

Introduction to second-generation sequencing

CMSC858B Spring 2012

Many slides courtesy of Ben Langmead

SECOND-GENERATION SEQUENCING

nature news

[nature news home](#) [news archive](#) [specials](#) [opinion](#) [features](#) [news blog](#) [events](#)

Access

This article is part of Nature's premium content.

Published online 15 October 2008 | *Nature* **455**, 847 (2008) | doi:10.1038/455847a

News

The death of microarrays?

High-throughput gene sequencing seems to be stealing a march on microarrays. Heidi Ledford looks at a genome technology facing intense competition.

Heidi Ledford

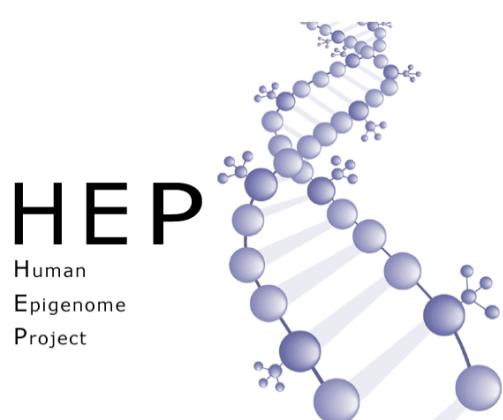
Faster, cheaper DNA sequencing technology is revolutionizing the burgeoning field of personal genomics. But it is having another, more subtle effect.

Tools

[Send to a Friend](#)



