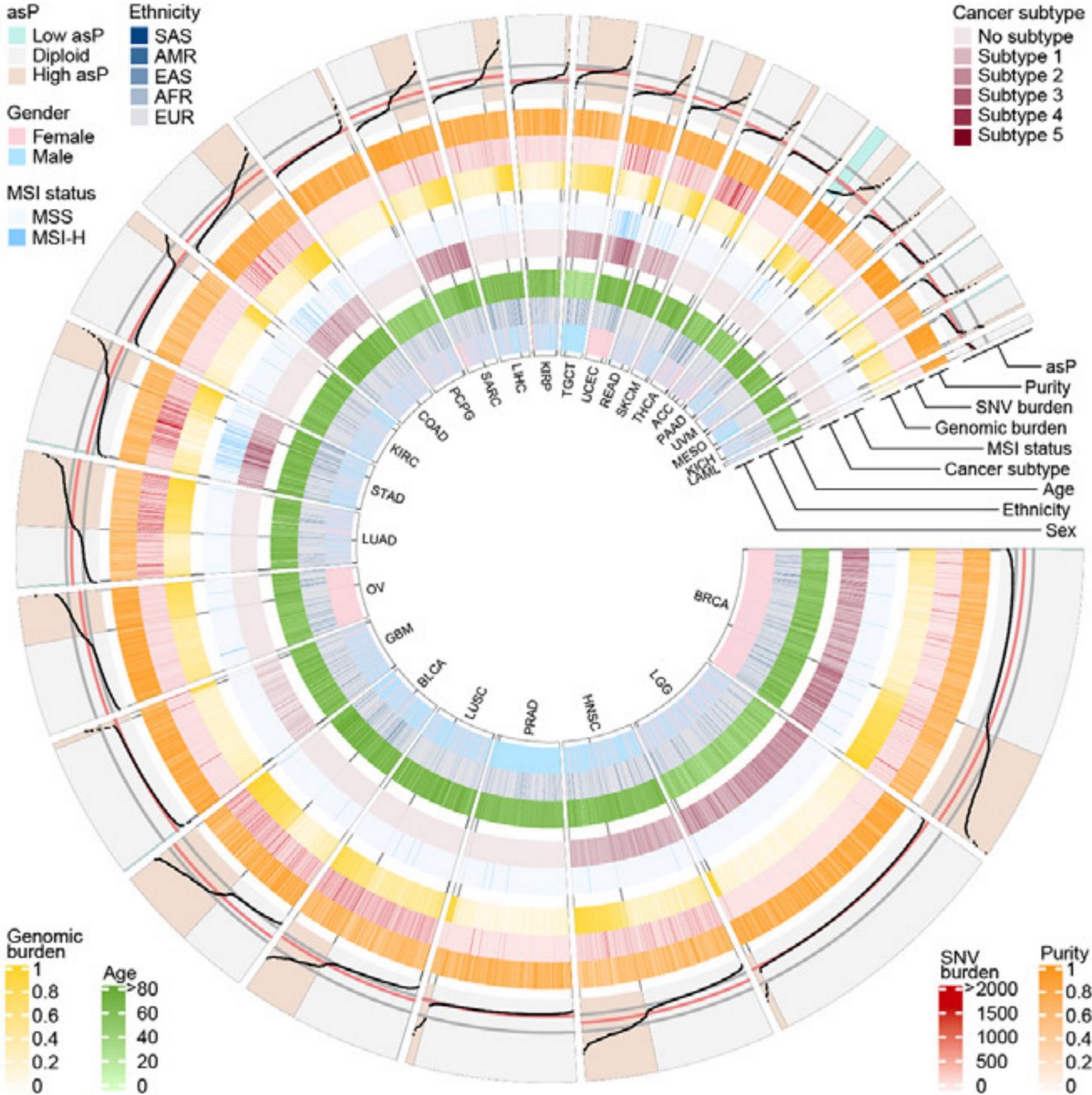
circos.initializeWithIdeogram()
circos.genomicLink.bed1, bed2)
    
```

bed1

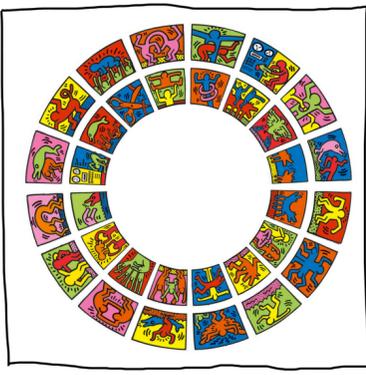
chr	start	end	value1
chr6	102324459	147617643	-0.50418830
chr17	65167455	77619820	-0.10264963
chr11	13366995	32331617	0.42482152
chr8	93343457	96256710	0.65620649
chr16	64403195	65047798	0.04966380

Summarising Entire Cohorts

~4800 patients



circlize implements and enhances circular visualization in R
Bioinformatics, 2014

