**HSTalks**

Dr. Giles Yeo – University of Cambridge, UK
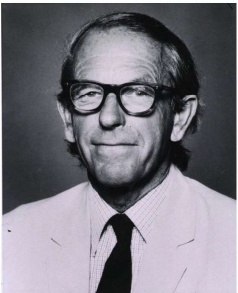


Genomics 101:
an Introduction
to Sequencing

**Dr. Giles Yeo**
Principle Research Associate
Department of Clinical Biochemistry
University of Cambridge, UK
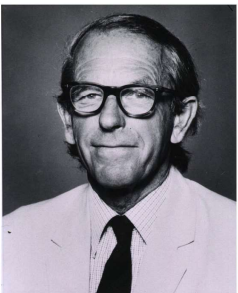
1

---

**Fred Sanger (1918 – 2013)**



**1958**

**1980**

**Sanger Institute**

2

---

**Fred Sanger (1918 – 2013)**



**1958**
Nobel Prize Chemistry
(Protein sequence of insulin)

**1980**
Nobel Prize Chemistry
(Nucleic acid sequencing)

The Sanger Institute is named
after Frederick Sanger

---

Dr. Giles Yeo – University of Cambridge, UK

## Sanger sequencing method

**Chain termination with a specific ddNTP (dideoxynucleotides)**

When a ddNTP binds, the DNA strand will stop extending

3

## Sanger sequencing method

**Chain termination with a specific ddNTP (dideoxynucleotides)**

Template: CGAGTCCTTAGGCATACA          dNTP and DNA polymerase
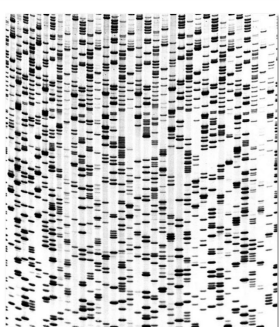
Primer: GCTCAG

ddT

Template: CGAGTCCTTAGGCATACA

GCTCAGGAAddT

GCTCAGGAATCCGddT

GCTCAGGAATCCGTddT

GCTCAGGAATCCGTATGddT

A 'targeted' sequencing approach; *i.e.* you know the surrounding sequence

Original sequencing used radioactive ddTs, so you could see it using X-ray

## Sanger sequencing method

**Chain termination with a specific ddNTP (dideoxynucleotides)**

Template: CGAGTCCTTAGGCATACA          dNTP and DNA polymerase

Primer: GCTCAG                                      +ddA    +ddG
                                                          +ddT    +ddC

ddT

Template: CGAGTCCTTAGGCATACA          A  T  G  C

GCTCAGGAAddT

GCTCAGGAATCCGddT

GCTCAGGAATCCGTddT

GCTCAGGAATCCGTATGddT

Dr. Giles Yeo – University of Cambridge, UK

**X-ray film of DNA sequencing**

4

**Fluorescent Sanger sequencing**

A A A A*
A A A A A
A* A A*

T T T T* T
T T T T*
T* T

G G G G
G* G G*
G G G*
G*G G

C C* C
C* C C C
C C C
C* C C

**Radioactivity was replaced with fluorescence**

5

**Fluorescent Sanger sequencing**

A A A A*
A A A A A
A* A A*

T T T T* T
T T T T*
T* T

G G G G
G* G G*
G G G*
G*G G

C C* C
C* C C C
C C C
C*

Dr. Giles Yeo – University of Cambridge, UK

## Fluorescent Sanger sequencing

A*      T*      G*      C*

Specific primer

ATAGTTAAGCGGGTT*
ATAGTTAAGCGGGT*
ATAGTTAAGCGGG*
ATAGTTAAGCGG*
ATAGTTAAGCG*
ATAGTTAAGC*
ATAGTTAAG*
ATAGTTAA*
ATAGTTA*
ATAGTT*
ATAGT*
ATAG*
ATA*
AT*
A*