HUMAN EPIGENOME PROJECT



HEP
Human
Epigenome
Project

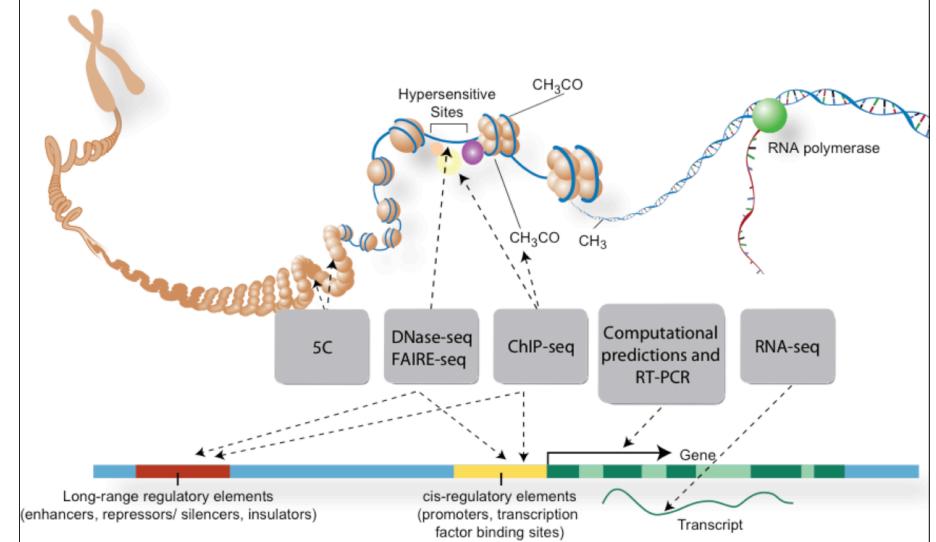


JOHNS HOPKINS
BLOOMBERG
SCHOOL OF PUBLIC HEALTH

3 Corrada Bravo 10/30/09

3

ENCODE project



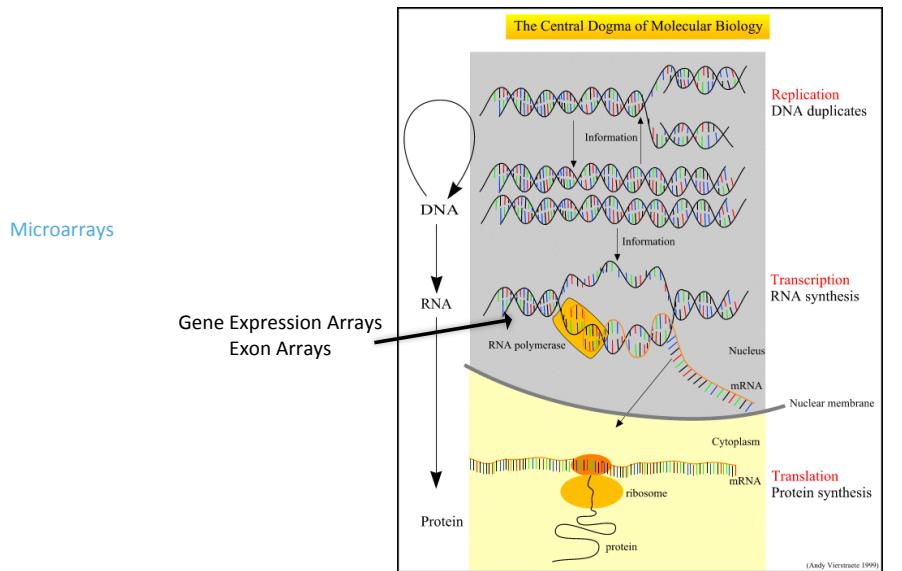
4

1000 GENOMES PROJECT

The screenshot shows the "1000 GENOMES PROJECT DATA RELEASE" section. It includes a heading "SNP data downloads and genome browser representing four high coverage individuals", a paragraph about the availability of SNP calls through FTP sites, and links to the EBI and NCBI FTP sites. It also mentions the "Data section" of the website and provides a "Quick start (pdf)" guide.

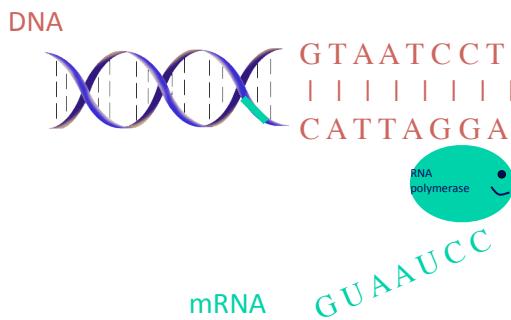
5

What Do They Measure?



6

Transcription

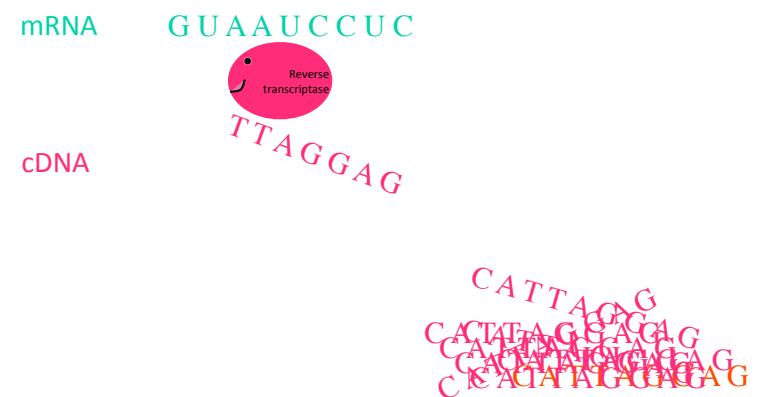


From DNA to mRNA

7

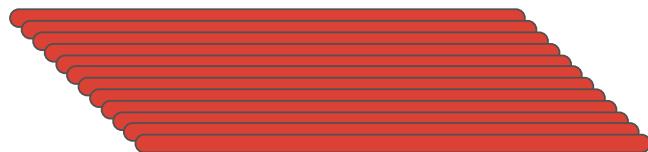
Reverse transcription

Clone cDNA strands, complementary to the mRNA



8

SEC-GEN SEQUENCING

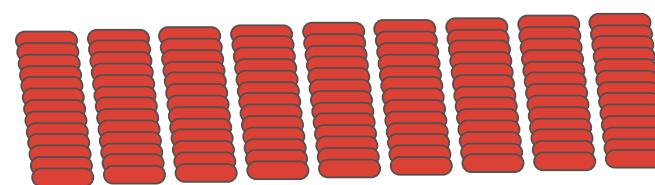


9 Corrada Bravo 10/30/09



9

SEC-GEN SEQUENCING



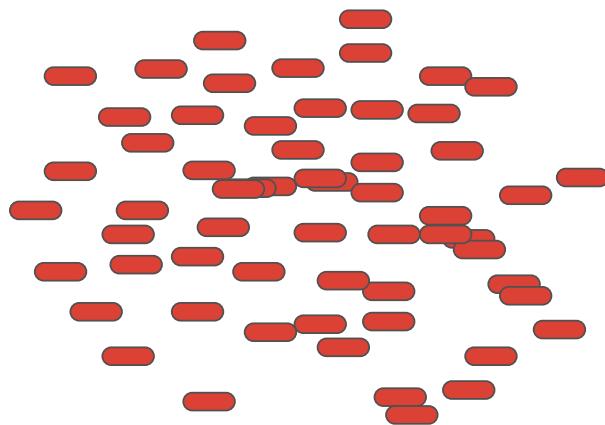
Fragmentation is random,
i.e., not equal-sized (but hard to draw)