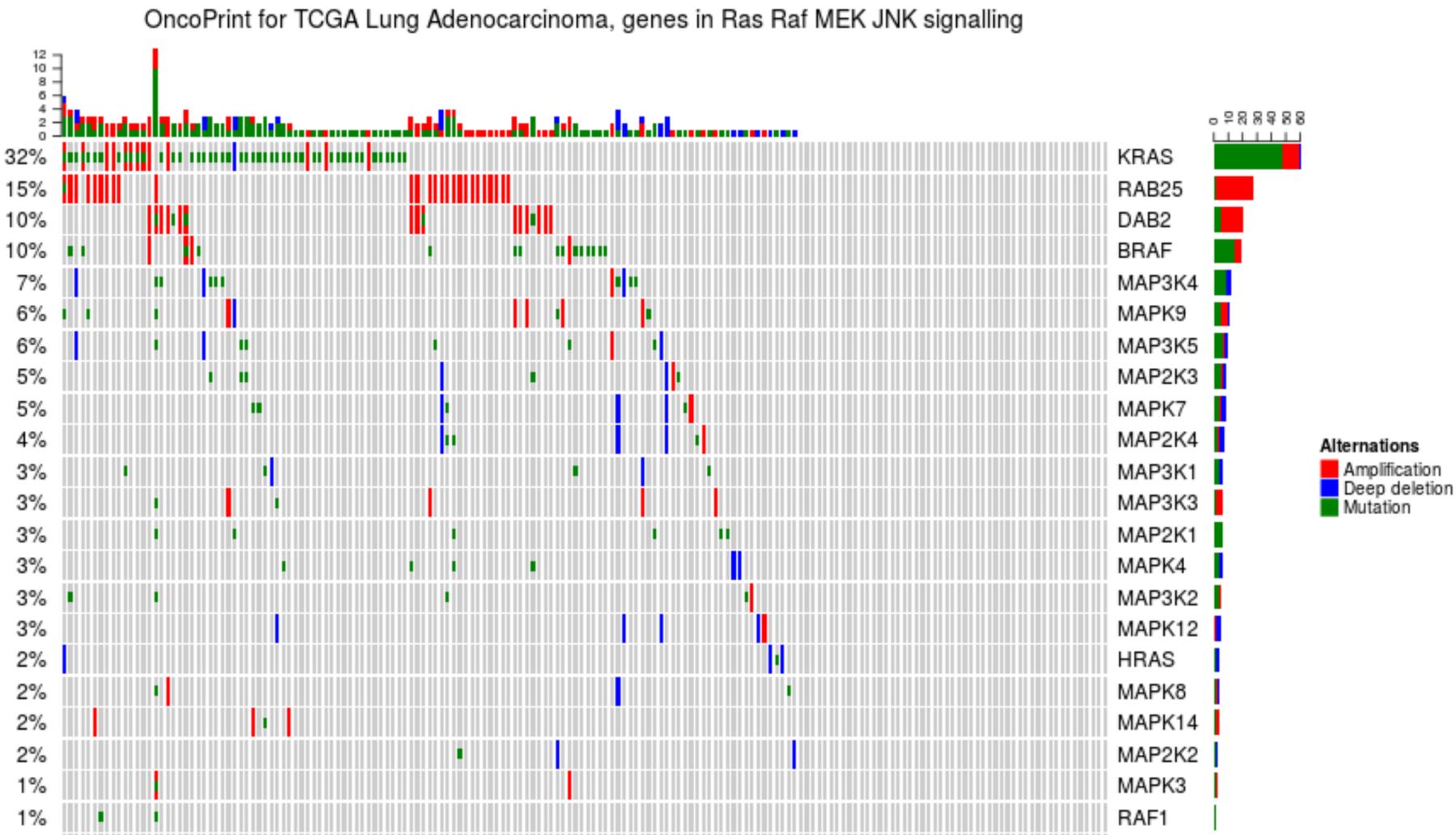


Zuguang Gu

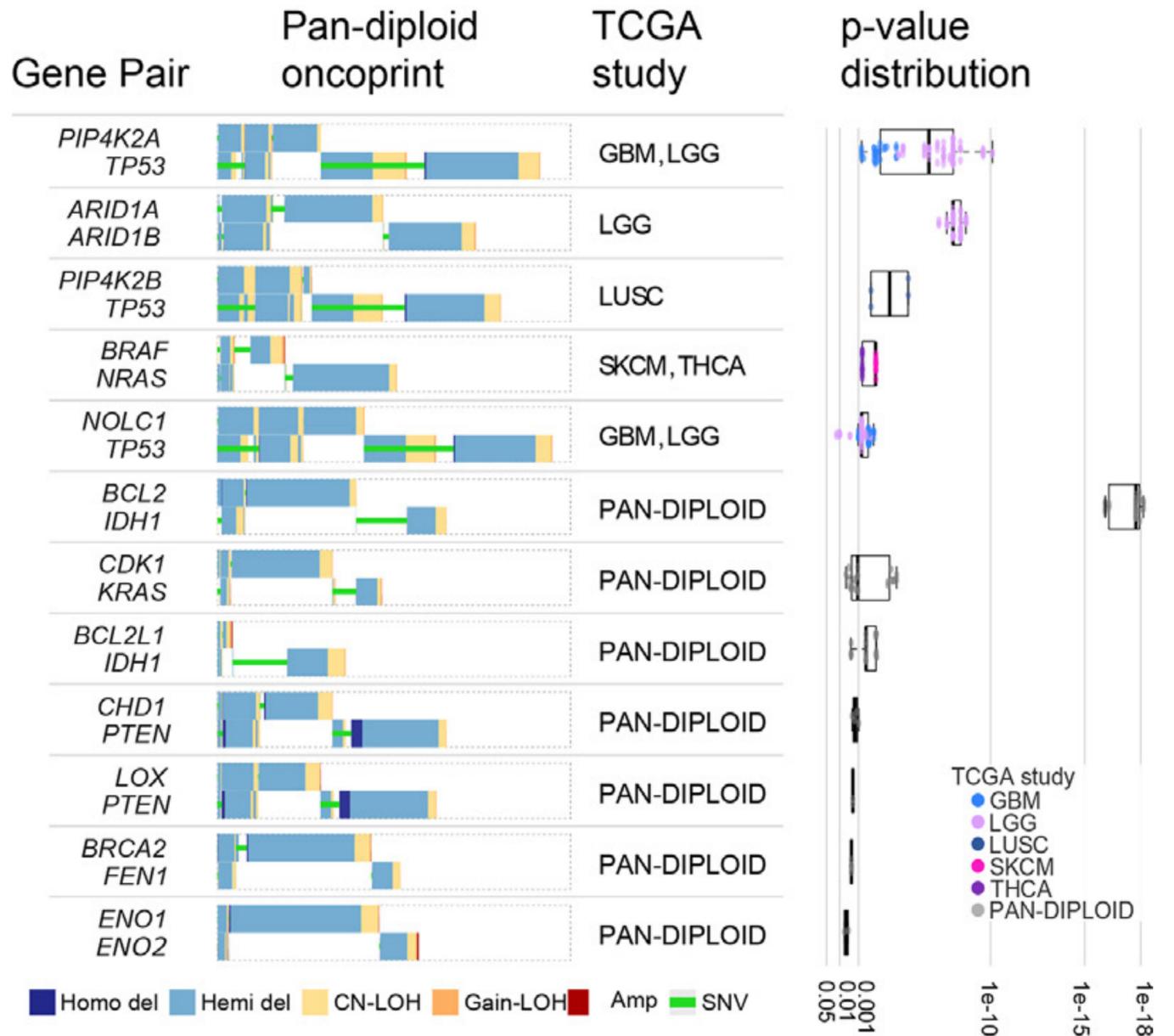
https://jokergoo.github.io/circlize_book/book/index.html

Summarising Entire Cohorts

Based on the ComplexHeatmap package

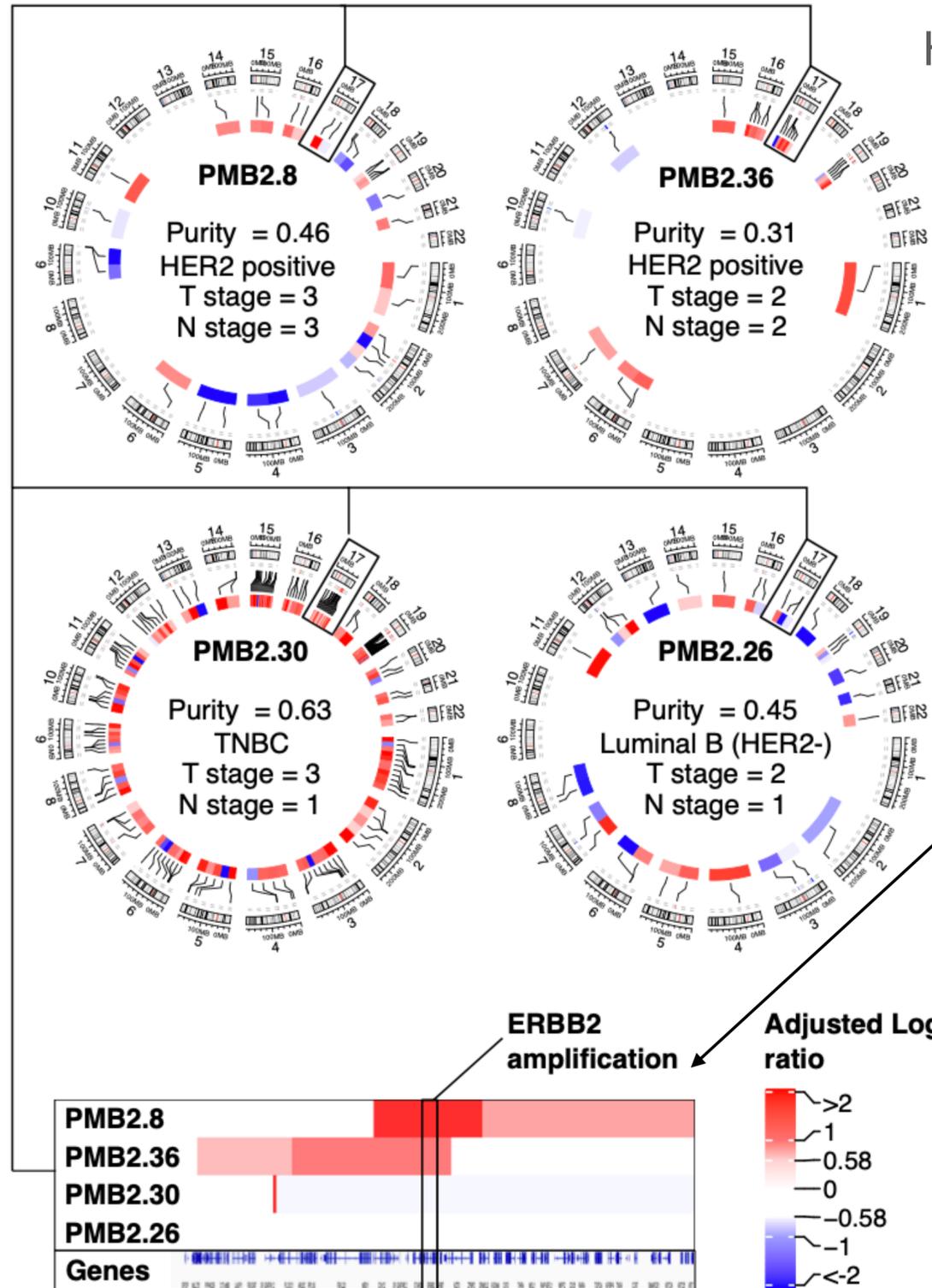


Highlighting mutual exclusivity



“Memo” Sorting, designed by B. Arman Aksoy to highlight mutual exclusivity

From Overview to Detail



Heatmap + Circos visualization

Image from IGV, a Java based alignment visualisation tool