Chromatogram

A T A G T T  A A G C G G G T T

## Fluorescent DNA image

6

## Point and in/del mutations

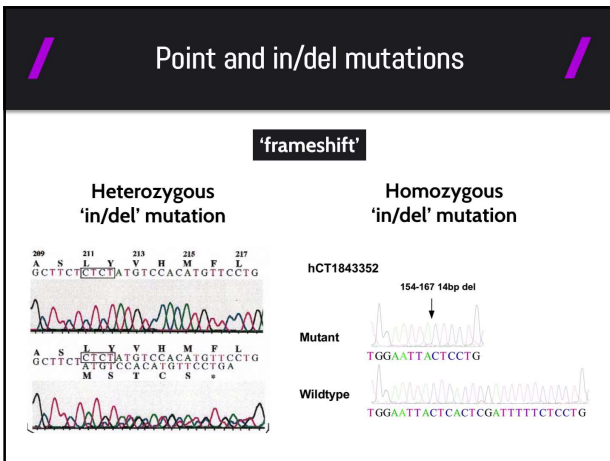'missense'

**Heterozygous 'point mutation'**

Y722C

718  719  720  721  722  723
Y    S    T    D    Y    Y
T A C A G C A C T G  A C T A C T A C

Y    S    T    D    C    Y
T A C A G T C A C T G  A C T    G C T A C

**Homozygous 'point mutation'**

R316Q

C C T C A G T T T    C C T C G G T T T

Affected child          Control

7

The screen versions of these slides have full details of copyright and acknowledgements

Dr. Giles Yeo – University of Cambridge, UK

**Point and in/del mutations**

'frameshift'

Heterozygous 'in/del' mutation | Homozygous 'in/del' mutation

hCT1843352

154-167 14bp del

Mutant — TGGAATTACTCCTG

Wildtype — TGGAATTACTCACTCGATTTTTCTCCTG

**Point and in/del mutations**

'frameshift'

Heterozygous 'in/del' mutation | Homozygous 'in/del' mutation

hCT1843352

154-167 14bp del

Mutant — TGGAATTACTCCTG

Wildtype — TGGAATTACTCACTCGATTTTTCTCCTG

Deleted section of DNA

**Cost of sequencing the human genome**

~ $3 billion USD

The first Human Genome Project
- Total project cost -

8

The screen versions of these slides have full details of copyright and acknowledgements
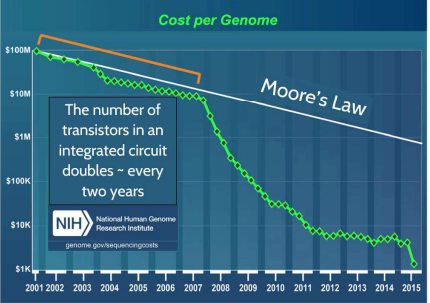
5

Dr. Giles Yeo – University of Cambridge, UK

## Cost of sequencing the human genome

~$10 million USD

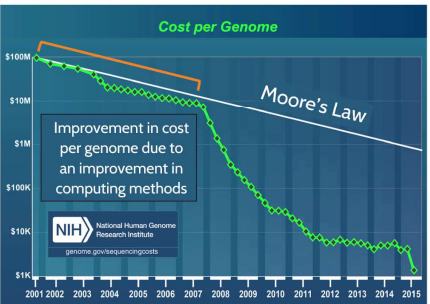Cost of sequencing a genome today using Sanger method and current technology

## Cost per genome

Cost per Genome

$100M

$10M

Moore's Law

$1M

The number of transistors in an integrated circuit doubles ~ every two years

$100K

NIH National Human Genome Research Institute

$10K

genome.gov/sequencingcosts

$1K

2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015

For further information, please see the tab of external links

9

## Cost per genome

Cost per Genome

$100M

$10M

Moore's Law

$1M

Improvement in cost per genome due to an improvement in computing methods

$100K

NIH National Human Genome Research Institute

$10K

genome.gov/sequencingcosts

$1K

2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015

For further information, please see the tab of external links

HSTalks

Dr. Giles Yeo – University of Cambridge, UK

## Cost per genome

Cost per Genome

$100M

$10M
Moore's Law

**Introduction of 'next-gen' sequencing**

$1M

$100K

$10K
NIH National Human Genome Research Institute
genome.gov/sequencingcosts

$1K
2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015

For further information, please see the tab of external links

---

**What is this 'next-gen' (now gen?) sequencing?**

10

---

## Characteristics of 'next-gen' sequencing

| Next-generation sequencing | Sanger sequencing |
| --- | --- |
| 'Polony sequencing' | Read lengths of 1 - 2 kbp (1000 - 2000 bp) |
| PCR colony | |
| Short reads of 50-250 bps | Needed to know what you were sequencing |
| Random 'shot-gun' sequencing | Targeted sequencing |

