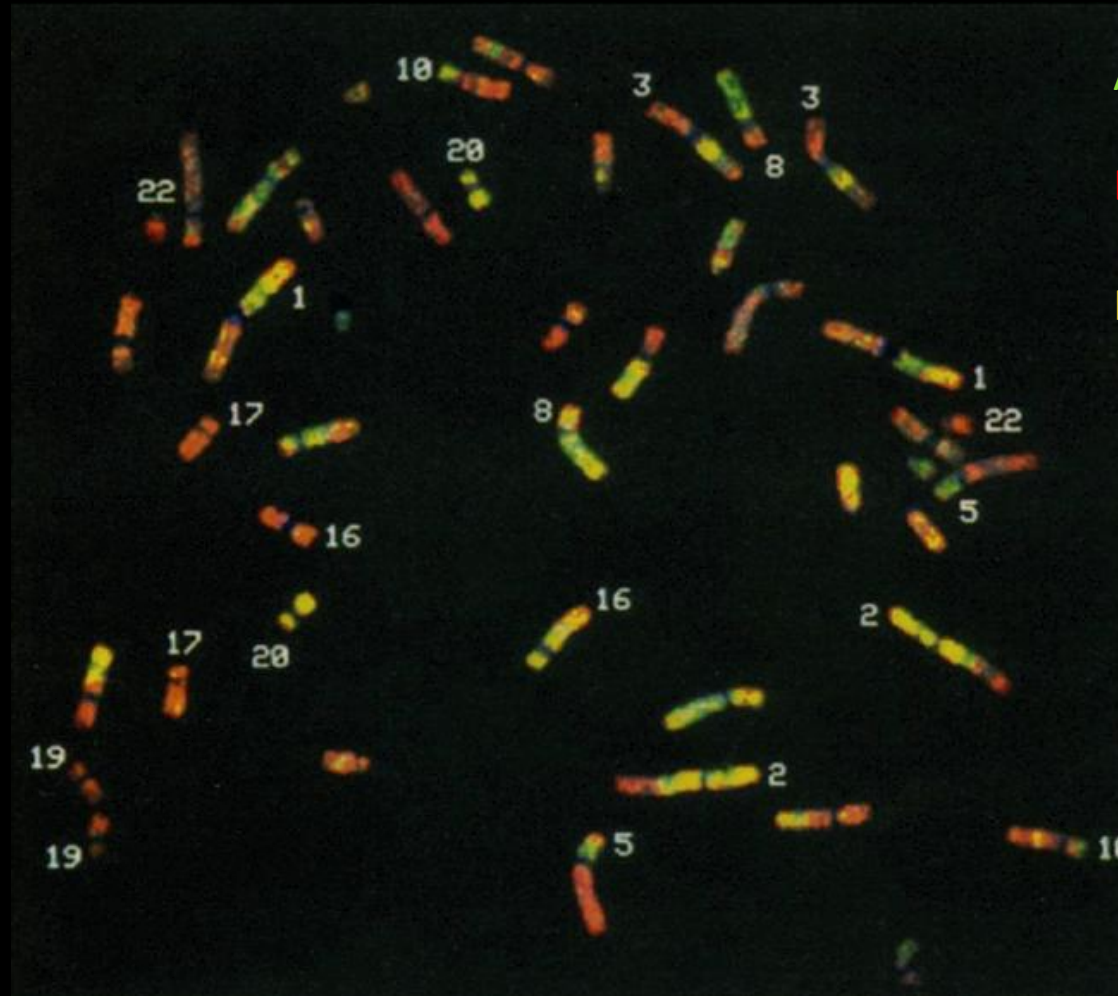# An algorithm to detect Copy Number Aberrations in cancer genomes of tumour specimens.

Arief Gusnanto, *Stefano Berri*,

Henry M. Wood and Pamela Rabbitts

# The cancer genome is often aneuploid

Amplifications

Deletions

Normal

Hartwell and Kastan. Science, 1994

# Detecting abnormalities

## Why?

- Molecular characterisation and classification of tumours

- Diagnostic, prognostic and predictive tool

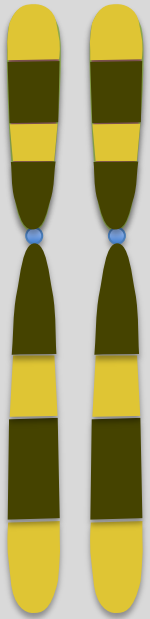- Understand the biology of cancer

## How?

- CGH

- aCGH (BAC or oligo)

- SNP microarray

- "NextGen" Sequencing
  - ✓ Tuneable resolution/cost
  - ✓ Re-use of data
  - ✓ Flexible platform
  - ✓ Technical independence Test – Control
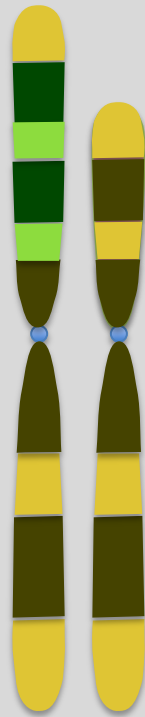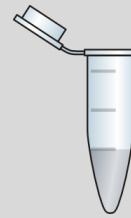  - ✓ Might become very cheap
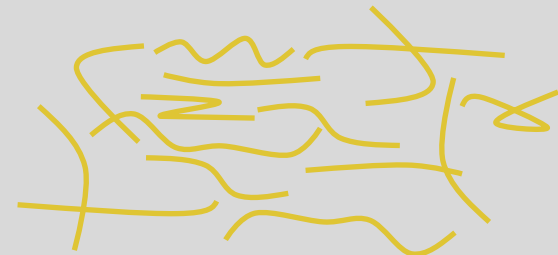
# Copy number by "NextGen" Sequencing

