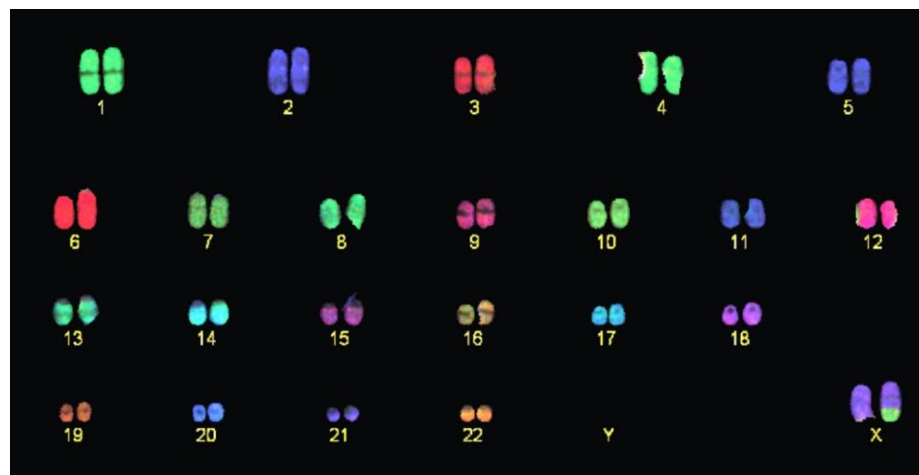# Introduction to
# next-generation sequencing
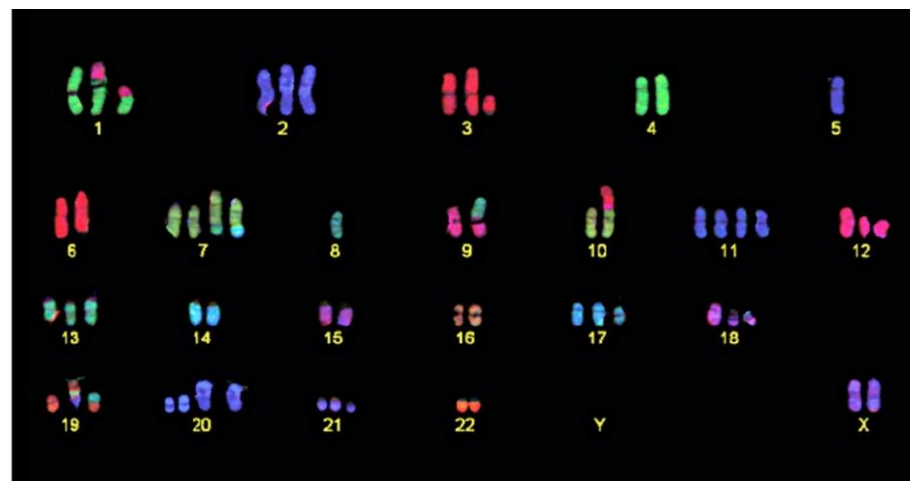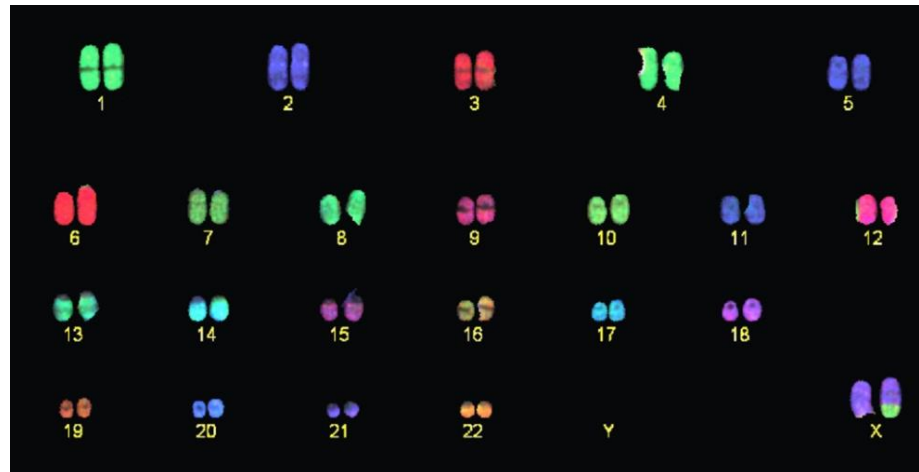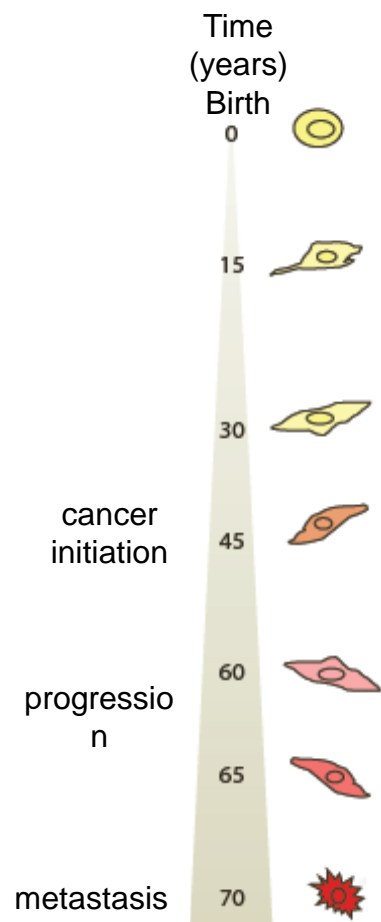
## Pre-IMPAKT 2015

Serena Nik-Zainal

Wellcome-Beit Fellow & WT Intermediate Clinical Research Fellow
Honorary Consultant Clinical Geneticist

wellcome trust
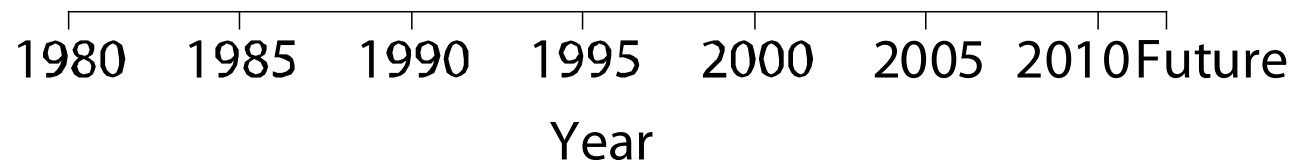**sanger**
institute

**Cambridge University Hospitals** **NHS**
NHS Foundation Trust

**International
Cancer Genome
Consortium**

Time
(years)
Birth
0

Time
(years)

Birth

0

15

30

cancer
initiation

45

60

progressio
n

65

metastasis 70

Part I: What's all the fuss about?

# MASSIVELY PARALLEL SEQUENCING

# Massively-parallel sequencing

1980　1985　1990　1995　2000　2005　2010 Future
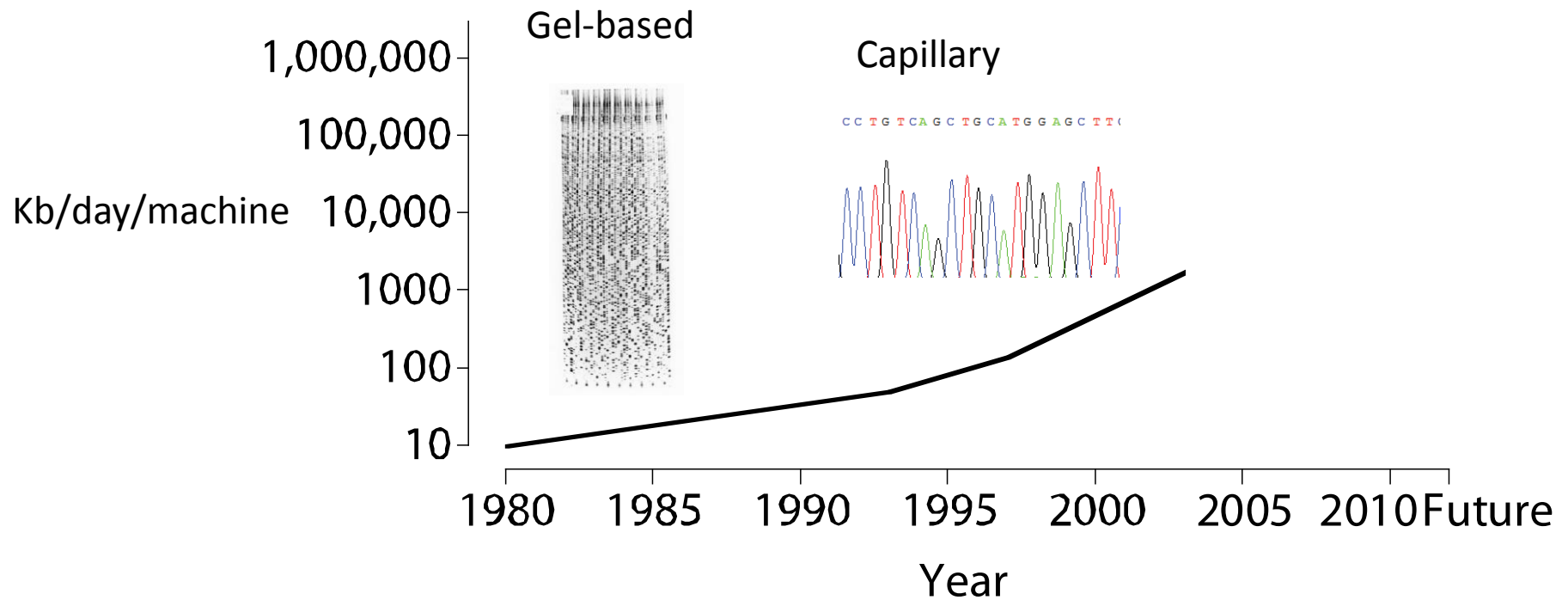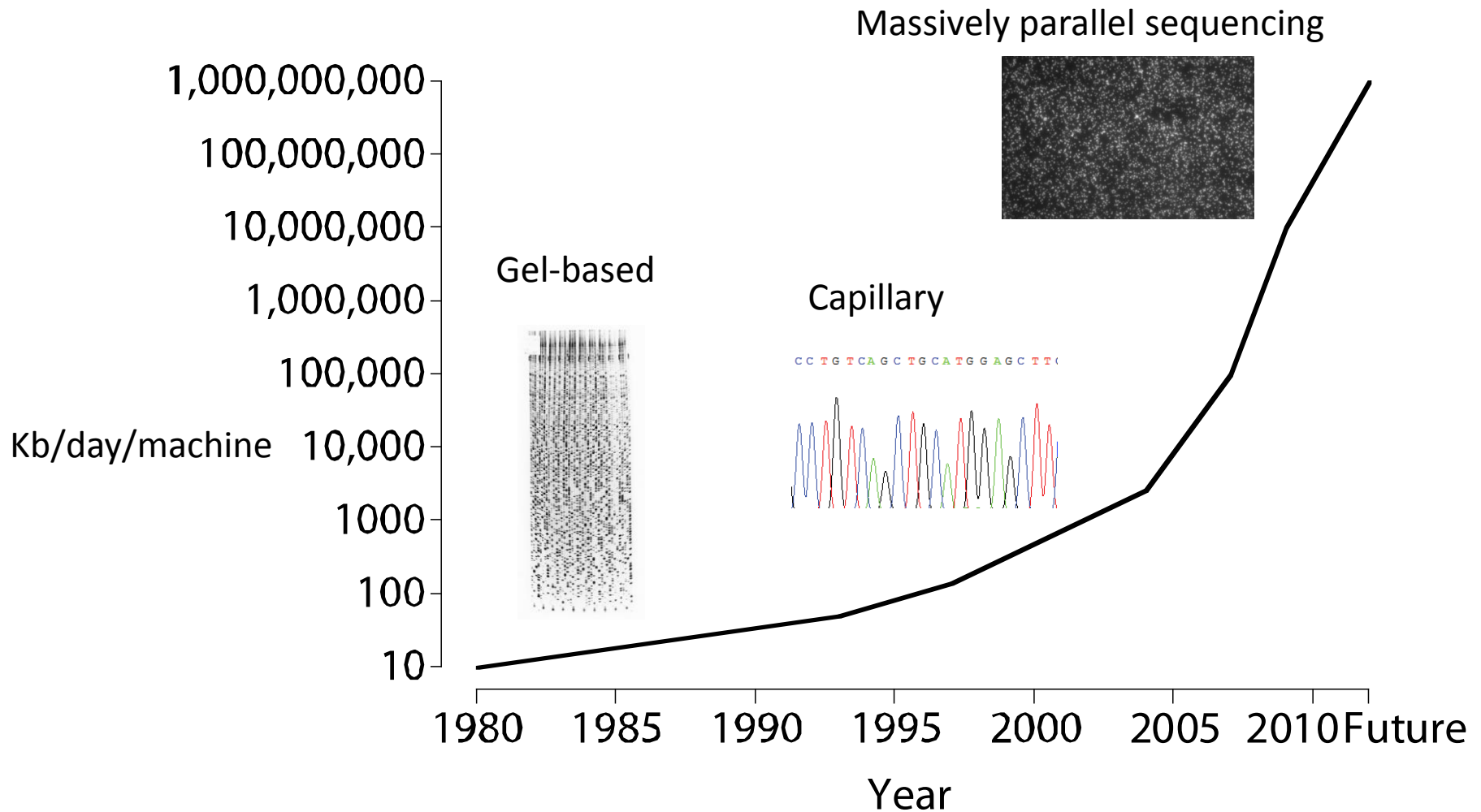
Year

# Massively-parallel sequencing

# Massively-parallel sequencing

# Massively-parallel sequencing


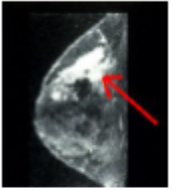
Massively parallel sequencing

Gel-based

Capillary

CC TG TCAGC TGCATG G AGC TT

Kb/day/machine

1,000,000,000

100,000,000

10,000,000

1,000,000

100,000

10,000

1000

100
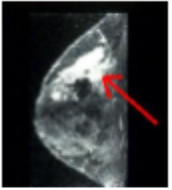
10

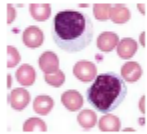1980    1985    1990    1995    2000    2005    2010 Future
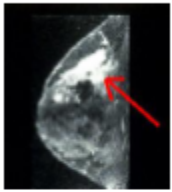
Year

DNA Samples

Tumour

DNA Samples

Tumour

Normal blood

DNA Samples

Tumour

Normal blood

Library preparation

500bp

DNA Samples

Tumour

Normal blood

Library preparation

500bp

Sequenced 100 bps          500 bp          Sequenced 100 bps
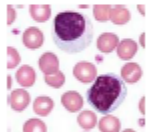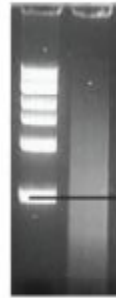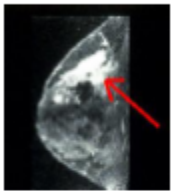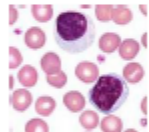
# Paired-end high-coverage next-generation sequencing experiment

## DNA Samples

Tumour

Normal blood

## Library preparation

500bp

## Sequencing
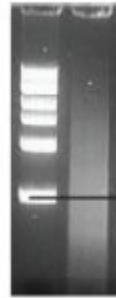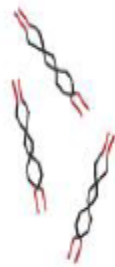
30X coverage

Sequenced 100 bps          500 bp          Sequenced 100 bps

# Massively-parallel sequencing

**sequencing experiment**          **genomic footprint**

- Whole genome sequencing          3,000,000,000 base pairs

# Massively-parallel sequencing

**sequencing experiment**           **genomic footprint**

- Whole genome sequencing     3,000,000,000 base pairs

- Exome sequencing                     50,000,000 base pairs

# Massively-parallel sequencing

| sequencing experiment | genomic footprint |
|---|---|
| • Whole genome sequencing | 3,000,000,000 base pairs |
| • Exome sequencing | 50,000,000 base pairs |
| • Targeted gene screens | 10,000,000 base pairs |