11

HSTalks

Dr. Giles Yeo – University of Cambridge, UK

## Characteristics of 'next-gen' sequencing

**Many platforms for next generation sequencing**
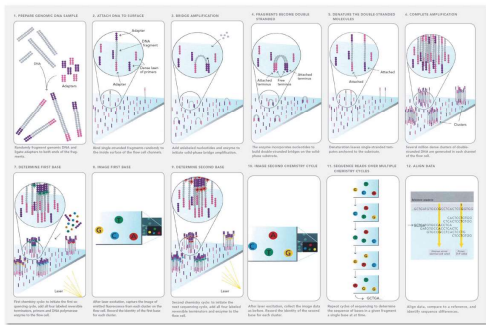
**Illumina**

Illumina Genome Analyser

**The basic principles apply to the various platforms available**
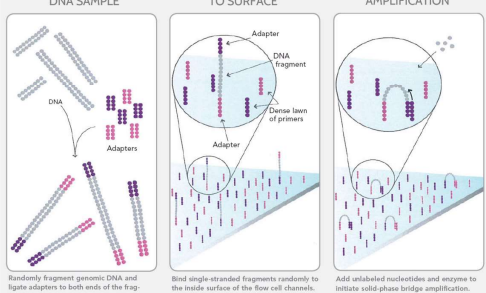
## 1. Bridge amplification

12

## 1. Bridge amplification

8

Dr. Giles Yeo – University of Cambridge, UK

## 1. Bridge amplification

1. PREPARE GENOMIC DNA SAMPLE

- Randomly fragment the DNA

- Create blunt ends

- Ligate the adapters to the fragments
  - Adapters are known sequences

Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

## 1. Bridge amplification

1. PREPARE GENOMIC DNA SAMPLE

2. ATTACH DNA TO SURFACE

3. BRIDGE AMPLIFICATION

Adapter

DNA fragment

Dense lawn of primers

Adapter

Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

## 1. Bridge amplification

3. BRIDGE AMPLIFICATION

- 'PCR-like': Add enzyme, nucleotides, buffer

- The DNA fragments are tethered to the glass slide

- The DNA fragments bind to an adjacent primer

- Bridge amplification can take place

Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

Dr. Giles Yeo – University of Cambridge, UK

HS**Talks**

Dr. Giles Yeo – University of Cambridge, UK







The screen versions of these slides have full details of copyright and acknowledgements

Dr. Giles Yeo – University of Cambridge, UK

## 4. Sequencing over multiple cycles



10. IMAGE SECOND CHEMISTRY CYCLE

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES

12. ALIGN DATA

## Image of fluorescent signals



▌ Fluorescent signals

▌ Randomly spaced

▌ However, they are a sufficient distance apart from each other

Voelkerding K.V. *et al.*, *Clinical Chemistry*. 2009; 55(4):641-58

16

## Flow cells



4 rows per 'flow cell'

**Illumina, NovaSeq flow cell**

▌ 2,500,000,000 'polonies' per row

▌ 300 base pairs per polony

**Therefore**

▌ 4 rows x 2,500,000,000 polonies x 300bp

= 3 TRILLION bases / run  (44 hrs)

▌ >99% accuracy

17

Dr. Giles Yeo – University of Cambridge, UK

## Flow cells

**Random massively parallel sequencing!**

**Problem?**

An IMMENSE amount of data

1000s of super high-resolution images per run!

Cheaper to repeat the sequencing than to store the data!

These servers & hard-drives (in which data is stored) are extremely heavy!

S4

## Flow cells

**Random massively parallel sequencing!**

**Problem?**

**Therefore, after extracting the data, the images are discarded**

amount of data

high-resolution per run!

the sequencing than to store the data!

These servers & hard-drives (in which data is stored) are extremely heavy!