Normal

Tumour

Tumour

Normal

# "NextGen" sequencing and reads mapping.

UNIVERSITY OF LEEDS

Tumour

Reference genome

Normal

# Counting number of sequences for each window

Tumour

2 3 3    2    2 3 2

Normal

2 2 2    2    2 3 2

Counting number of sequences for each window

Distribution of read counts. Simulated Data, 3M reads

# Ratio Test/Control

Copy number from simulated Data

# Ratio Test/Control

UNIVERSITY OF LEEDS



**Copy number from simulated Data**

SmoothSeg. Huang et al, 2007

Copy Number

Chromosome 1, Mbp

**Copy number from simulated Data**
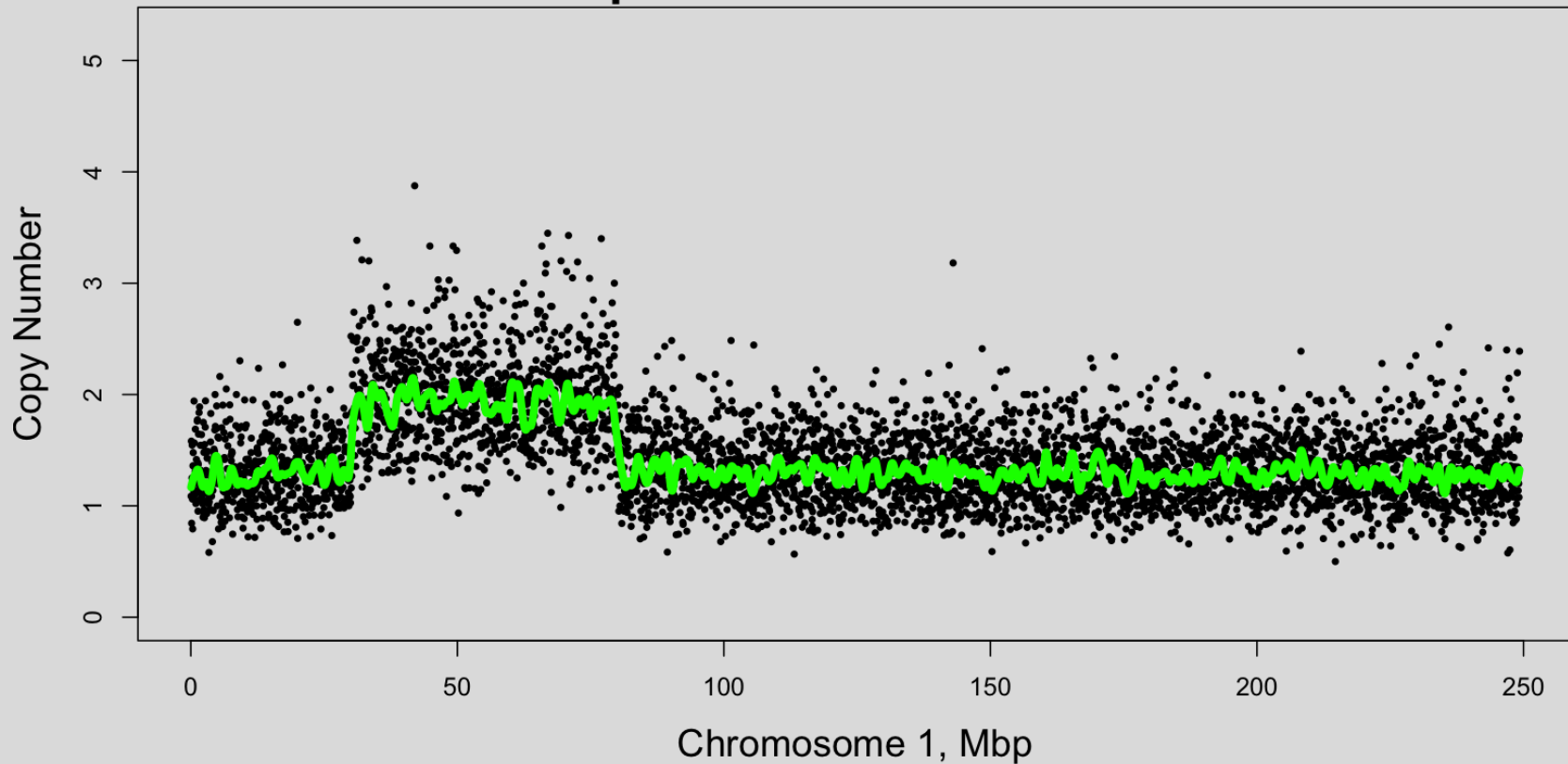