# Massively-parallel sequencing

| sequencing experiment | genomic footprint |
|---|---|
| • Whole genome sequencing | 3,000,000,000 base pairs |
| • Exome sequencing | 50,000,000 base pairs |
| • Targeted gene screens | 10,000,000 base pairs |

# Basic principle for calling somatic mutations

Millions of short reads

# Basic principle for calling somatic mutations



Millions of short reads

**Reference genome**

Alignment to reference

# Basic principle for calling somatic mutations



Millions of short reads

**Reference genome**

Alignment to reference

Tumour

Call all variants

# Basic principle for calling somatic mutations



Millions of short reads

**Reference genome**

Alignment to reference

Tumour

Normal

Call all variants

# Basic principle for calling somatic mutations



Millions of short reads

**Reference genome**

Alignment to reference

Tumour

Normal

Call all variants

Somatic variants

# Bioinformatics

- Data processing

Sequencing

# Bioinformatics

- Data processing

Sequencing

↓

Alignment

# Bioinformatics

- Data processing

Sequencing

↓

Alignment

↓

Calling mutations

Substitutions    Indels    Rearrangements    Copy number aberrations

# Genomic abnormalities

exon                    intron                    UTR

# Genomic abnormalities



exon          intron          UTR

TTCACG

# Genomic abnormalities



exon     intron     UTR

TTCACG

TTTACG

substitution

# Genomic abnormalities



exon          intron          UTR

TTCACG

TTTACG          TTT-CG

substitution    deletion/
                insertion

# Genomic abnormalities



exon          intron          UTR

TTCACG

TT**T**ACG          TTT-CG

substitution          deletion/
                      insertion

base pair resolution

# Genomic abnormalities



exon    intron    UTR

duplication

TTCACG

TTTACG    TTT-CG

substitution    deletion/
insertion

base pair resolution

# Genomic abnormalities



exon · intron · UTR
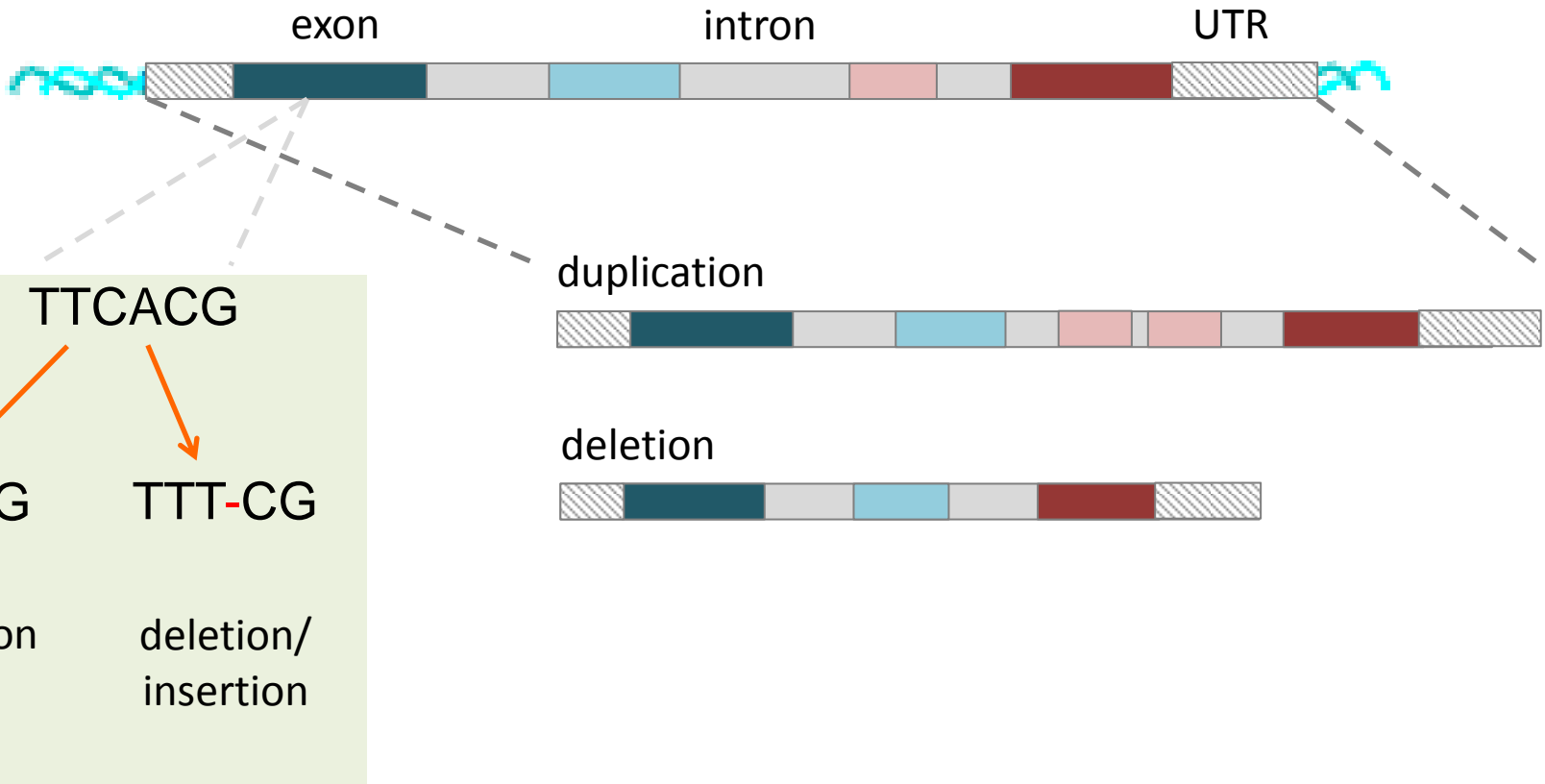
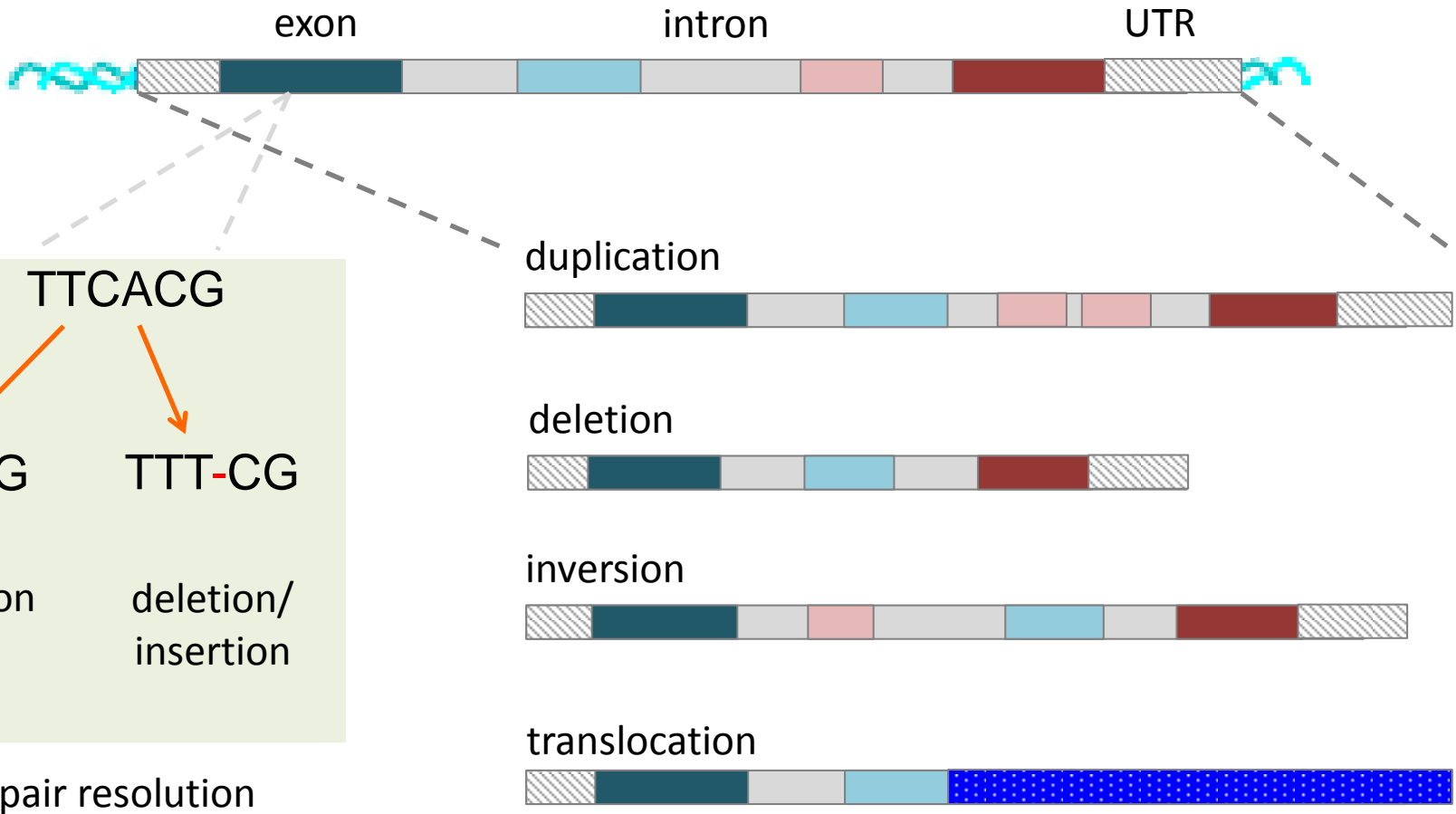duplication

deletion

TTCACG

TTTACG → substitution

TTT-CG → deletion/insertion

base pair resolution
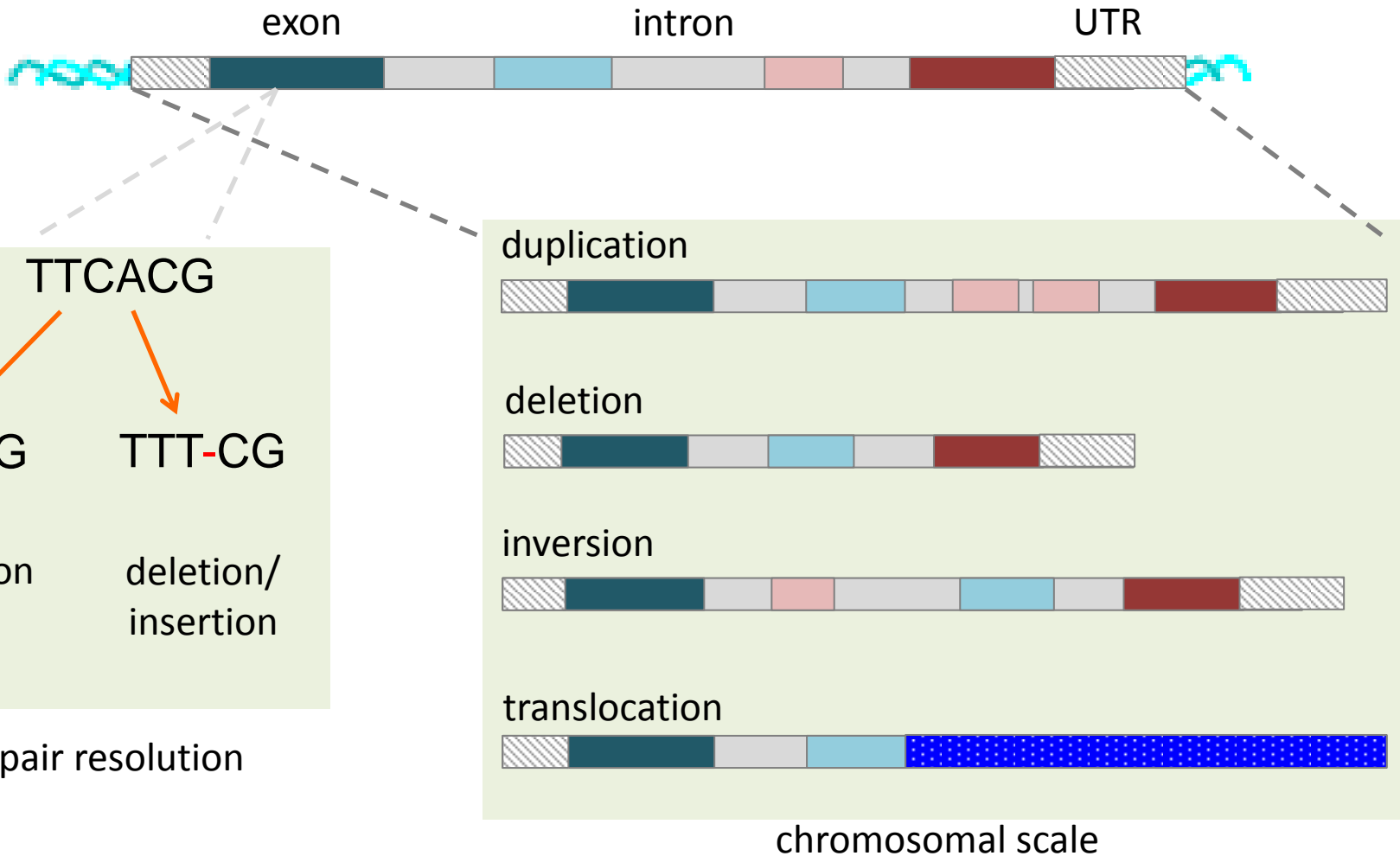
# Genomic abnormalities