## 'Paired-end' reads increase accuracy

18

**There are two primers on the DNA fragment**

**Sequence from both ends**

Paired-End Reads

Alignment to the Reference Sequence

Read 1

Read 2

Reference

Repeats

**Figure : Paired-End Sequencing and Alignment**—Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

Dr. Giles Yeo – University of Cambridge, UK

## 'Paired-end' reads increase accuracy

| 'Paired-end' reading | Increases the accuracy of sequencing |
| | Each base pair is represented twice |

Paired-End Reads

*Better mapping*

Read 1

Read 2

Alignment to the Reference Sequence

Reference

Repeats

*Better accuracy*

Figure : **Paired-End Sequencing and Alignment**—Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

## Multiplexing

**Multiplexing** allows for an increased number of samples

19

## Multiplexing

A
Library Preparation

B
Pool

Index 1
(CATTCG)

Index 2
(AACTGA)

— Library 1 Barcode
— Library 2 Barcode
— Sequencing Reads
— DNA Fragments
— Reference Genome

Add a DNA index while making the library

Figure: Library Multiplexing Overview

Dr. Giles Yeo – University of Cambridge, UK

## Multiplexing

**A**
Library Preparation

**B**
Pool

Index 1
(CATTCG)

Index 2
(AACTGA)

— Library 1 Barcode
— Library 2 Barcode
— Sequencing Reads
— DNA Fragments
— Reference Genome

Figure: Library Multiplexing Overview

Mix the no. of samples down one flow cell

## Multiplexing

**A**
Library Preparation

**B**
Pool

**C**
Sequence

Index 1
(CATTCG)

Index 2
(AACTGA)

Sequence Output
to Data File

CATTCGACGGATCG
AACTGAGTCCGATA
AACTGATCGGATCC
CATTCGTGGCAGTC
AACTGAACCTGATG
AACTGAGATTACAA
CATTCGCAGTTCATT
CATTCGAACTTCGA

— Library 1 Barcode
— Library 2 Barcode
— Sequencing Reads
— DNA Fragments
— Reference Genome

Figure: Library Multiplexing Overview

Can do all of the sequencing together

## Multiplexing

**A**
Library Preparation

**B**
Pool

**C**
Sequence

**D**
Demultiplex

**E**
Align

Index 1
(CATTCG)

Index 2
(AACTGA)

Sequence Output
to Data File

CATTCGACGGATCG
AACTGAGTCCGATA
AACTGATCGGATCC
CATTCGTGGCAGTC
AACTGAACCTGATG
AACTGAGATTACAA
CATTCGCAGTTCATT
CATTCGAACTTCGA

CATTCGACGGATCG
CATTCGTGGCAGTC
CATTCGCAGTTCATT
CATTCGAACTTCGA

AACTGAGTCCGATA
AACTGATCGGATCC
AACTGAACCTGATG
AACTGAGATTACAA

— Library 1 Barcode
— Library 2 Barcode
— Sequencing Reads
— DNA Fragments
— Reference Genome

Figure: Library Multiplexing Overview

Deconvolute the sequence in the analysis

Dr. Giles Yeo – University of Cambridge, UK

**'Whole exome sequencing'**

- **Exome**: all exons within a genome

- **Target enrichment**
  The coding region of the genome
  represents ~2% of the total genome

- **How do you enrich for this 2%?**

20

**'Whole exome sequencing'**

You can pull out the fragments
of DNA you want to sequence

Hybridization in Solution | Solid Phase Hybridization

**'Whole exome sequencing'**

You can pull out the fragments
of DNA you want to sequence

**Solid Phase Hybridization**

Fix the 'bait'
onto a glass slide

Wash the fragmented
genome over the slide

Individual primers can bind
to the desired fragment

Dr. Giles Yeo – University of Cambridge, UK

## 'Whole exome sequencing'

**You can pull out the fragments of DNA you want to sequence**

| Hybridization in Solution | Solid Phase Hybridization |
|---|---|

**Droplet PCR**

Merge the primers with bubbles of DNA

## Mosaic mutations

**The cost of sequencing is being reduced**

21

## Mosaic mutations

**An Activating Mutation of *AKT2* and Human Hypoglycemia**

K. Hussain[1,*], B. Challis[2,*], N. Rocha[2,*], F. Payne[3], M. Minic[2], A. Thompson[3], A. Daly[3], C. Scott[3], J. Harris[2], B.J.L. Smillie[2], D.B. Savage[2], U. Ramaswami[4], P. De Lonlay[5], S. O'Rahilly[2], I. Barroso[2,3], and R.K. Semple[2]

*Science.* 2011 October 28; 334(6055): 474. doi:10.1126/science.1210878.