UNIVERSITY OF LEEDS



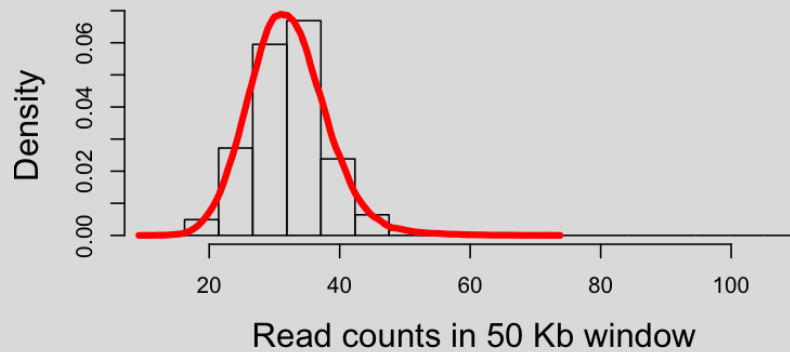Tumour

Reference genome

Normal

# Total number of reads varies.

**Copy number from simulated Data.
Unequal number of total reads**
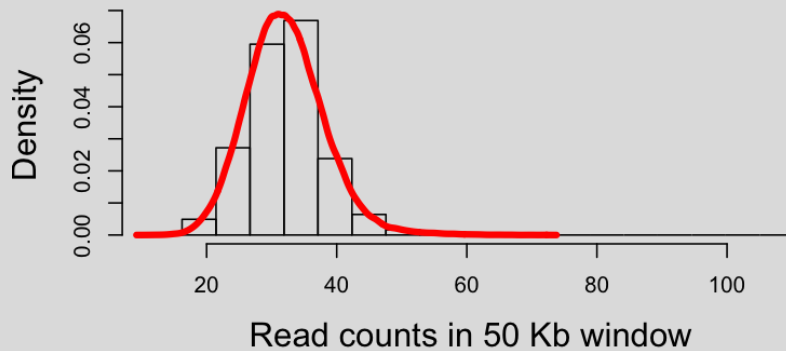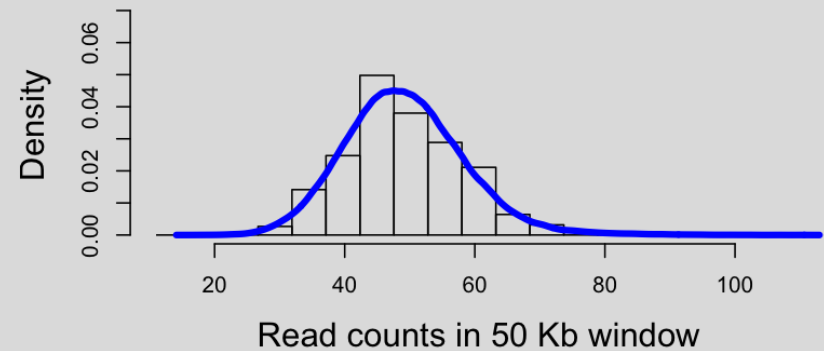
# Normalization. A crucial step

# Normalization. A crucial step

# Normalization. A crucial step

**Copy number from simulated data after median normalization.**

# The cancer genome is often aneuploid

Amplifications
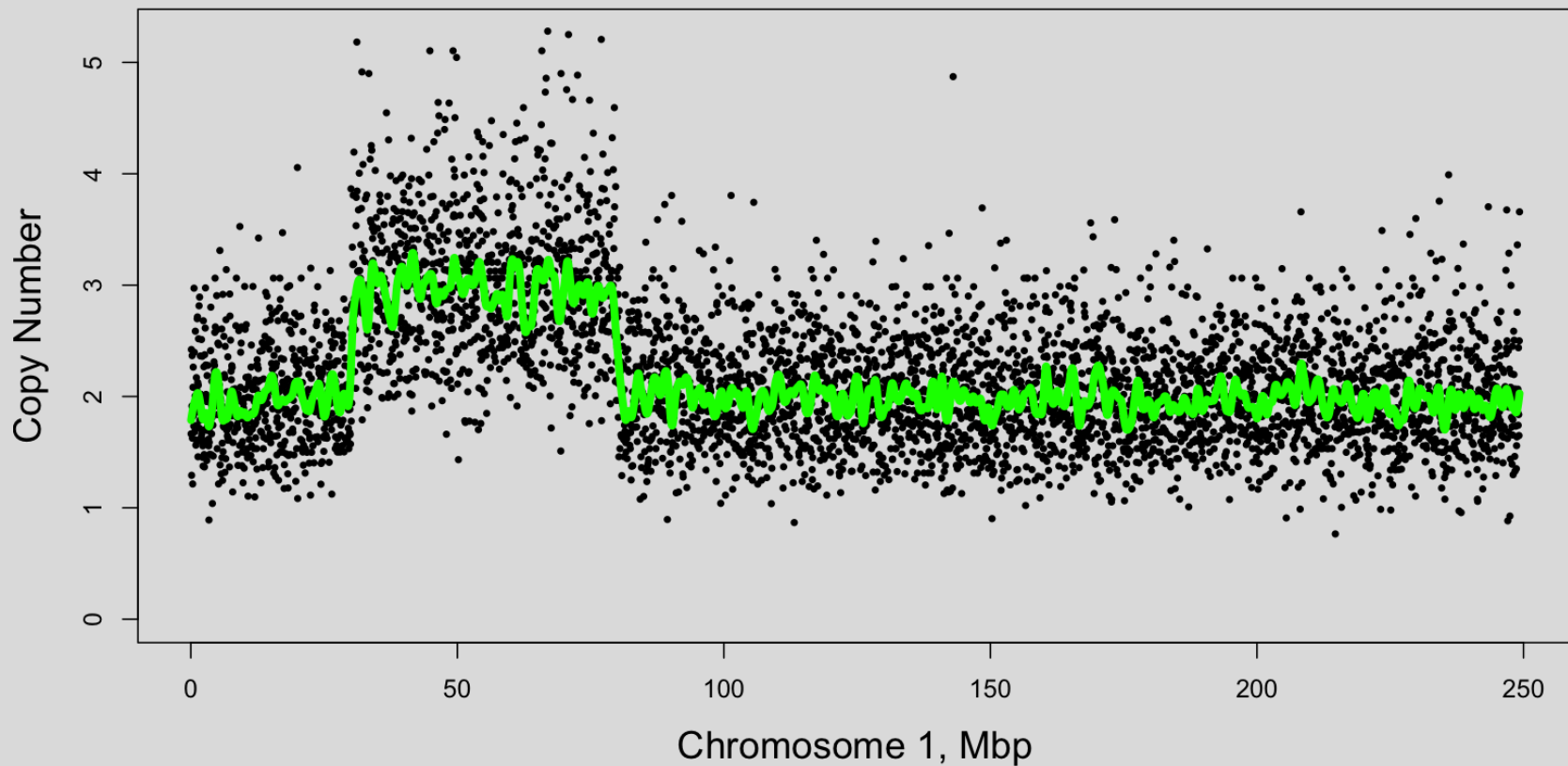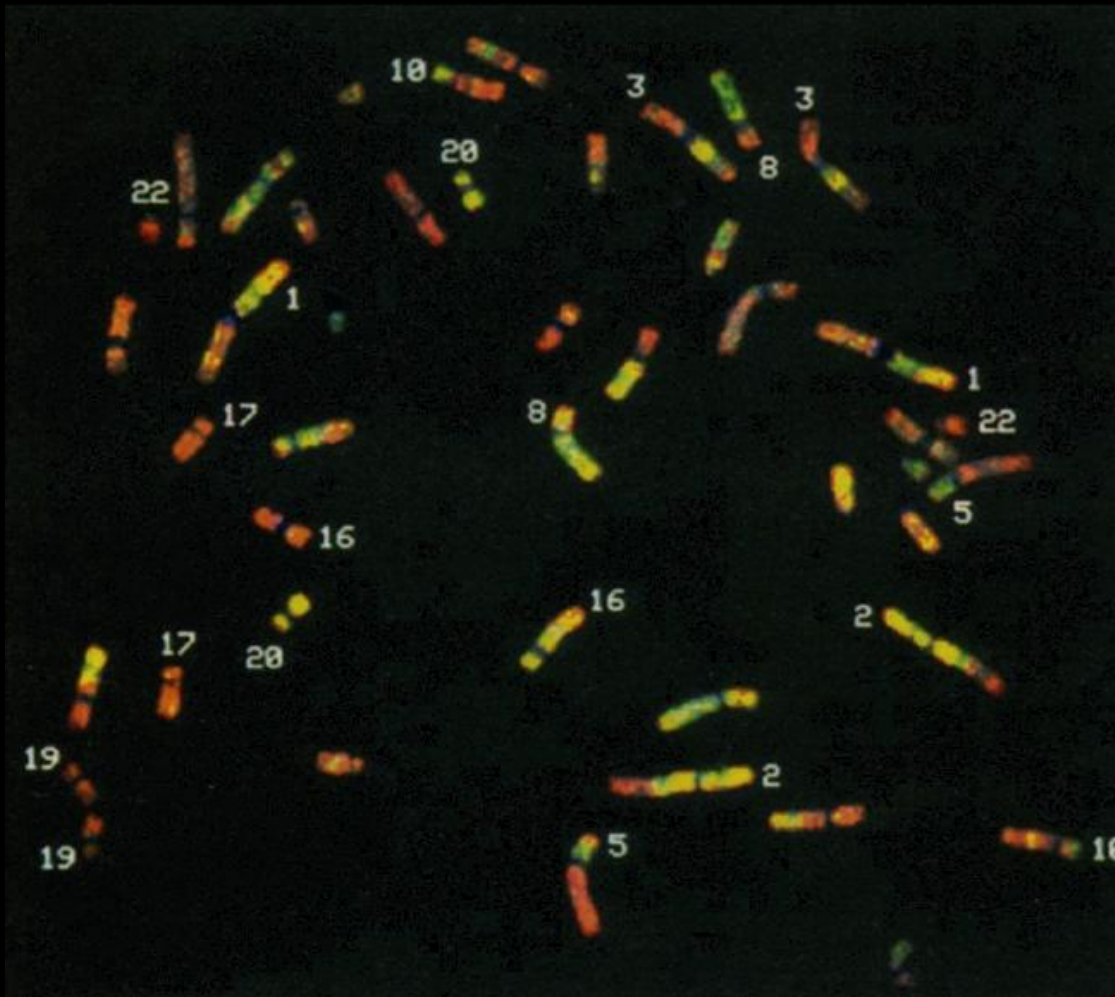
Deletions