# Genomic abnormalities



exon          intron          UTR

TTCACG

TTTACG          TTT-CG

substitution          deletion/
                      insertion

base pair resolution

duplication

deletion

inversion

translocation

# Bioinformatics

- ## Data processing

Sequencing

↓

Alignment

↓

Calling mutations

Substitutions    Indels    Rearrangements    Copy number aberrations
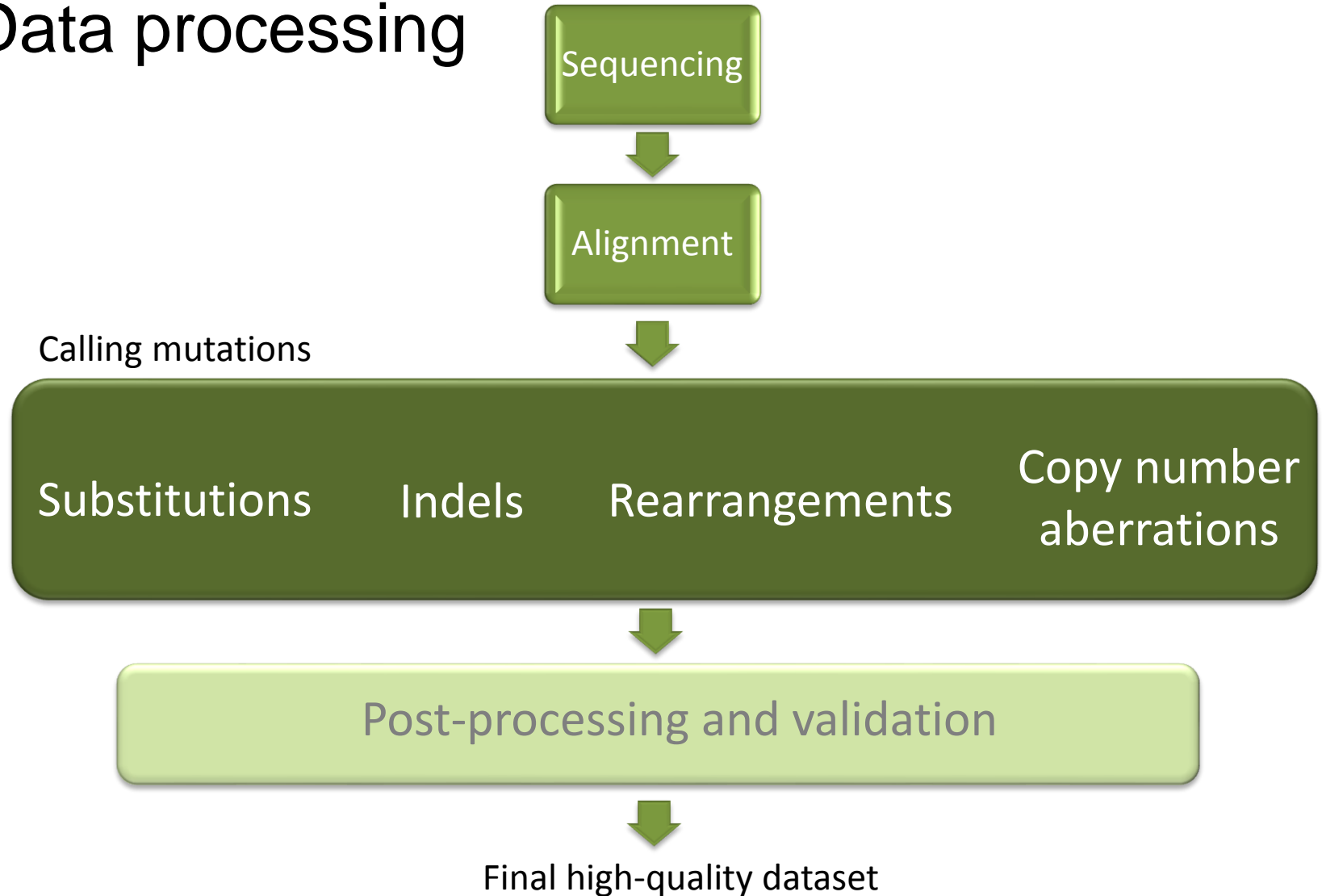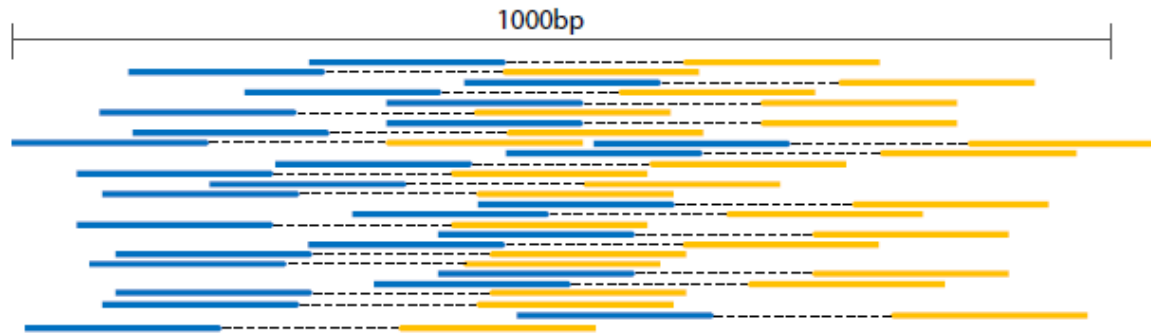
# Limitations

- DNA
  - Quality of DNA (Fresh frozen, FFPE)
  - Ploidy of DNA
  - Normal cell "contamination"

- Sequencing
  - Variation in coverage
  - Systematic sequencing artefacts

- Reference genome
  - Poorly-defined parts of the genome
  - Repeats

- Mutation-calling
  - Sensitivity
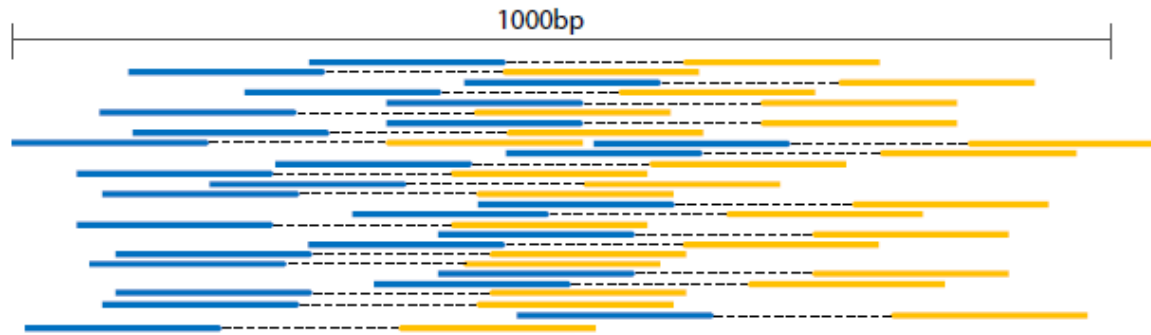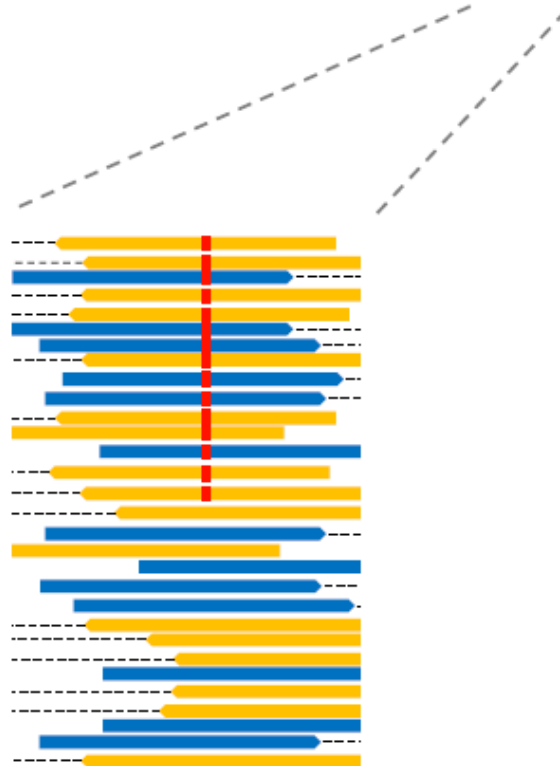  - Specificity

# Bioinformatics

- ## Data processing

Align billions of NGS sequenced read pairs back to the reference genome
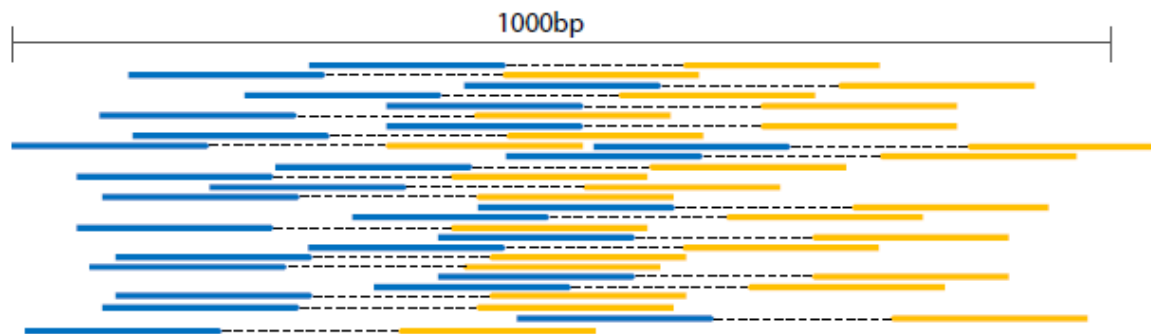
1000bp

Align billions of
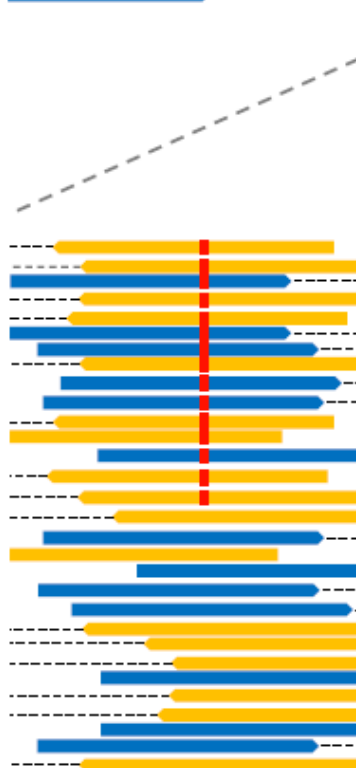NGS sequenced
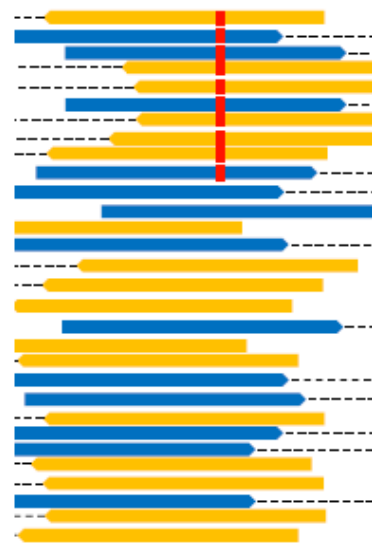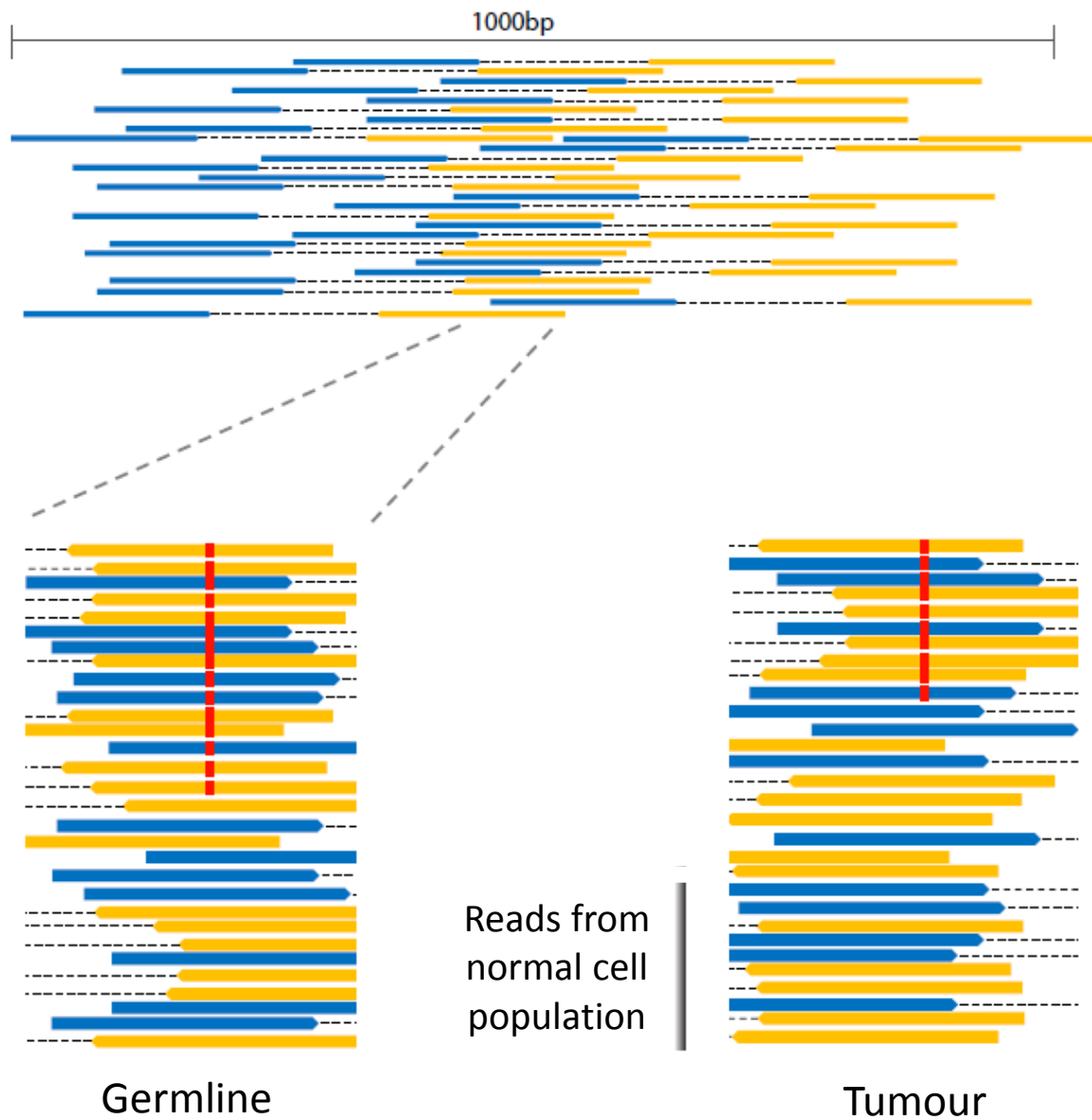read pairs back
to the reference
genome