Received: 26 April 2023 | Revised: 19 July 2023 | Accepted: 11 August 2023

DOI: 10.1002/jex2.108

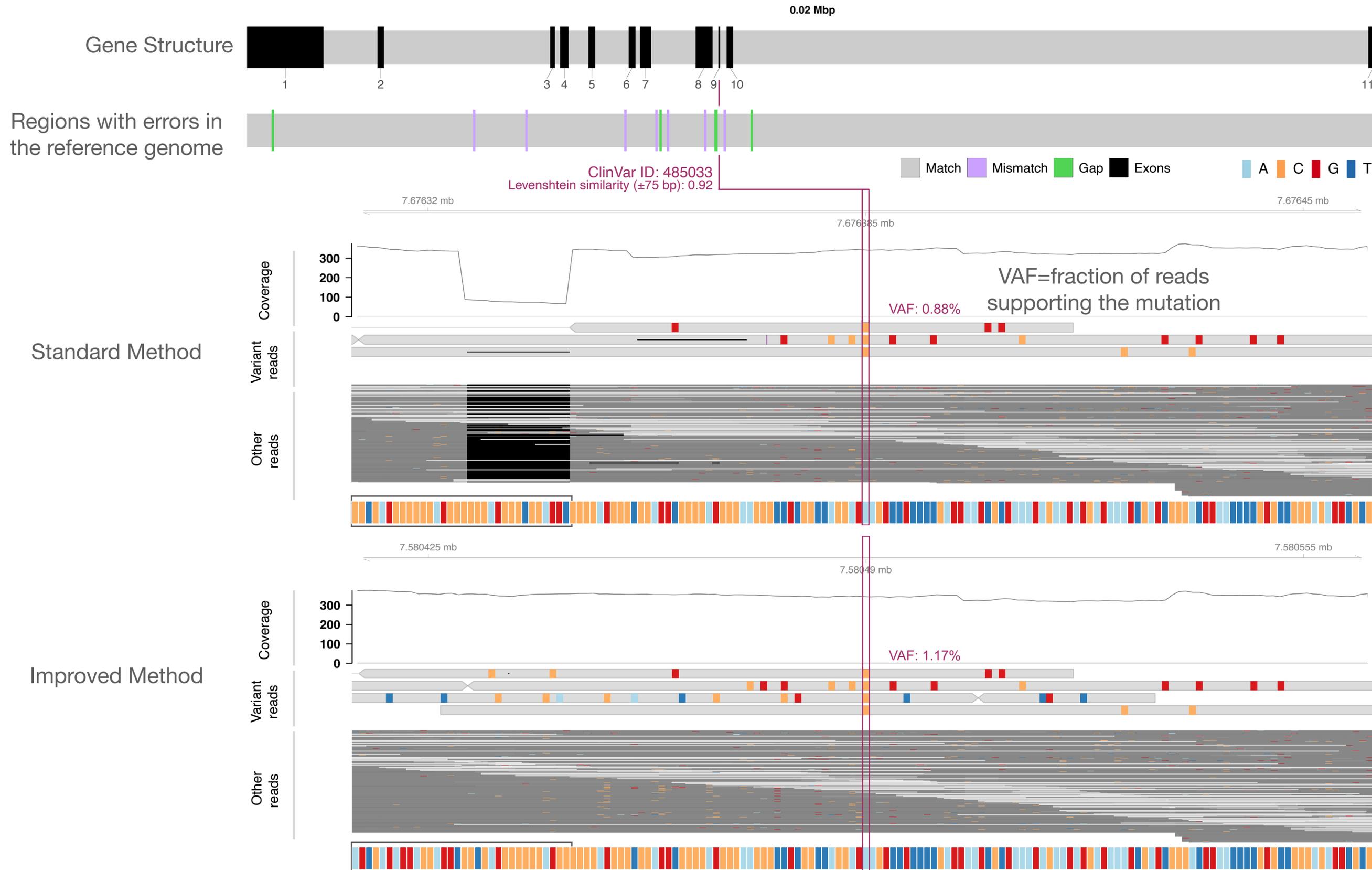


RESEARCH ARTICLE

Integrating extracellular vesicle and circulating cell-free DNA analysis using a single plasma aliquot improves the detection of HER2 positivity in breast cancer patients

Vera Mugoni¹ | Yari Ciani¹ | Orsetta Quaini¹ | Simone Tomasini¹ | Michela Notarangelo¹ | Federico Vannuccini¹ | Alessia Marinelli¹ | Elena Leonardi² | Stefano Pontalti³ | Angela Martinelli¹ | Daniele Rossetto¹ | Isabella Pesce¹ | Sheref S. Mansy¹ | Mattia Barbareschi² | Antonella Ferro³ | Orazio Caffo³ | Gerhardt Attard⁴ | Dolores Di Vizio⁵ | Vito Giuseppe D'Agostino¹ | Caterina Nardella¹ | Francesca Demichelis¹

Focusing on single bases



Based on packages:
ggplot2
ggrepel
Gviz
seqvisr
GenomicAlignments

Designed by Ilaria Cherchi,
PhD student

Unpublished data
do not post

Conclusions

- Bioinformatics have high requirements in terms of visualisation. Starting from huge amounts of data, we need to show broad overview of results but also precise details.
- DNA is a 1D entity of 3.3B points. At the same time it's a dynamic 3D physical object. Each DNA base is interesting on its own but also in relationship with the others.
- R provides access to visualisation packages that are pivotal for our comprehension of biology and for the dissemination of our results.

Visualisation packages used in this presentation:

ggplot2, ComplexHeatmap, Circlize, patchwork, Gviz, ggrepel, seqvisr, GenomicAlignments

ACKNOWLEDGEMENTS

**Francesca Demichelis Lab:
Liquid Biopsies Team**

Prof. Francesca Demichelis
 Francesco Orlando
 Caterina Nardella
 Orsetta Quaini
 Federico Vannuccini
 Alessia Marinelli
 Marta Paoli
Ilaria Cherchi

Previous Members:

Vera Mugoni
 Gian Marco Franceschini
 Thomas Cantore
 Davide Prandi
Tarcisio Fedrizzi
 Giacomo D'Amato

All the other members of the group!

Core Facilities, CIBIO Department, University of Trento
 Next generation Sequencing (NGS)



Many thanks to all patients and families!