Normal
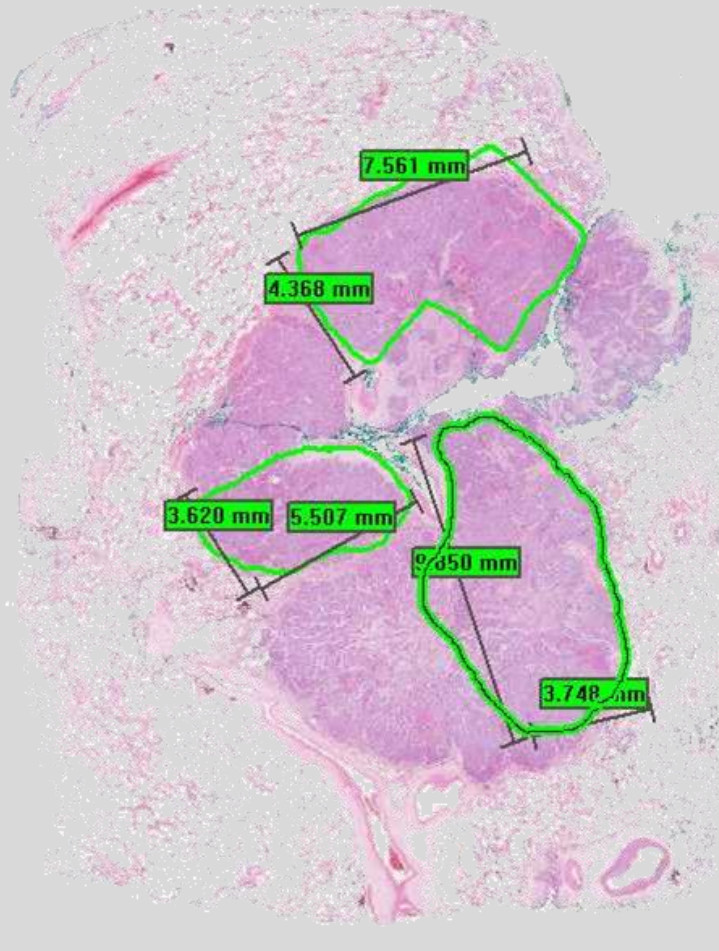
Many amplifications and deletions!