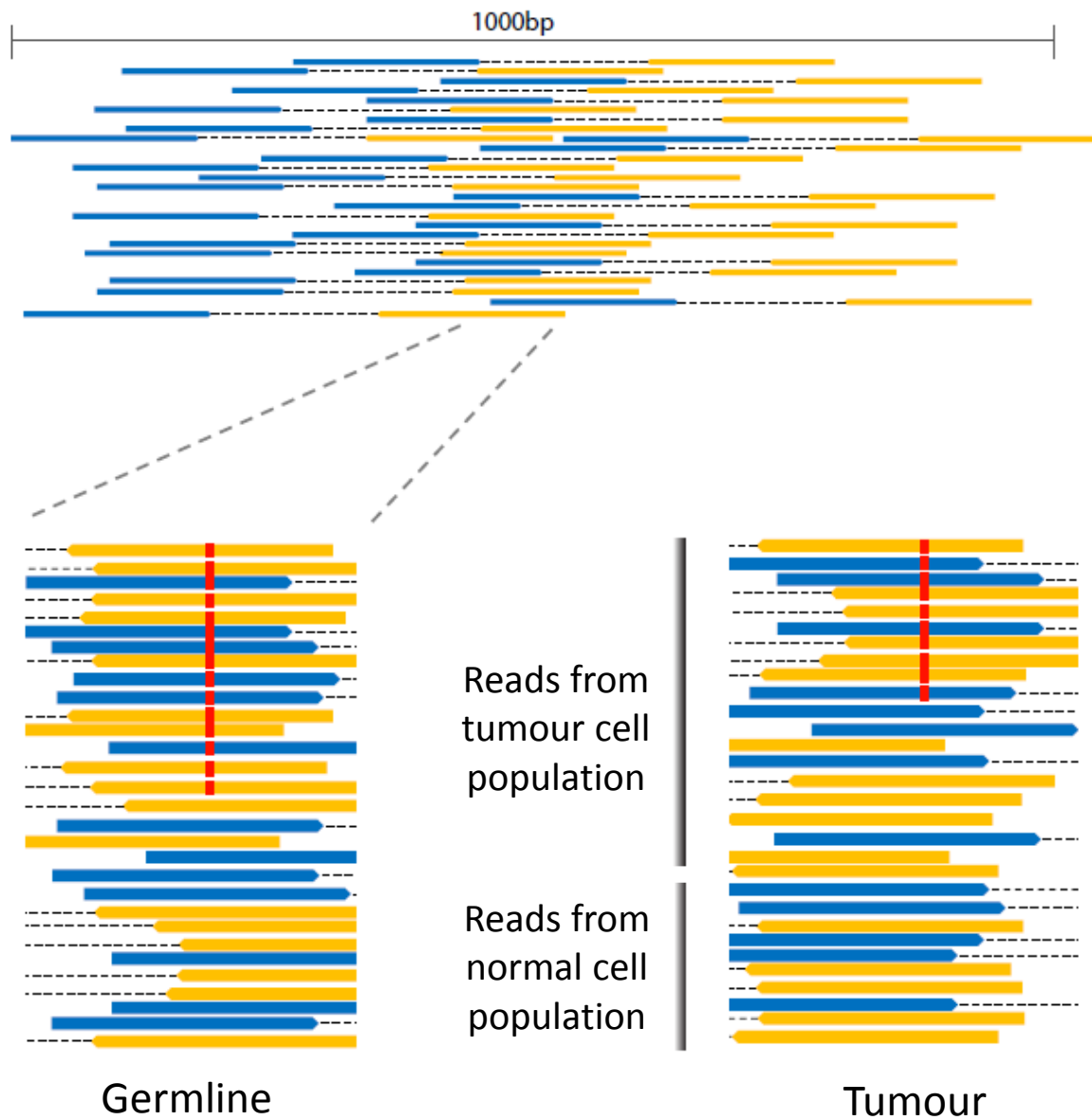Fraction
of reads
50%

Germline

1000bp

Fraction
of reads
50%

Germline

Tumour

1000bp

Fraction
of reads
50%

Reads from
normal cell
population

Germline

Tumour

1000bp

Fraction
of reads
50%

Reads from
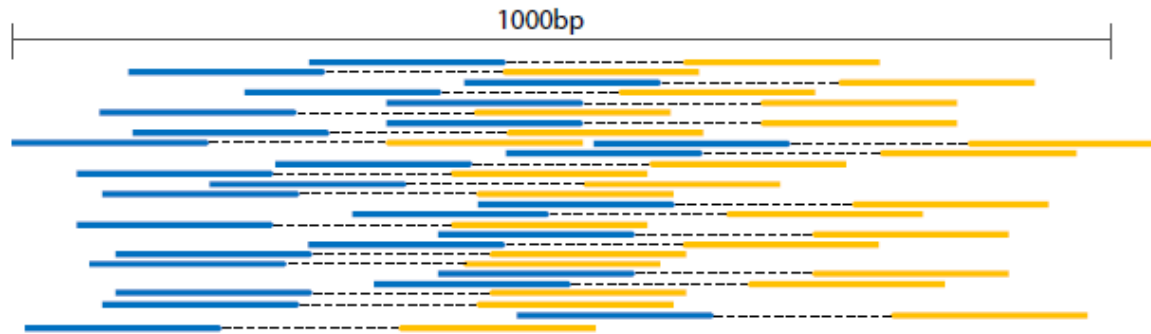tumour cell
population

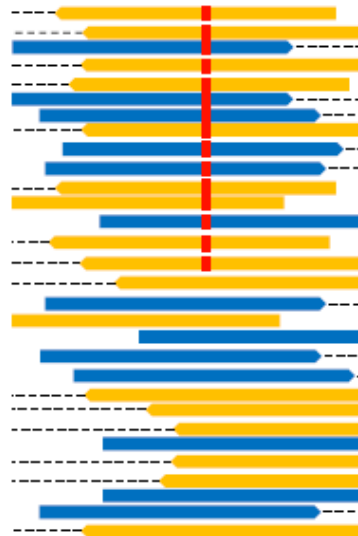Reads from
normal cell
population

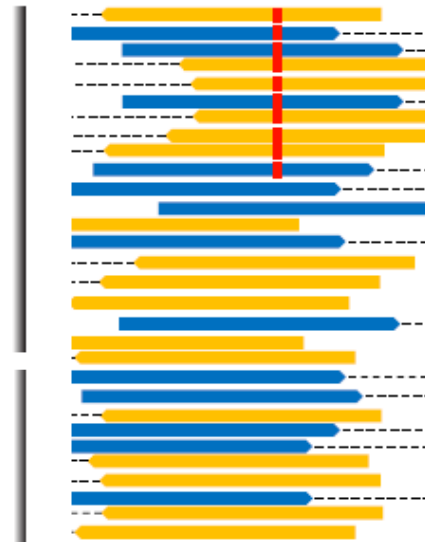Germline

Tumour

1000bp

Fraction of reads 50%

Reads from tumour cell population

Reads from normal cell population

Fraction of reads 35%

Germline

Tumour

# Summary I

- Advances in sequencing chemistry has led to a vast increase in scale and speed of sequencing, permitting unprecedented access to all parts of the human genome

- This technology is digital, providing quantifiable information for every mutation seen

- A huge amount of compute is required for processing and for storage of raw data

- A considerable amount of computational expertise is required to ensure high quality datasets with high sensitivity and high specificity

Part II: Making the most out of NGS data

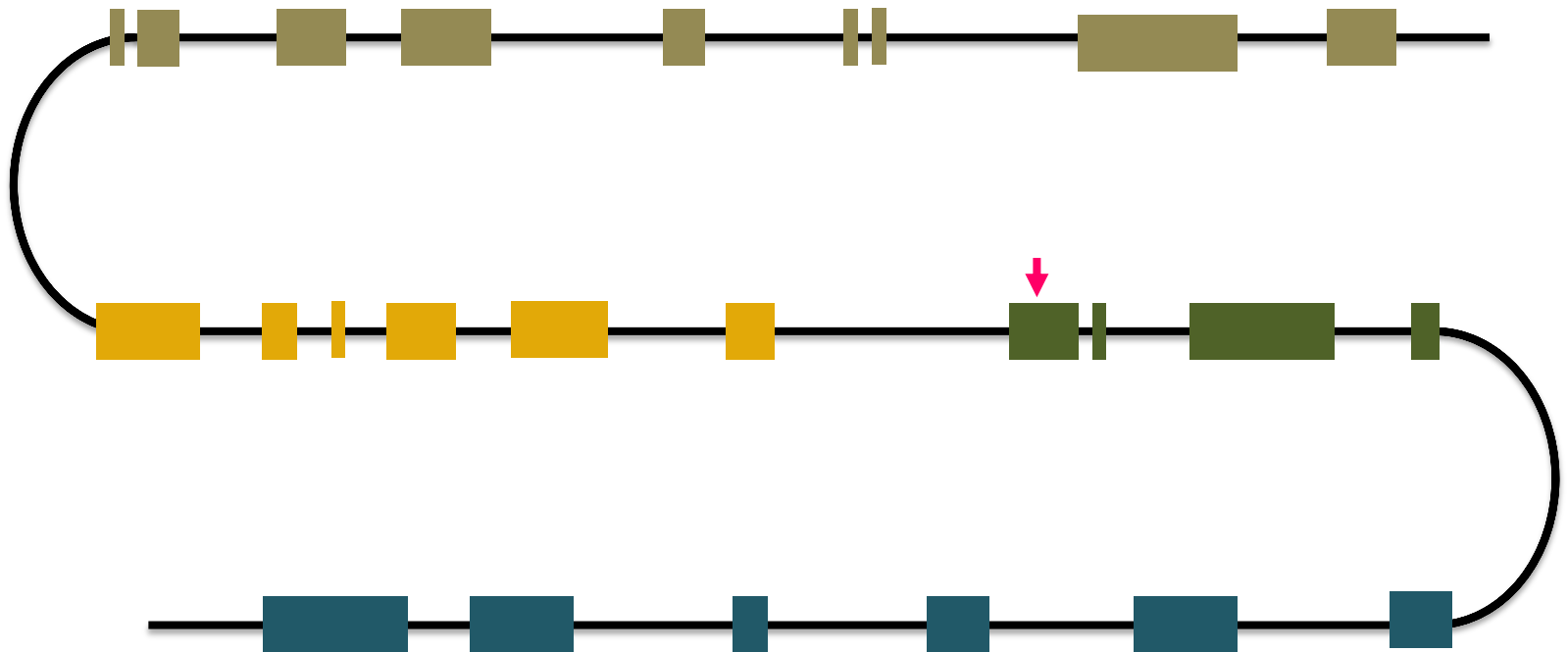# DATA ANALYSIS
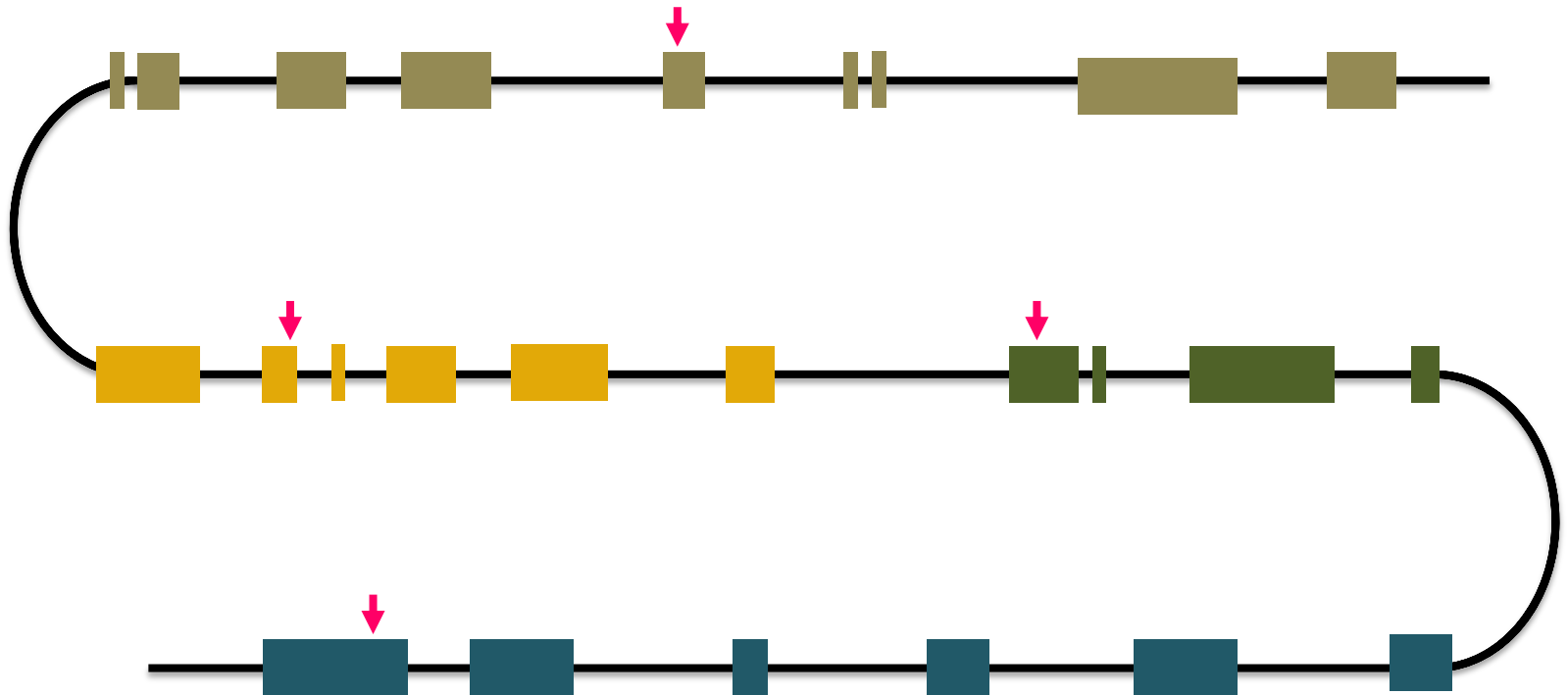
# Bioinformatics

- Data processing

# Bioinformatics

- Data processing

- Downstream analysis

  - Cancer genes
  - Somatic mutation signatures
  - Cancer evolution

# DRIVERS
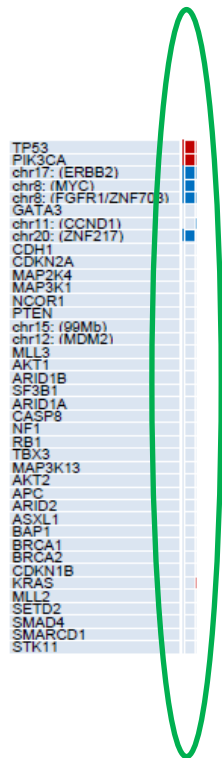
# DRIVERS

# Downstream analysis: Cancer genes I

- ## Genomic scenario

  - *ERBB2* Amplification
    - (Breast Cancer)

  - *BCR-ABL*
    - (CML)

  - *EGFR*
    - *(NSCLC)*

  - *EML4-ALK*
    - *(NSCLC)*

  - *KRAS*-negative
    - (colorectal cancer)

  - *BRAF(V600E)*
    - *(Metastatic Melanoma)*

- ## Targeted drug

  - Herceptin & Lapatinib

  - Imatinib (and others)