▌ Digital sequencing

▌ You can count the copies of every fragment

  ▌ Each of the base pairs is sequenced hundreds of times

Dr. Giles Yeo – University of Cambridge, UK

## Mosaic mutations

- Can pick up **mosaic mutations**
  - This is difficult with Sanger sequencing
  - **Next Gen Sequencing:** picks up a small percentage of mutant fragments of DNA

*NOT every cell carries the mutation*

Mosaic overgrowth with fibroadipose hyperplasia is caused by somatic activating mutations in *PIK3CA*

Marjorie J Lindhurst[1,16], Victoria E R Parker[2,16], Felicity Payne[3], Julie C Sapp[1], Simon Rudge[4], Julie Harris[2], Alison M Witkowski[1], Qifeng Zhang[4], Matthijs P Groeneveld[2], Carol E Scott[3], Allan Daly[3], Susan M Huson[5], Laura L Tosi[6], Michael L Cunningham[7], Thomas N Darling[8], Joseph Geer[9], Zoran Gucev[10], V Reid Sutton[11], Christos Tziotzios[12], Adrian K Dixon[13], Timothy Helliwell[14], Stephen O'Rahilly[2,15], David B Savage[2,15], Michael J O Wakelam[4], Inès Barroso[2,3], Leslie G Biesecker[1] & Robert K Semple[2,15]

VOLUME 44 | NUMBER 8 | AUGUST 2012 NATURE GENETICS

## Identifying genetic drivers in a specific cancer

Somatic mutations in *ATP1A1* and *CACNA1D* underlie a common subtype of adrenal hypertension

Elena A B Azizan[1,12], Hanne Poulsen[2,12], Petronel Tuluc[3,12], Junhua Zhou[1,12], Michael V Clausen[2], Andreas Lieb[3], Carmela Maniero[1], Sumedha Garg[4], Elena G Bochukova[4], Wanfeng Zhao[5], Lalarukh Haris Shaikh[1], Cheryl A Brighton[1], Ada E D Teo[1], Anthony P Davenport[1], Tanja Dekkers[6], Bas Tops[7], Benno Küsters[7], Jiri Ceral[8], Giles S H Yeo[4], Sudeshna Guha Neogi[4], Ian McFarlane[4], Nitzan Rosenfeld[10], Francesco Marass[10], James Hadfield[9], Wojciech Margas[11], Kanchan Chaggar[11], Miroslav Solar[8], Jaap Deinum[6], Annette C Dolphin[11], I Sadaf Farooqi[4,12], Joerg Striessnig[3,12], Poul Nissen[2,12] & Morris J Brown[1,12]

**Identifying the genetic driver in a specific cancer**

- Laser capture microdissection
- Compare the adenoma with the wild-type tissue *via* next generation sequencing of whole exomes

22

## Identifying genetic drivers in a specific cancer

Somatic mutations in *ATP1A1* and *CACNA1D* underlie a common subtype of adrenal hypertension

Elena A B Azizan[1,12], Hanne Poulsen[2,12], Petronel Tuluc[3,12], Junhua Zhou[1,12], Michael V Clausen[2], Andreas Lieb[3], Carmela Maniero[1], Sumedha Garg[4], Elena G Bochukova[4], Wanfeng Zhao[5], Lalarukh Haris Shaikh[1], Cheryl A Brighton[1], Ada E D Teo[1], Anthony P Davenport[1], Tanja Dekkers[6], Bas Tops[7], Benno Küsters[7], Jiri Ceral[8], Giles S H Yeo[4], Sudeshna Guha Neogi[4], Ian McFarlane[9], Nitzan Rosenfeld[10], Francesco Marass[10], James Hadfield[9], Wojciech Margas[11], Kanchan Chaggar[11], Miroslav Solar[8], Jaap Deinum[6], Annette C Dolphin[11], I Sadaf Farooqi[4,12], Joerg Striessnig[3,12], Poul Nissen[2,12] & Morris J Brown[1,12]