Hartwell and Kastan. Science, 1994

# Patient's tumour samples

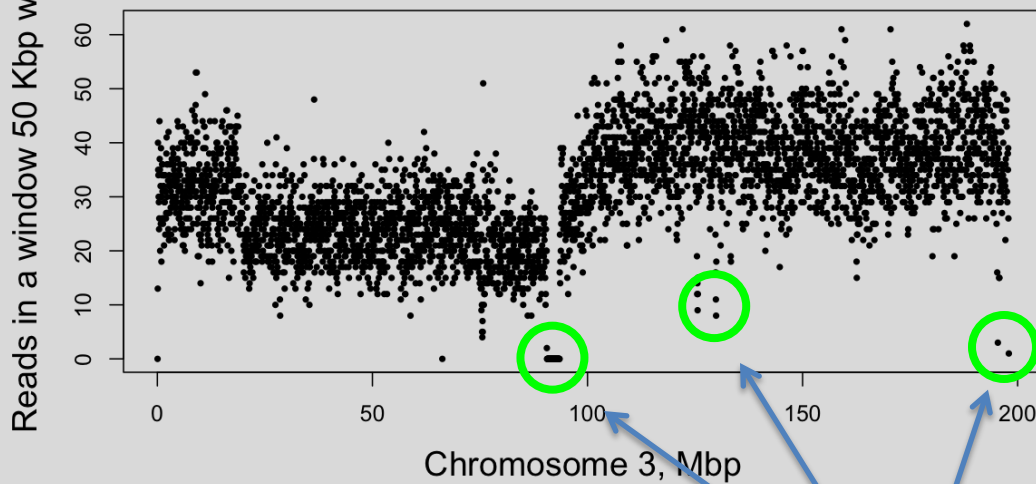- Contamination with stroma, inflammatory cells...

Lung tumour

# The real samples. A lot noisier

UNIVERSITY OF LEEDS



Distribution of read counts. Patient's specimens
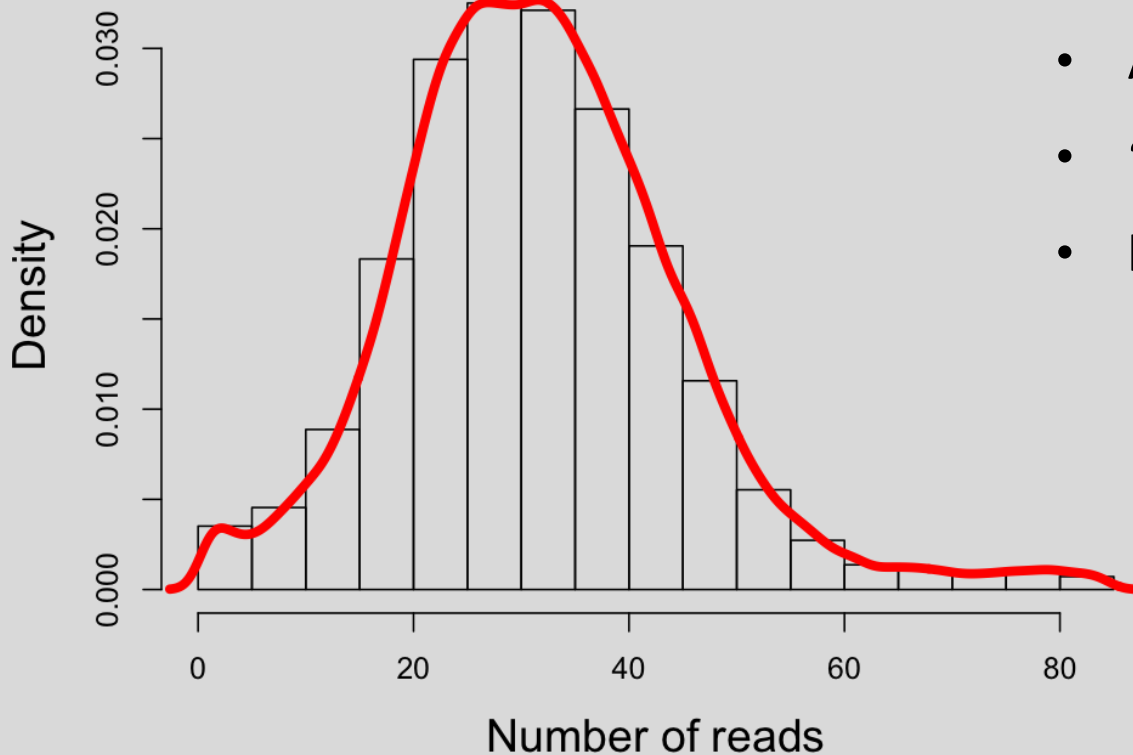
Aligning artefacts

- Some sequences cannot be aligned (repeated regions)

- GC content bias

- Unequal number of total reads.

- Extra noise of unknown origin
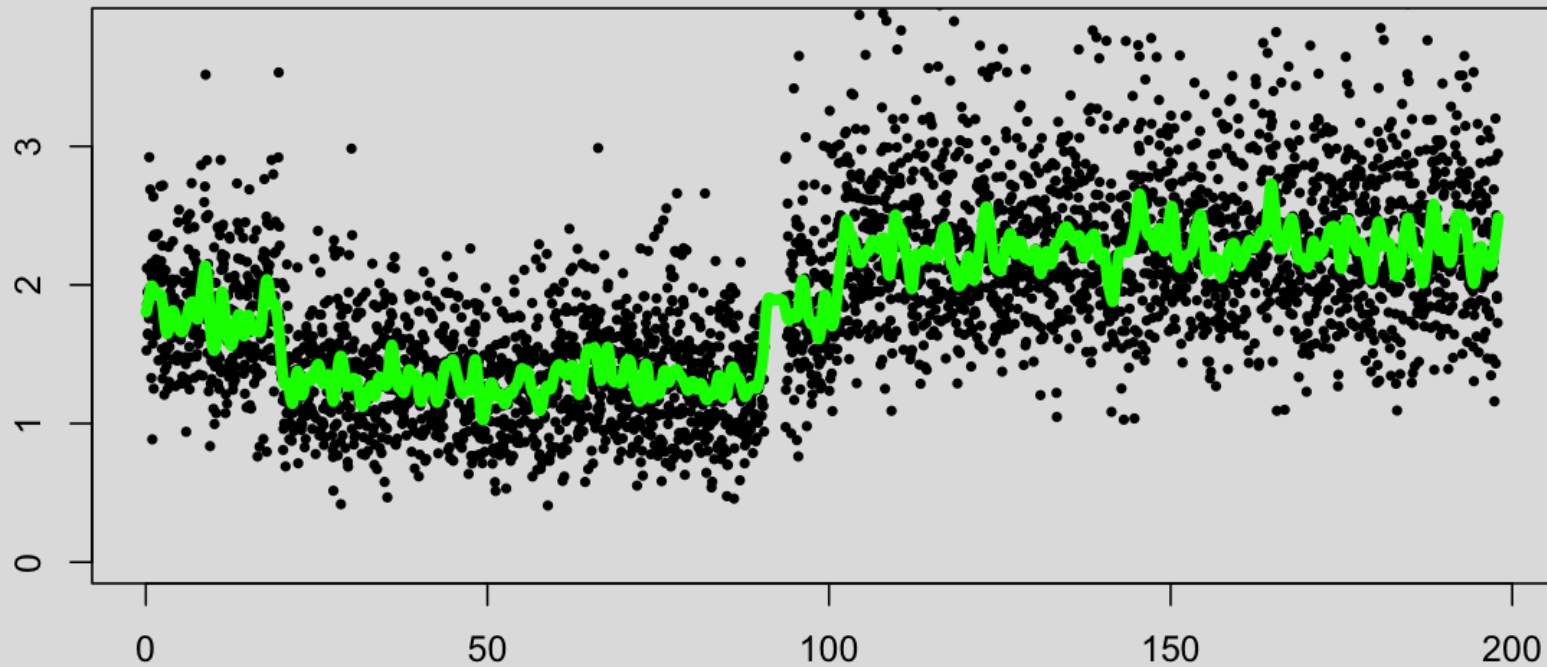
# The median might be meaningless