  - Erlotinib, Gefitinib
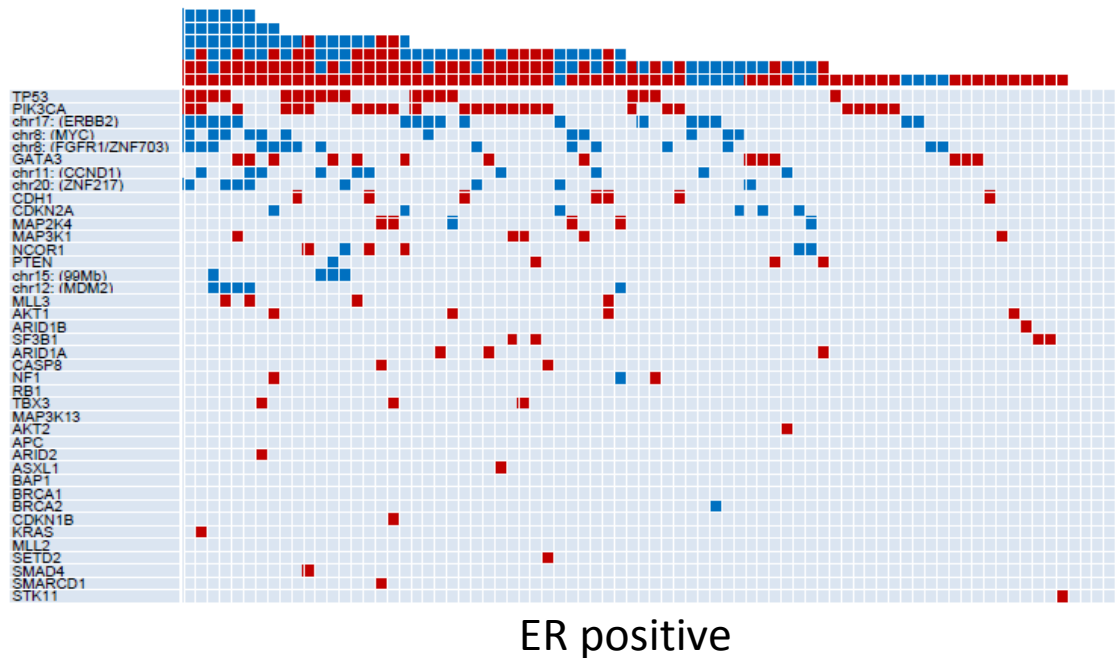
  - Crizotinib

  - Cetuximab

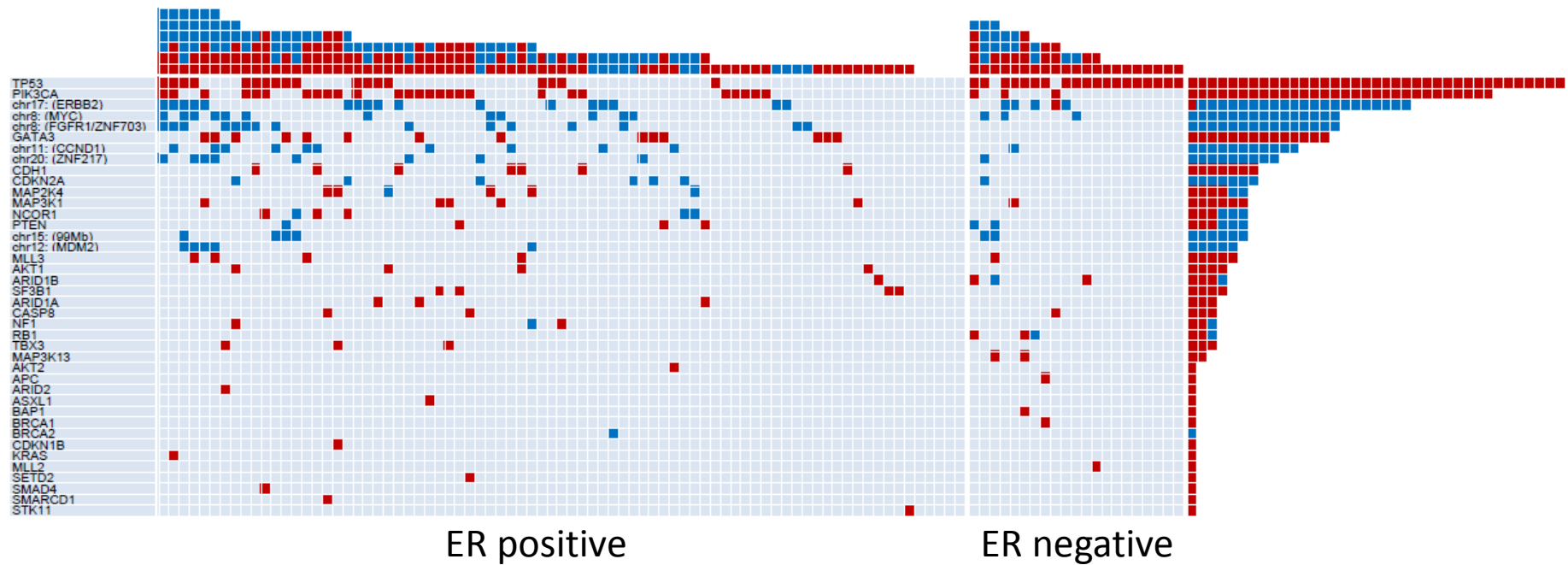  - Vemurafenib

# Downstream analysis: Cancer genes I



TP53
PIK3CA
chr17: (ERBB2)
chr8: (MYC)
chr8: (FGFR1/ZNF703)
GATA3
chr11: (CCND1)
chr20: (ZNF217)
CDH1
CDKN2A
MAP2K4
MAP3K1
NCOR1
PTEN
chr15: (99Mb)
chr12: (MDM2)
MLL3
AKT1
ARID1B
SF3B1
ARID1A
CASP8
NF1
RB1
TBX3
MAP3K13
AKT2
APC
ARID2
ASXL1
BAP1
BRCA1
BRCA2
CDKN1B
KRAS
MLL2
SETD2
SMAD4
SMARCD1
STK11

■ Copy number changes

■ Point mutations

# Downstream analysis: Cancer genes I



ER positive

■ Copy number changes

■ Point mutations

# Downstream analysis: Cancer genes I



ER positive

ER negative

Copy number changes

Point mutations

Stephens et al 2012

The NEW ENGLAND JOURNAL of MEDICINE

ORIGINAL ARTICLE

## Somatic SF3B1 Mutation in Myelodysplasia with Ring Sideroblasts

E. Papaemmanuil, M. Cazzola, J. Boultwood, L. Malcovati, P. Vyas, D. Bowen,
A. Pellagatti, J.S. Wainscoat, E. Hellstrom-Lindberg, C. Gambacorti-Passerini,
A.L. Godfrey, I. Rapado, A. Cvejic, R. Rance, C. McGee, P. Ellis, L.J. Mudie,
P.J. Stephens, S. McLaren, C.E. Massie, P.S. Tarpey, I. Varela, S. Nik-Zainal,
H.R. Davies, A. Shlien, D. Jones, K. Raine, J. Hinton, A.P. Butler, J.W. Teague,
E.J. Baxter, J. Score, A. Galli, M.G. Della Porta, E. Travaglino, M. Groves, S. Tauro,
N.C. Munshi, K.C. Anderson, A. El-Naggar, A. Fischer, V. Mustonen, A.J. Warren,
N.C.P. Cross, A.R. Green, P.A. Futreal, M.R. Stratton, and P.J. Campbell
for the Chronic Myeloid Disorders Working Group of the International
Cancer Genome Consortium

**nature genetics**

## Frequent mutation of the major cartilage collagen gene COL2A1 in chondrosarcoma

Patrick S Tarpey[1,8], Sam Behjati[1,2,8], Susanna L Cooke[1], Peter Van Loo[1,3], David C Wedge[1], Nischalan Pillay[4,5],
John Marshall[1], Sarah O'Meara[1], Helen Davies[1], Serena Nik-Zainal[1], David Beare[1], Adam Butler[1], John Gamble[1],
Claire Hardy[1], Jonathon Hinton[1], Ming Ming Jia[1], Alagu Jayakumar[1], David Jones[1], Calli Latimer[1],
Mark Maddison[1], Sancha Martin[1], Stuart McLaren[1], Andrew Menzies[1], Laura Mudie[1], Keiran Raine[1],
Jon W Teague[1], Jose M C Tubio[1], Dina Halai[4], Roberto Tirabosco[4], Fernanda Amary[4], Peter J Campbell[1,6,7],
Michael R Stratton[1], Adrienne M Flanagan[4,5] & P Andrew Futreal[1]

# DRIVERS

DRIVERS      AND      PASSENGERS

# Downstream analysis II: Using passengers

BRCA1 null        5

BRCA2 null        4

ER+, HER2-        5

ER+, HER2+        2

ER-, HER2+        2

ER-, HER2-        3

_____

Total        21 whole-genome sequenced breast cancers

Nik-Zainal et al, Cell, 2012a

# Downstream analysis II: Using passengers

BRCA1 null          5

BRCA2 null          4                          Somatic substitutions          183,916

ER+, HER2-          5

ER+, HER2+          2                          Somatic indels          2,869

ER-, HER2+          2

ER-, HER2-          3                          Somatic rearrangements          1,192

_____

Total          21 whole-genome sequenced breast cancers

Nik-Zainal et al, Cell, 2012a

time (years)

0

15

30

45

60

65

70

time (years)

0

15

30

45

60

65

70

time (years)

0

15

30

45

60

65

70

time (years)

multiple mutational processes added together

time (years)

sequenced cancer genome

# Mutation signatures in human cells

NMF

1

2

3

Lee & Seung 1999, Brunet et al 2004

Ludmil B. Alexandrov

# Downstream analysis II: Mutation signatures
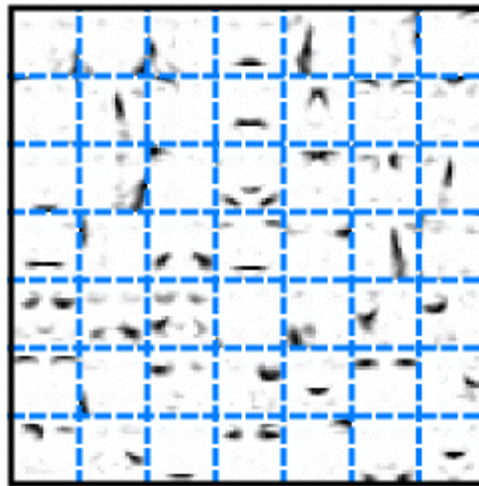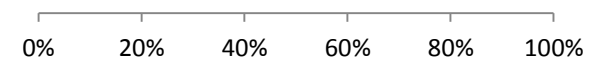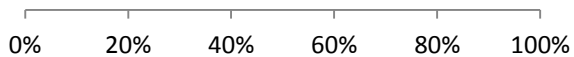
NMF