**Identifying the genetic driver in a specific cancer**

**Somatic mutation driven disease**

Dr. Giles Yeo – University of Cambridge, UK

## Prenatal diagnosis

Noninvasive prenatal diagnosis of fetal chromosomal aneuploidy by massively parallel genomic sequencing of DNA in maternal plasma

Rossa W. K. Chiu[a,b], K. C. Allen Chan[a,b], Yuan Gao[c,d], Virginia Y. M. Lau[a,b], Wenli Zheng[a,b], Tak Y. Leung[e], Chris H. F. Foo[f], Bin Xie[c], Nancy B. Y. Tsui[a,b], Fiona M. F. Lun[a,b], Benny C. Y. Zee[f], Tze K. Lau[e], Charles R. Cantor[g,1], and Y. M. Dennis Lo[a,b,1]

20458–20463 | PNAS | December 23, 2008 | vol. 105 | no. 51

| In the past, to identify if a fetus had a mutation in a specific gene | Amniocentesis |
| | High chance of miscarriage |
| | Sampling where the fetus is |

23

## Prenatal diagnosis

Noninvasive prenatal diagnosis of fetal chromosomal aneuploidy by massively parallel genomic sequencing of DNA in maternal plasma

Rossa W. K. Chiu[a,b], K. C. Allen Chan[a,b], Yuan Gao[c,d], Virginia Y. M. Lau[a,b], Wenli Zheng[a,b], Tak Y. Leung[e], Chris H. F. Foo[f], Bin Xie[c], Nancy B. Y. Tsui[a,b], Fiona M. F. Lun[a,b], Benny C. Y. Zee[f], Tze K. Lau[e], Charles R. Cantor[g,1], and Y. M. Dennis Lo[a,b,1]

20458–20463 | PNAS | December 23, 2008 | vol. 105 | no. 51

| The fetus releases DNA, which can be found in the mother's blood stream | Therefore, using high-throughput sequencing... |
| | ... we can carry out non-invasive prenatal diagnosis |

## Prenatal diagnosis

Noninvasive prenatal diagnosis of fetal chromosomal aneuploidy by massively parallel genomic sequencing of DNA in maternal plasma

ONLY 8 million reads per case are required to get enough sequence reads to achieve this!

Dr. Giles Yeo – University of Cambridge, UK

---

**Identifying aneuploidy in chromosome 21**

S4

4 rows per 'flow cell'

2,500,000,000
'polonies' per row

It means that 312 samples
can be multi-plexed per row!

24

---

**Identifying aneuploidy in chromosome 21**

**Aneuploidy in chromosome 21**

Chromosome

21                    X

z-score

euploid male fetus          T21 male fetus
euploid female fetus        T21 female fetus

Case no.          Case no.

■ Are there any additional
copies of chromosome 21?

■ Are there discrepancies in
the no. of chromosomes
that are present?

---

**Identifying aneuploidy in chromosome 21**

**Aneuploidy in chromosome 21**

Chromosome

21                    X

z-score

euploid male fetus          T21 male fetus
euploid female fetus        T21 female fetus

Case no.          Case no.

There is an overrepresentation
of chromosome 21

Can identify fetuses with
aneuploidy in chromosome 21

---

The screen versions of these slides have full details of copyright and acknowledgements

Dr. Giles Yeo – University of Cambridge, UK

**RNA sequencing (RNAseq)**

GENETICS ➡ TRANSCRIPTOME

25

**RNA sequencing (RNAseq)**

RNA extracted,
converted to cDNA,
& sequenced

Sequences mapped onto genome

Condition A

Upregulation
of transcript

Condition B

with drug

**RNA sequencing (RNAseq)**

Sequences mapped onto genome

There are more copies
of the transcript
in condition B...

↓

...therefore, the drug
likely upregulated
this transcript

Upregulation
of transcript

The screen versions of these slides have full details of copyright and acknowledgements

Dr. Giles Yeo – University of Cambridge, UK

Dr. Giles Yeo – University of Cambridge, UK

**Summary**

**THANK YOU VERY MUCH**