Patient's sample

- Asymmetric distribution

- "Flat top"

- long tail

# Median normalisation

Median-normalized data. Patient's specimens

UNIVERSITY OF LEEDS



Median-normalized data. Patient's specimens

# Discrete data normalisation



**Patient's sample, segmented data**

(Density vs Ratio Test/Normal)

# Discrete data normalisation

Patient's sample, segmented data

1 (2)
1.5 (3)
0.5 (1)
2 (4)
0 (0)

Ratio (Copy num)

Density

Ratio Test/Normal

# Discrete data normalisation
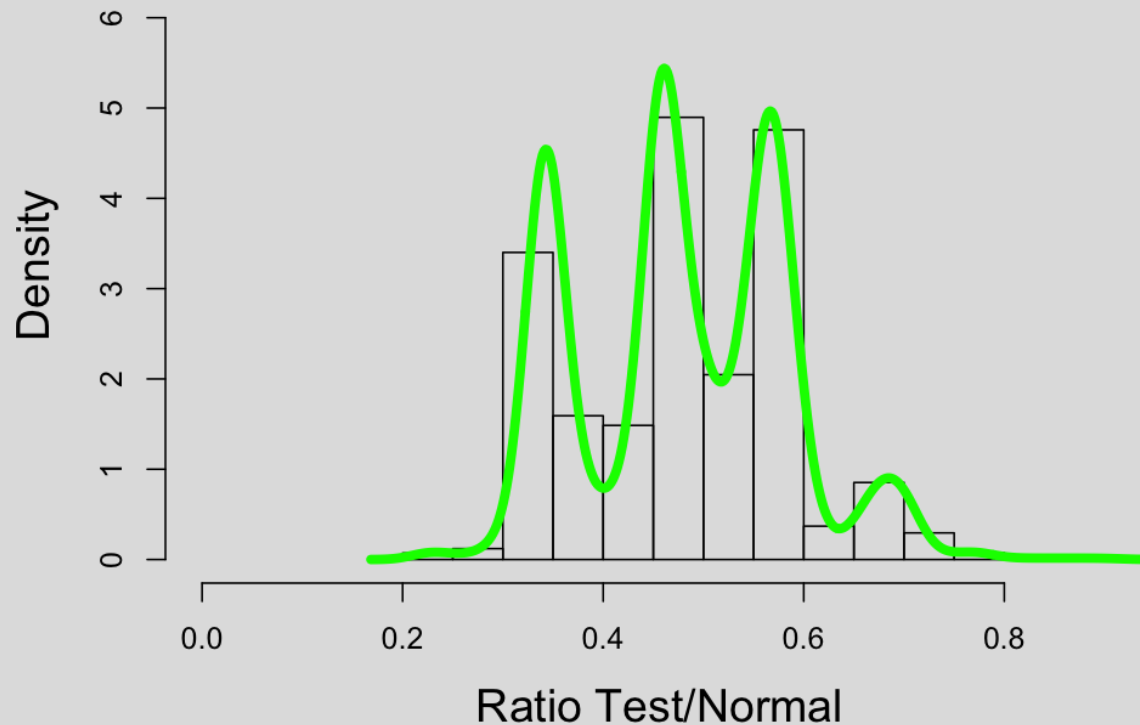
Patient's sample, segmented data

1 (2)