Lee & Seung 1999, Brunet et al 2004

1

2

3

# Downstream analysis II: Mutation signatures



NMF

Lee & Seung 1999, Brunet et al 2004

1

2

3

# Downstream analysis II: Mutation signatures



NMF

Lee & Seung 1999, Brunet et al 2004

| nose | eyes | mouth |

0%   20%   40%   60%   80%   100%

# Downstream analysis II: Mutation signatures



NMF

Lee & Seung 1999, Brunet et al 2004

# Mutation signatures in human cancers
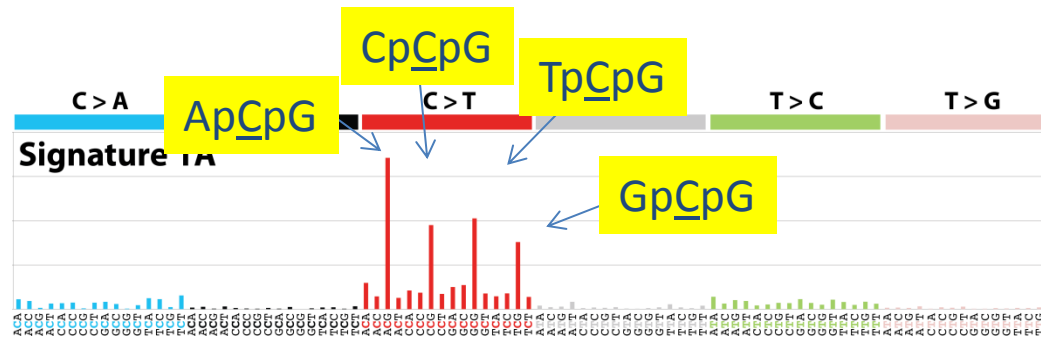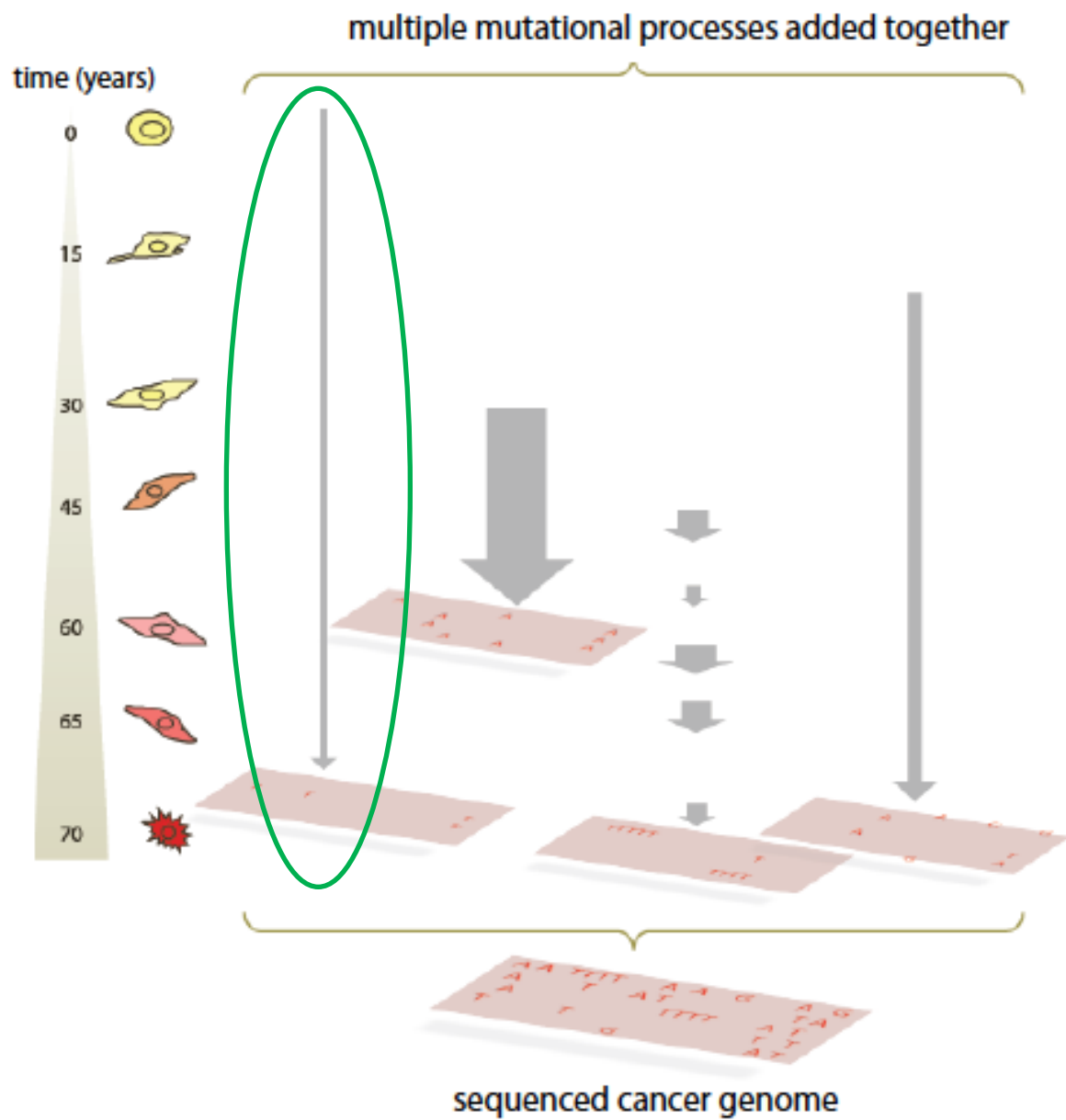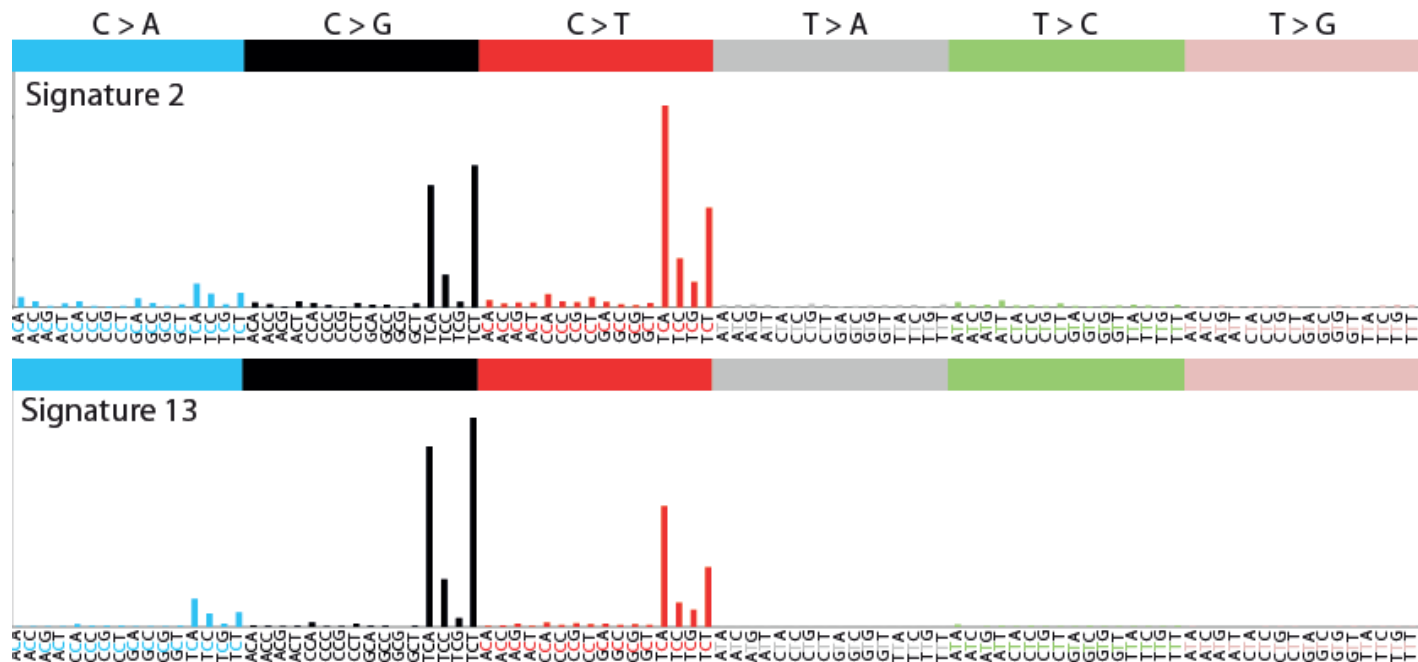
# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types
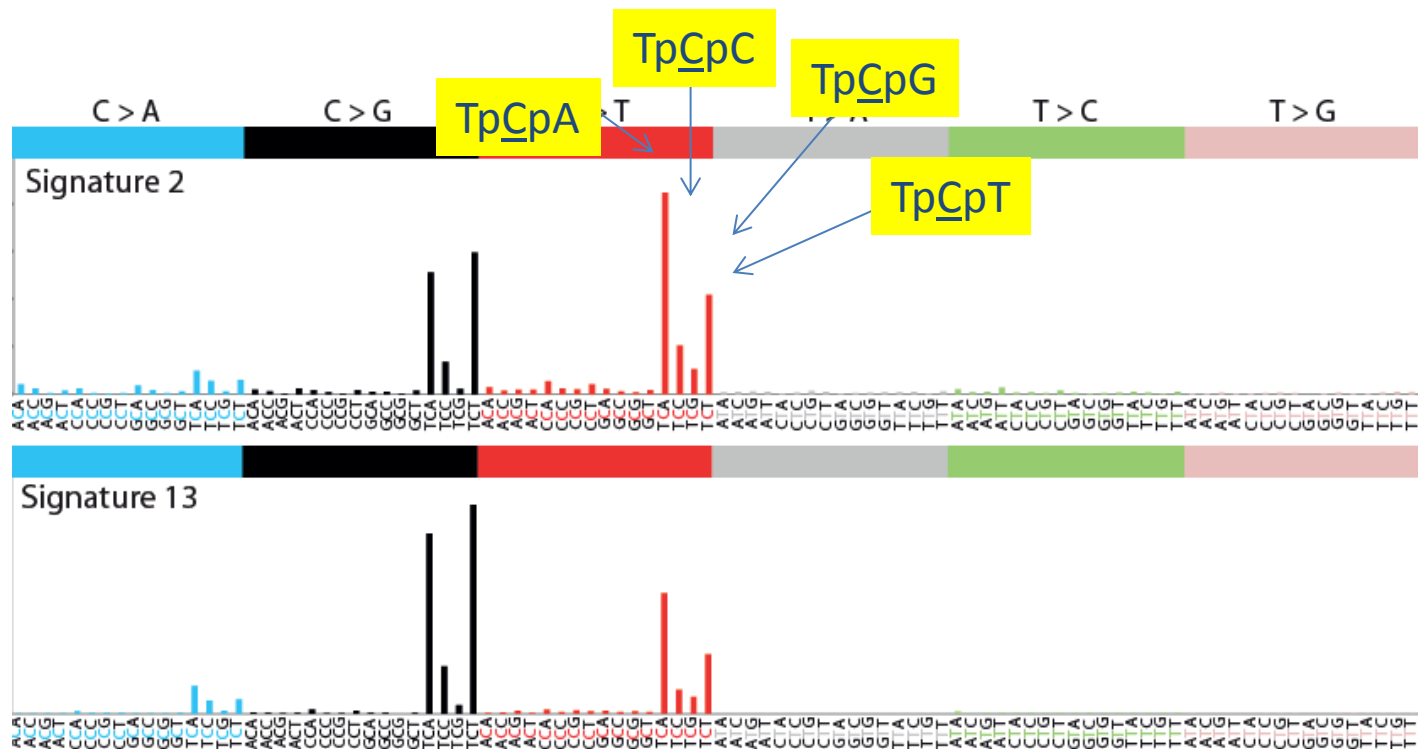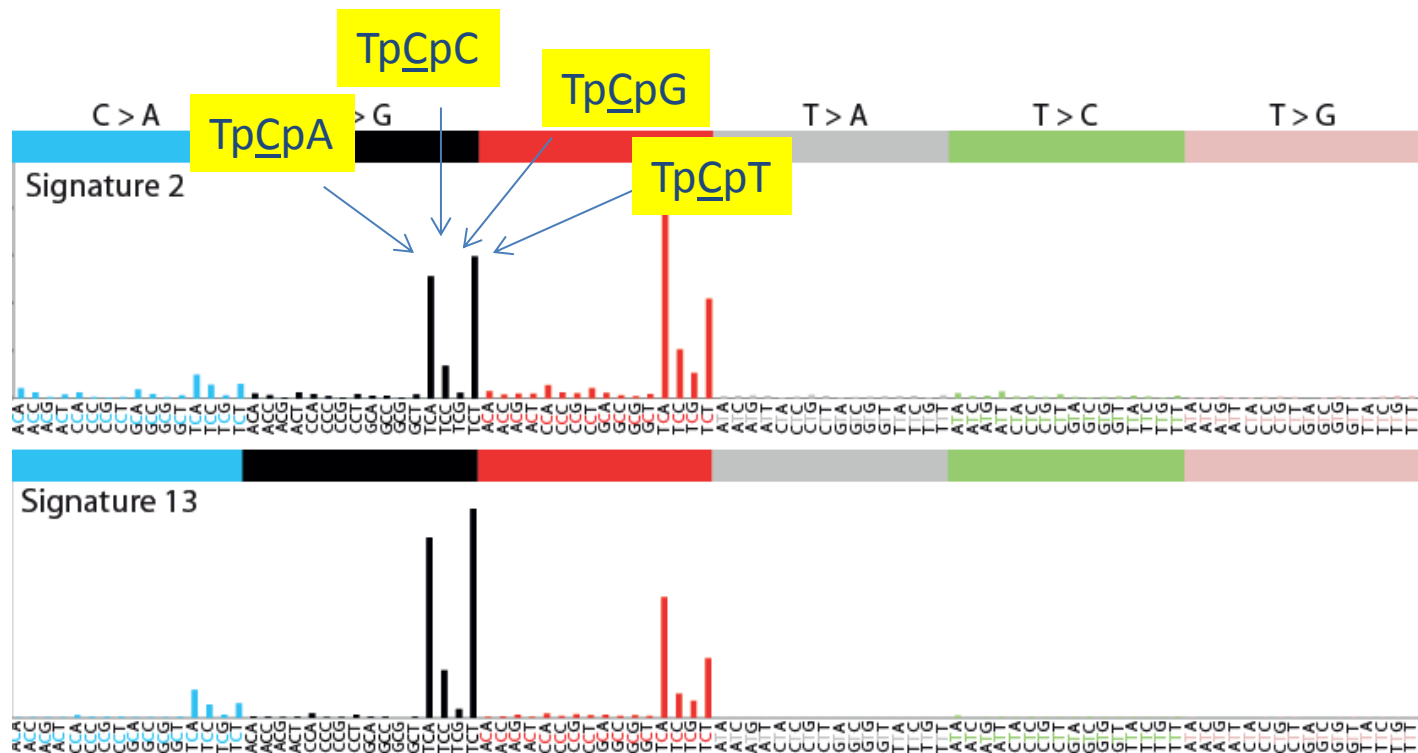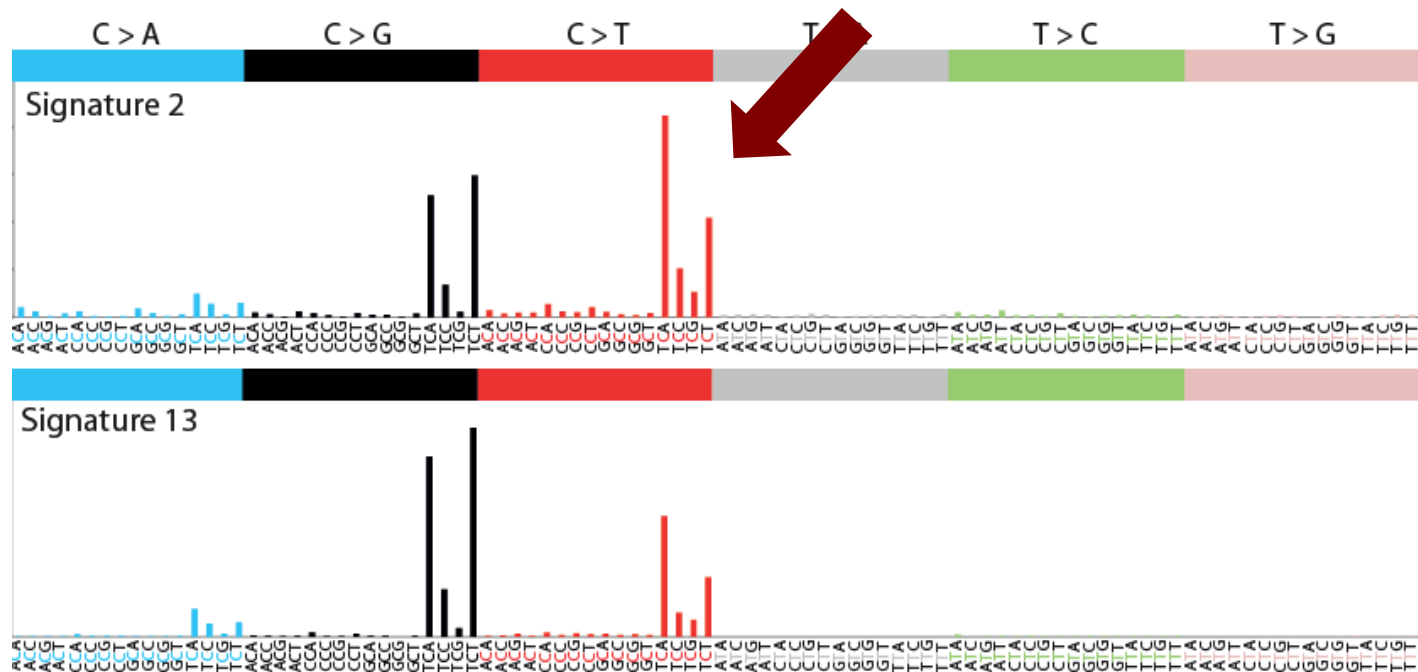
# Mutation signatures present in many cancer types

multiple mutational processes added together

time (years)

sequenced cancer genome

# Mutation signatures present in many cancer types
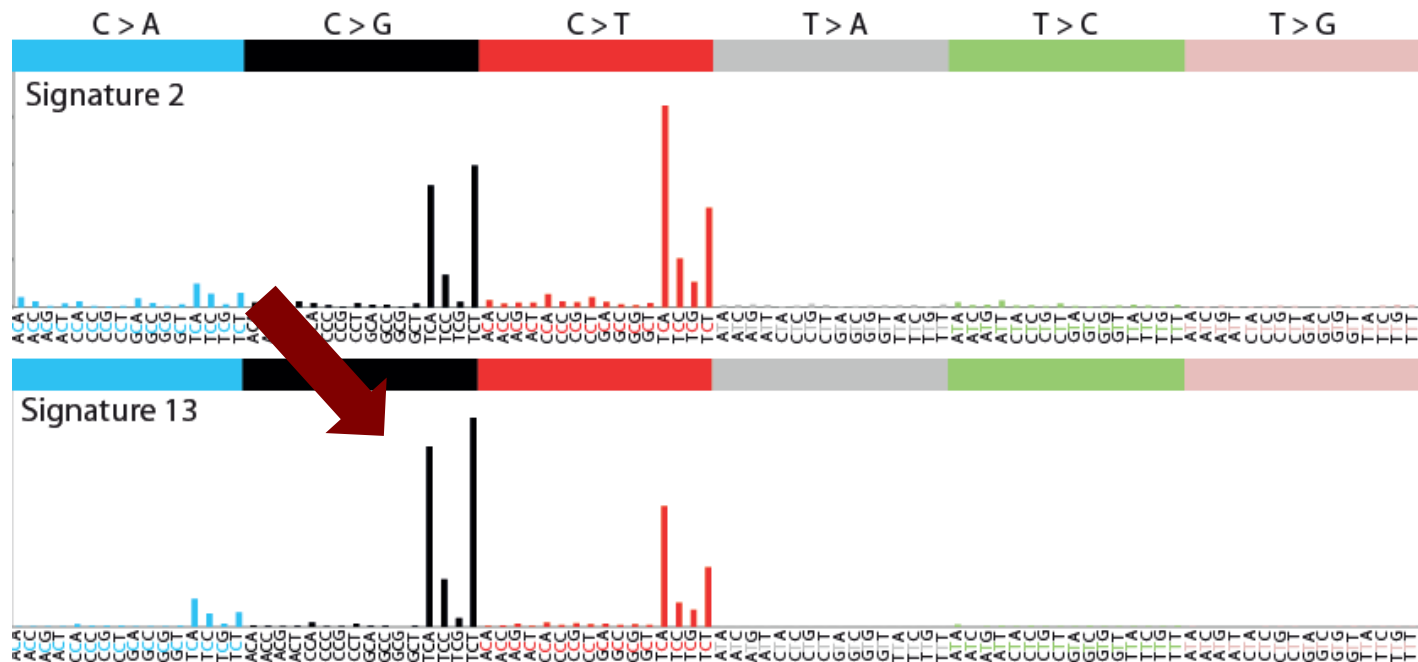


deamination of methylated cytosines

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# Mutation signatures present in many cancer types

# What is the biological explanation for the mutagenic process underlying Signature 2/13?

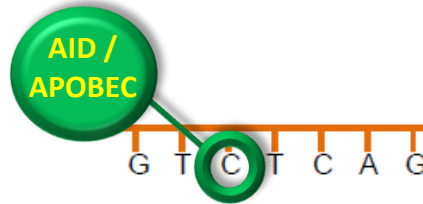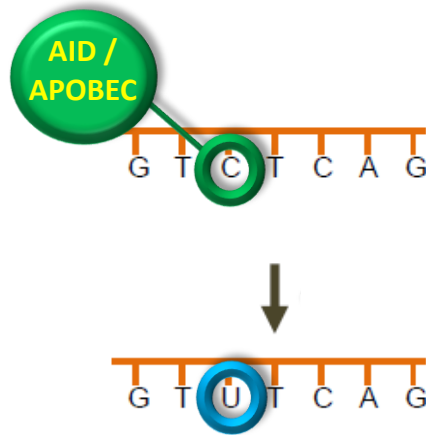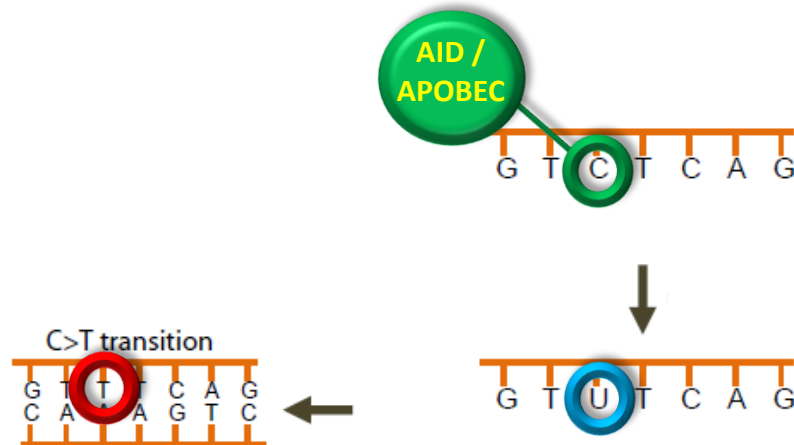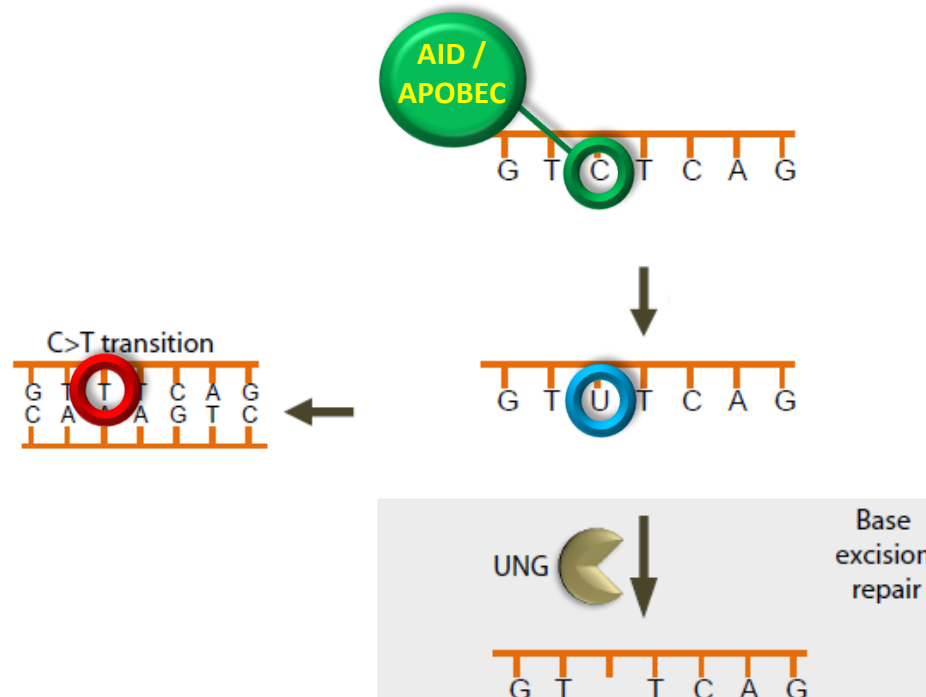# What is the biological explanation for the mutagenic process underlying Signature 2/13?