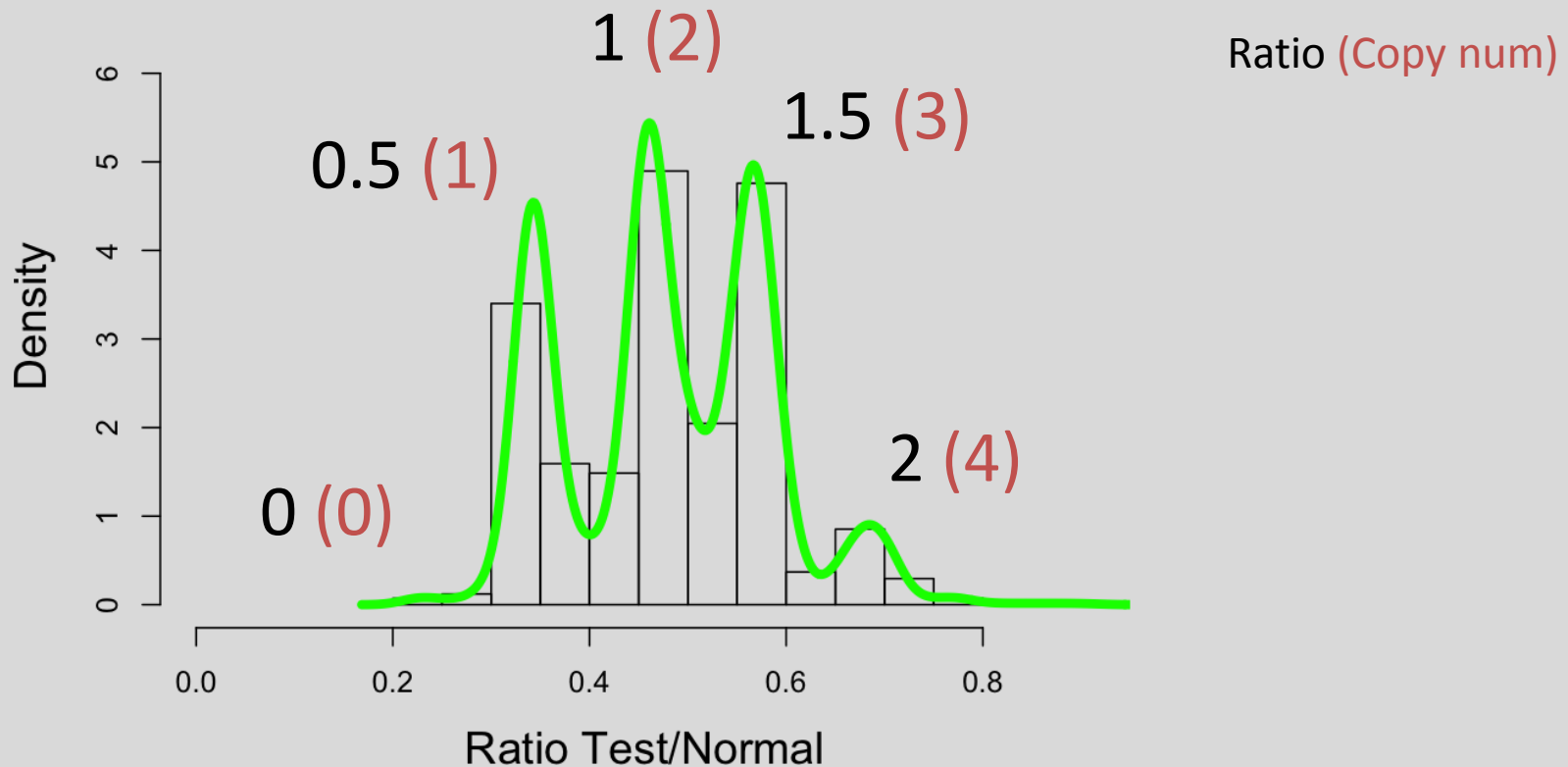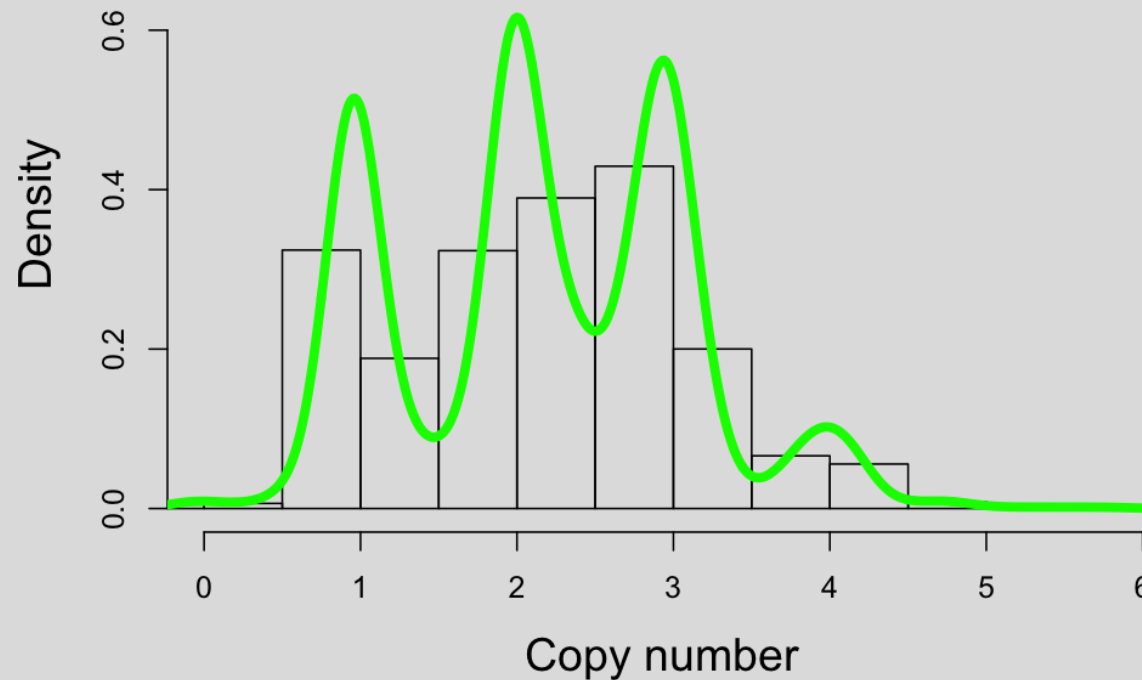Ratio (Copy num)
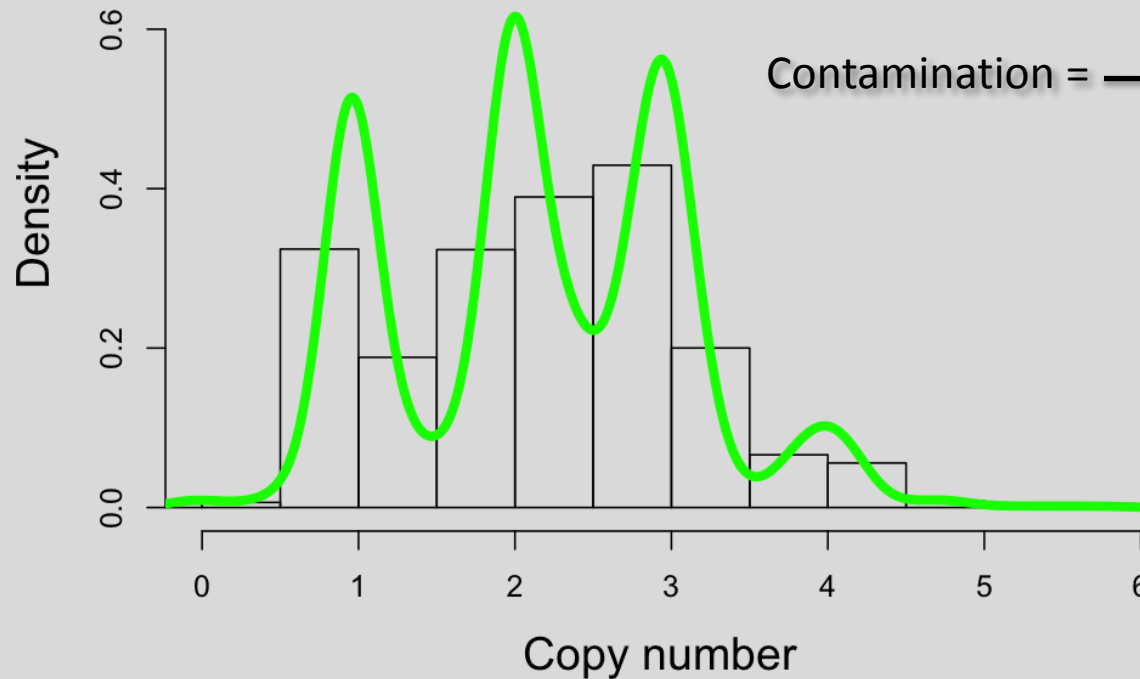
# Discrete data normalisation



Patient's sample, segmented data

# Discrete data normalisation



**Patient's sample, segmented data**

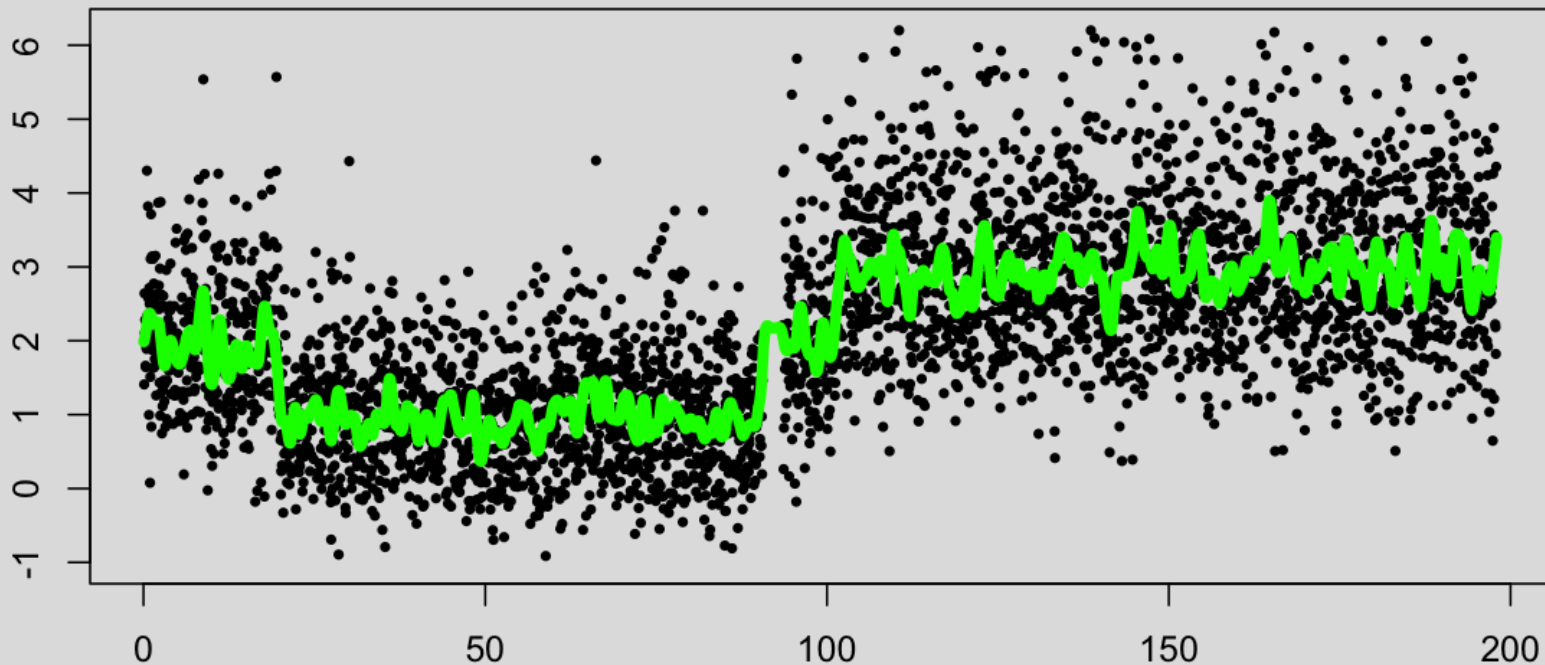$$\text{Contamination} = \frac{\text{Tumour tissue}}{\text{Total tissue}} = 51\%$$

Discrete normalization. Patient's specimens

# Conclusions

- Develop a novel normalisation method for "NextGen" data that can cope with

  ✓ Highly abnormal genomes

  ✓ Tumour samples contaminated by normal cells

- We can estimate contamination percentage.

UNIVERSITY OF LEEDS

- Contamination is allowed, but otherwise the tumour should be homogeneous.

- Process might require human supervision when calling discrete states.

# Acknowledgements