- Deamination of cytosine by one of the family of AID/APOBEC enzymes?

- The family includes
  AID
  APOBEC1
  APOBEC2
  APOBEC3A-H
  APOBEC4

# What is the biological explanation for the mutagenic process underlying Signature 2/13?
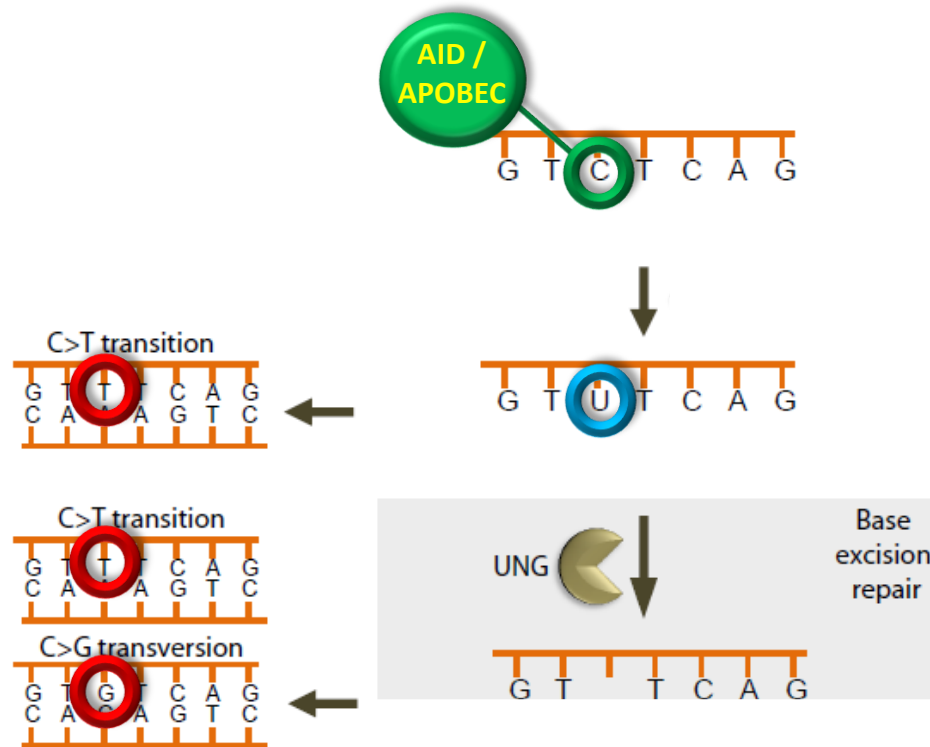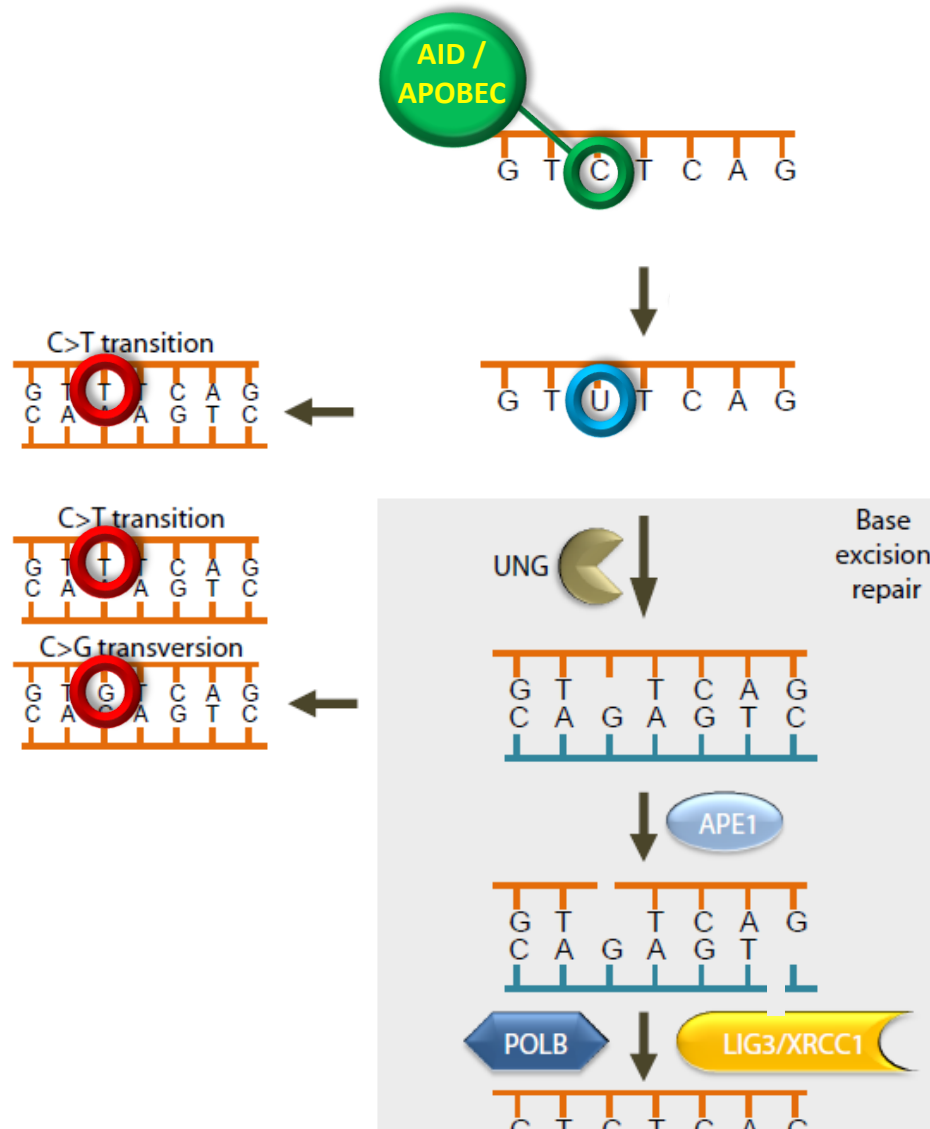
G T C T A G

# What is the biological explanation for the mutagenic process underlying Signature 2/13?

# What is the biological explanation for the mutagenic process underlying Signature 2/13?

# What is the biological explanation for the mutagenic process underlying Signature 2/13?

# What is the biological explanation for the mutagenic process underlying Signature 2/13?

# Which member(s) of the family is responsible for Signature 2/13?

AID

APOBEC1

APOBEC2

APOBEC3A

APOBEC3B

APOBEC3C

APOBEC3DE

APOBEC3F

APOBEC3G

APOBEC3G

APOBEC3H

APOBEC4

# Mutation signatures present in many human cancers

# Mutation signatures due to environmental exposures



tobacco

# Mutation signatures due to environmental exposures



ultraviolet radiation

# Mutation signatures due to environmental exposures
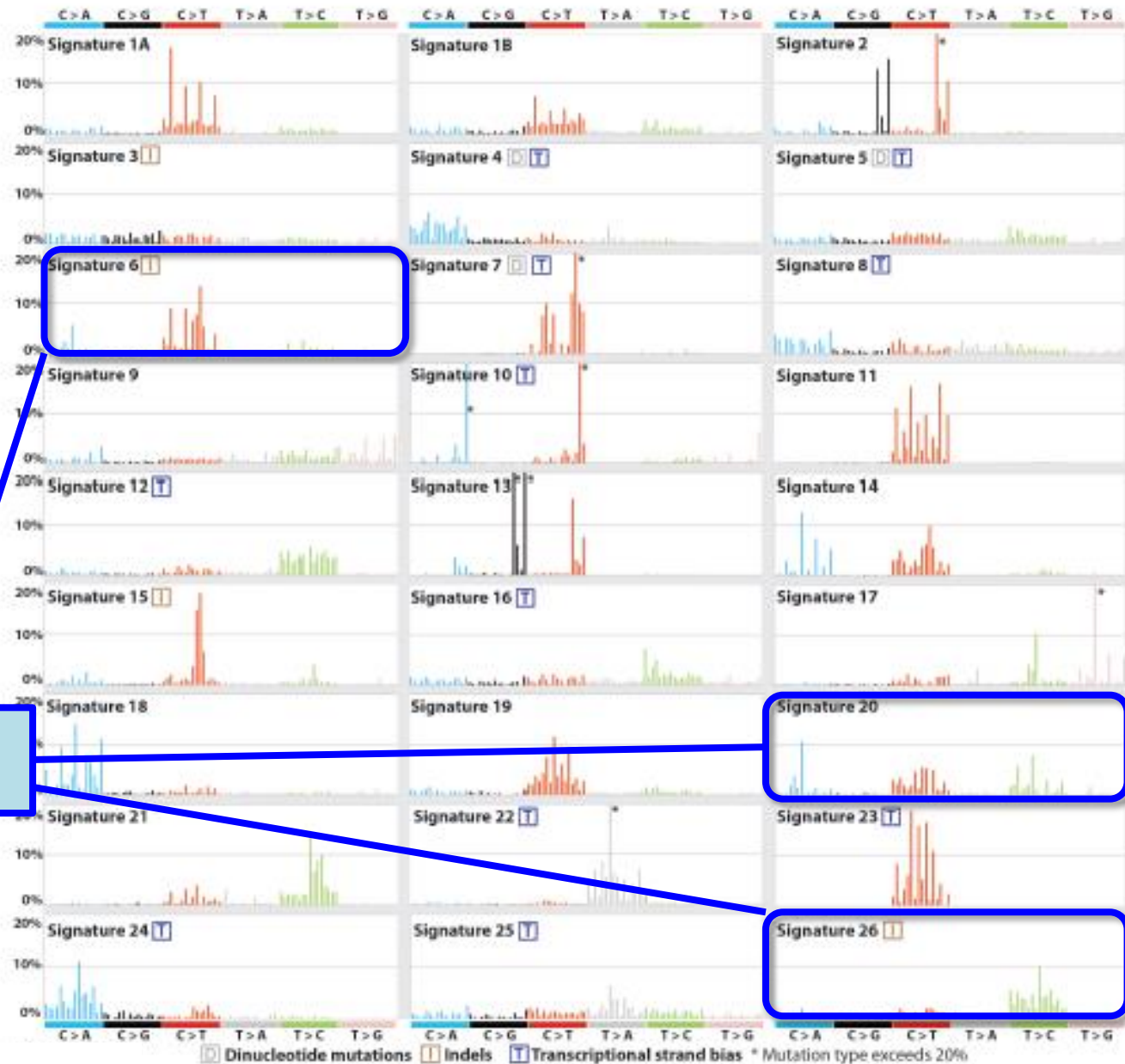


aristolochic acid

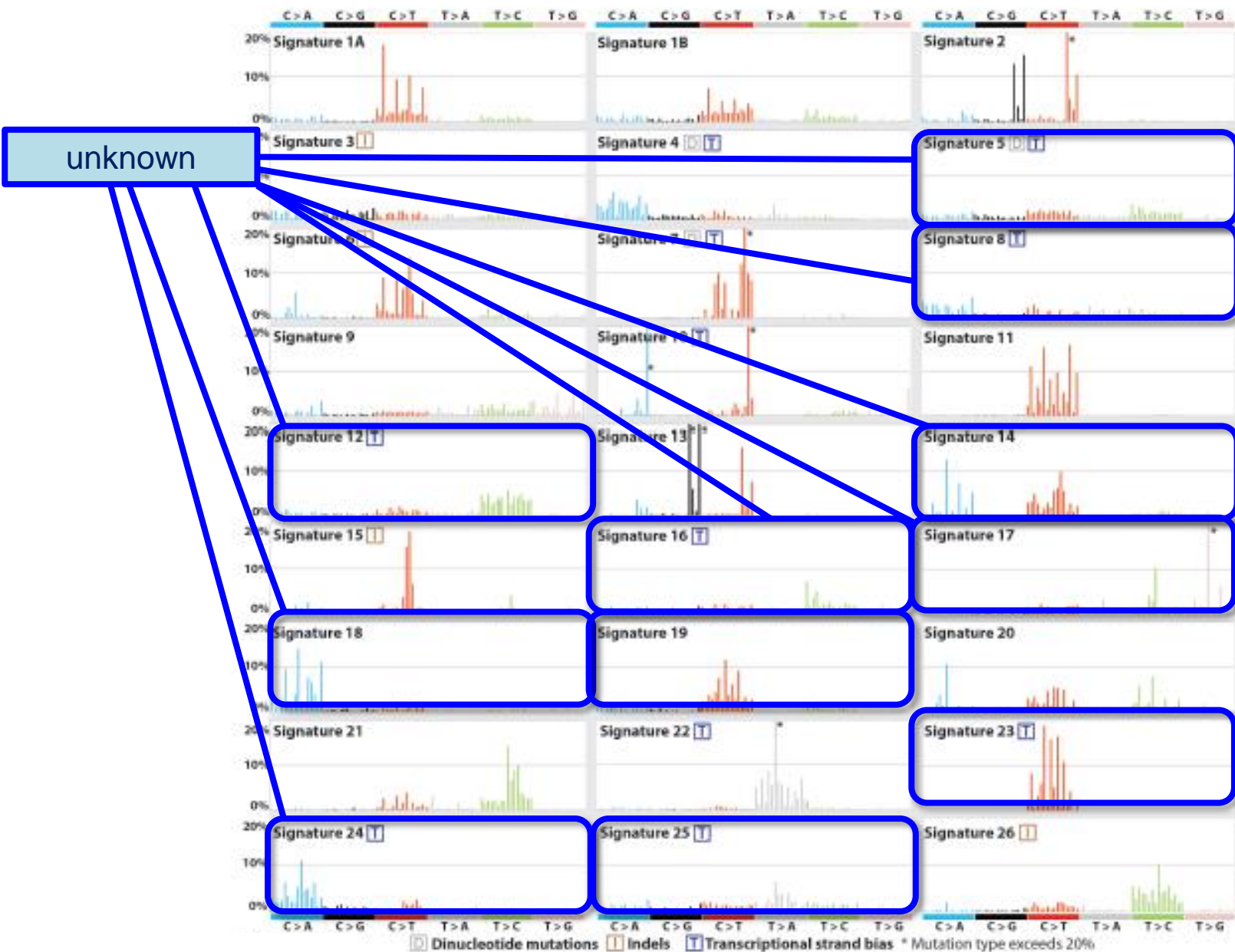# Mutation signatures due to environmental exposures



chemo-therapeutic agent

# Mutation signatures due to defective DNA repair



*BRCA1* and *BRCA2* mutations : defective homologous recombination based DSBR

# Mutation signatures due to defective DNA repair



defective mismatch repair

# Mutation signatures of unknown aetiology
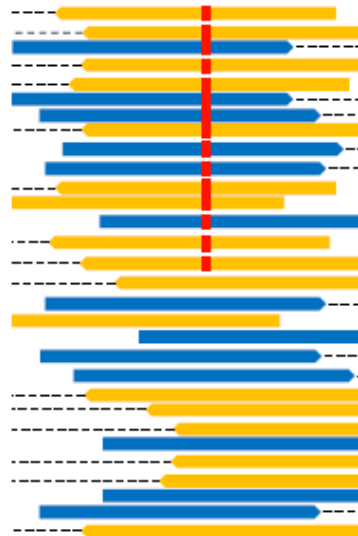
# Mutation signatures in breast cancer

# Summary II

- Modern sequencing technology permits unprecedented access to all parts of the human genome

- The enormous datasets that we can now glean through these new approaches contain a vast amount of information

- We need to ask questions of these datasets in order to extract maximum information from them
  - Ascertain all the "drivers" in a cancer
  - Use all the "passengers" to inform us about cancer biology through mutation signatures

# CONSTRUCTING EVOLUTIONARY TREES IN HUMAN BREAST CANCER

1000bp

Fraction of reads 50%

Reads from tumour cell population

Fraction of reads 35%

Reads from normal cell population

Germline

Tumour
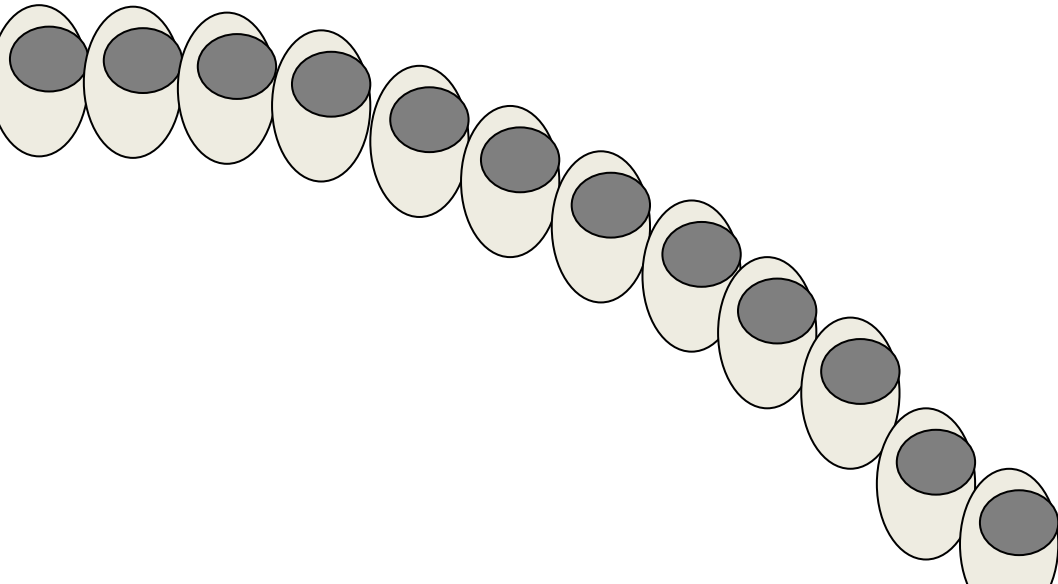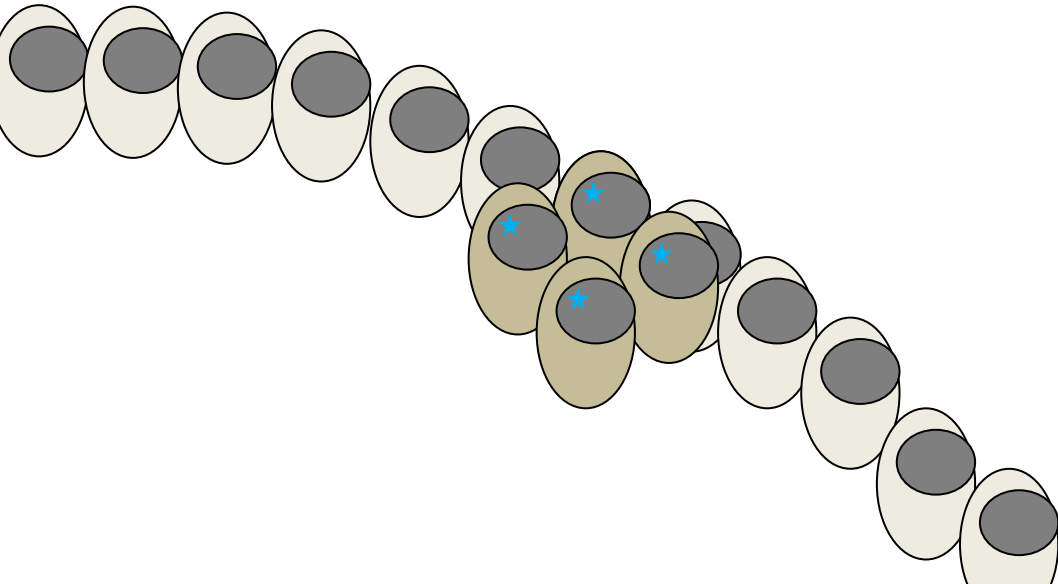
# Downstream analyses III: Cancer evolution



Ploidy: 1.80, aberrant cell fraction: 73%, goodness of fit: 94.7%

Nik-Zainal et al, Cell, 2012b

**100% cells**
26,700 mutations
10 CN changes
PIK3CA, TP53, GATA3,
SMAD4, NCOR1 muts

**18% cells**
11,000 muts

**14% cells**
Loss of 6, 8,
9, & 21

**Fertilised
egg**

**68% cells**
Del 13

**65% cells**
Del 13 &
15,600 muts

Nik-Zainal et al, Cell, 2012b

PD4120a

clonal

early          late

PD4120a

PD4120a

clonal          subclonal

early          late          subclonal

PD4120a

clonal

subclonal

early late

subclonal

PD4005a PD4107a PD4116a

PD4006a PD4109a

signature 1A
signature 2
signature 3
signature 8
signature 13

clonal

subclonal

PD4120a

early

late

subclonal

PD4005a

Early clonal

PD4107a

PD4116a

PD4006a

Early clonal

PD4109a

signature 1A
signature 2
signature 3
signature 8
signature 13

clonal | subclonal

PD4120a

early | late | subclonal

PD4005a

Early clonal
Late clonal

PD4107a

PD4116a

PD4006a

Early clonal
Late clonal

PD4109a

signature 1A
signature 2
signature 3
signature 8
signature 13

clonal                                    subclonal

early            late                              subclonal

PD4120a

PD4005a

| | | |
|---|---|---|
| Early clonal | | |
| Late clonal | | |
| Subclonal | | |

PD4107a

| | | |
|---|---|---|
| Early clonal | | |
| Late clonal | | |
| Subclonal | | |

PD4116a

| | | |
|---|---|---|
| Early clonal | | |
| Late clonal | | |
| Subclonal | | |

PD4006a

| | | |
|---|---|---|
| Early clonal | | |
| Late clonal | | |
| Subclonal | | |

PD4109a

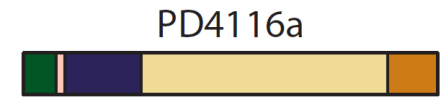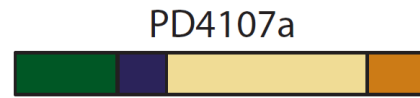| | | |
|---|---|---|
| Early clonal | | |
| Late clonal | | |
| Subclonal | | |

- signature 1A
- signature 2
- signature 3
- signature 8
- signature 13

# Summary III

- Exploiting the digital features of NGS technology, we can delve deep into the biology of tumours to gain insights into cancer evolution

- Using the totality of base subsitution mutations as well as copy number information, we can integrate this data in order to draw up phylogenetic structures of each patient's cancer
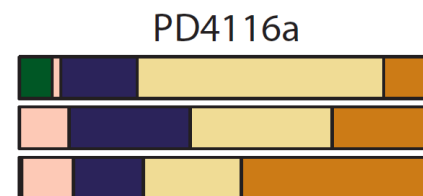
- We can identify the main cancer clone as well as subclonal populations in cancers.

- Not only can we place cancer genes within the phylogenetic tree of individual cancers, we can identify the signatures within different parts of a tree structure and examine how those signatures change in time

# WHAT DOES THE FUTURE HOLD?

# Improving genomic profiling

- Cancer genes

# Improving genomic profiling

- Cancer genes
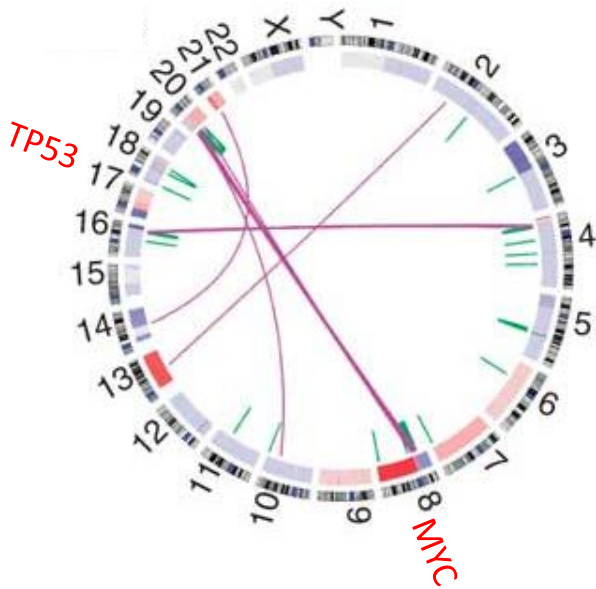- Comprehensive genomic characterisation

# Improving genomic profiling

- Cancer genes
- Comprehensive genomic characterisation
- Signatures



Signature 1

pre-treatment clonal

Drug sensitivities:
drug A

TP53

MYC

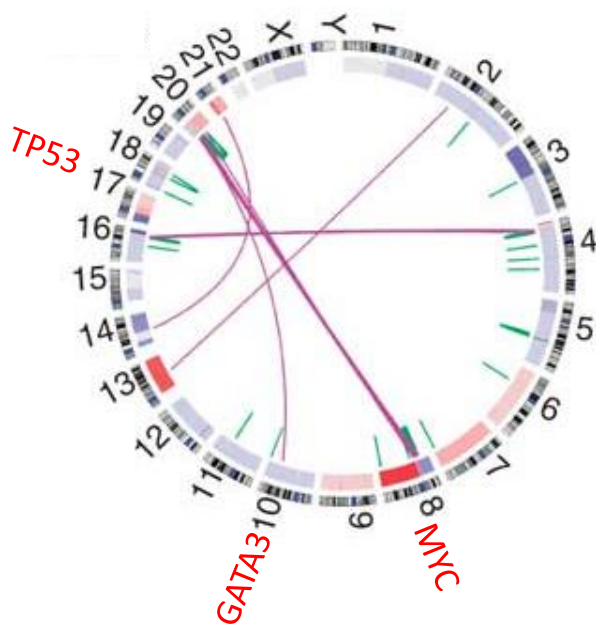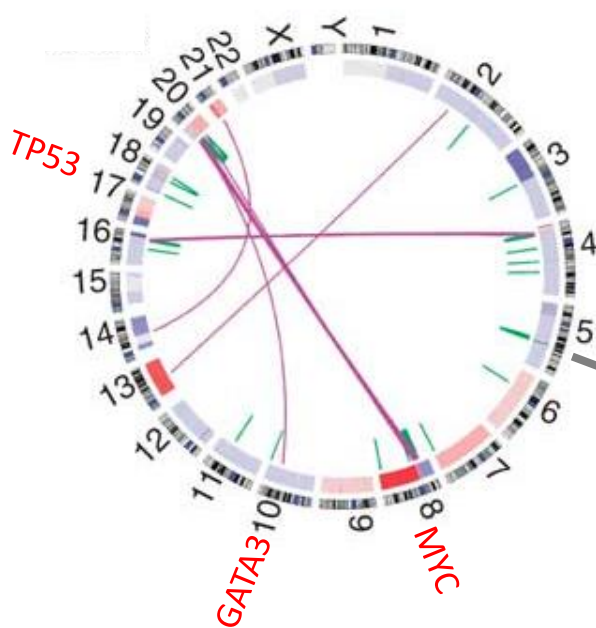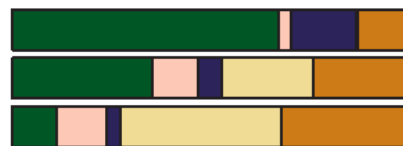# Improving genomic profiling

- Cancer genes
- Comprehensive genomic characterisation
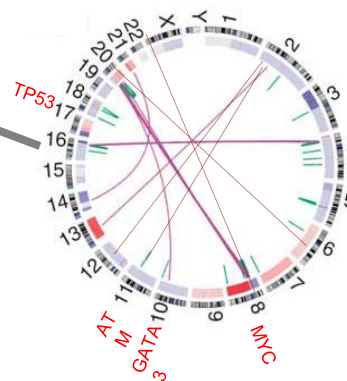- Signatures
- Subclonal populations

# Improving genomic profiling

- Cancer genes
- Comprehensive genomic characterisation
- Signatures
- Subclonal populations

# Final Summary

- The increased speed and scale of NGS technology allows the collection of vast amounts of genomic information about each person's cancer genome

- Crosstalk between clinicians, biologists and mathematicians/statisticians is required in order to extract the value-added information that is buried in cancer genomic data

- We need to have an awareness that there are still difficulties in processing and analysis data (reproducibility).

- The challenge is to design trials that best use the improved ability to stratify patients using genomic information

- Notwithstanding, there is a future to look forward to which is altogether more individual to each patient

**BREAST CANCER WORKING GROUP**

# The Cancer Genome Atlas

**Sanger Institute**
Ludmil Alexandrov
Peter van Loo
David Wedge
Patrick Tarpey
Keiran Raine
Helen Davies
Manasa Ramakrishna
Dominik Glodzik
Xueqing Zou
Sancha Martin
Andy Futreal
Ultan McDermott
Peter Campbell
Michael R Stratton

Harold Swerdlow (some NGS slides

**Breast Cancer Working Group**
Sam Aparicio
Alan Ashworth
Ake Borg
Anne-Lise Borresen-Dale
Carlos Caldas
Doug Easton
Diana Eccles
Ian Ellis
Jorunn Eyfjord
John Foekens
Louise Jones
Jocelyne Jacquemier
Jorge Reis-Filho
Sunil Lakhani
Mike Lee
Larry Norton

Angelo Paradiso
Martine Piccart
Jorge Reis-Filho
Andrea Richardson
Anne Salomon
Christos Sotiriou
Paul Spellman
Henk Stunnenberg
Fred Sweep
Benita Tan
Gilles Thomas
Andy Tutt
Laura Van t' Veer
Marc Van de Vijver

**WELLCOME-BEIT MEMORIAL**