10 Corrada Bravo 10/30/09



10

SEC-GEN SEQUENCING



11 Corrada Bravo 10/30/09



11

SECOND-GENERATION SEQUENCING

- “Ultra high throughput” DNA sequencing
- 3 gigabases / day vs.
- 3 gigabases / 13 years (human genome project, more or less)

12 Corrada Bravo 10/30/09



12

PLATFORMS



Semiconductor Sequencing for Life™

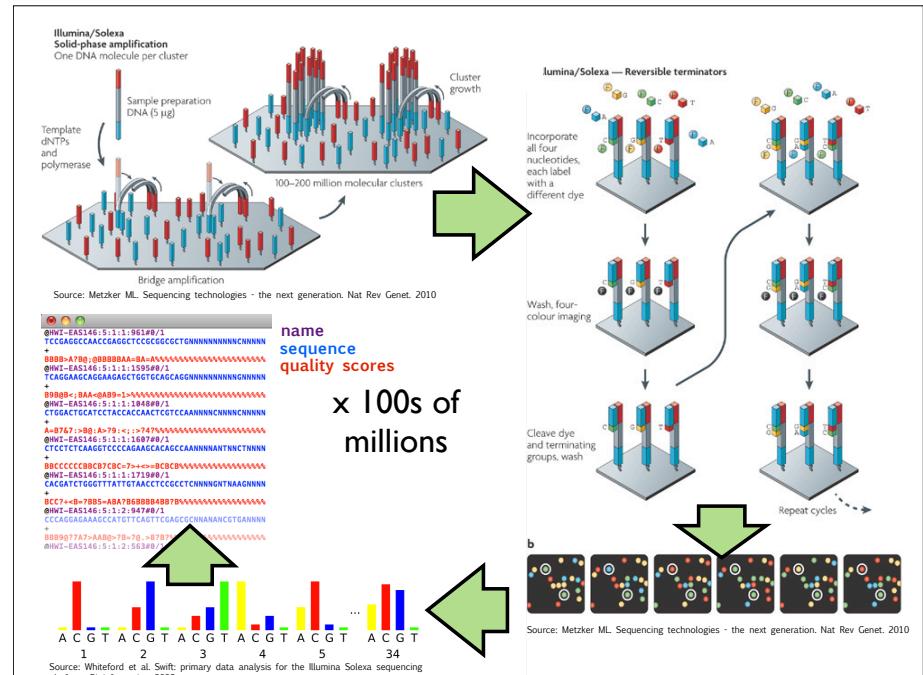
ion torrent

454
SEQUENCING



- Millions of short DNA fragments (~100 bp) sequenced in parallel

13



14

Sequencing throughput



GA II
1.6 billion bp per day
(2008)



GA IIx
5 billion bp per day
(2009)



HiSeq 2000
25 billion bp per day
(2010)

Images: www.illumina.com/systems
Numbers: www.politigenomics.com/next-generation-sequencing-informatics
Dates: Illumina press releases

15

Sequencing throughput



GA II
1.6 billion bp per day
(2008)



GA IIx
5 billion bp per day
(2009)



HiSeq 2500
60 billion bp per day
(2012)

Images: www.illumina.com/systems
Numbers: www.politigenomics.com/next-generation-sequencing-informatics
Dates: Illumina press releases

16

Sequencing throughput

End of 2009

Mid 2010

Late 2011/Early 2012

New SOLID™ Systems		
SOLID 3+ System	SOLID™ 4 System	SOLID™ 4hg System*
Up to 50 Gb	Up to 100 Gb	Up to 300 Gb
System Accuracy	99.94%	99.99%
Cost/Genome***	As low as \$6,000	As low as \$3,000
Read Length	• Fragment: 50 bp • Mate-pair: 2 x 50 bp • Paired-end: 50 x 25 bp	• Fragment: 75 bp • Mate-pair: 2 x 75 bp • Paired-end: 75 x 35 bp
Multiplexing	• 48 RNA barcodes • 96 DNA barcodes	• 96 RNA barcodes • 96 DNA barcodes
Run Time	• 3 days for 35 bp • 11 days for 50 x 25 bp • 12 days for 2 x 50 bp	• 3 days for 35 bp • 12 days for 75 x 35 bp • 14 days for 2 x 75 bp

* These systems are under development and the specifications are subject to change.
** Expected throughput.
*** Prices reflect US list at optimal running efficiencies. Regional pricing may vary.

Source: www3.appliedbiosystems.com/cms/groups/mcb_marketing/documents/generaldocuments/cms_061241.pdf

17

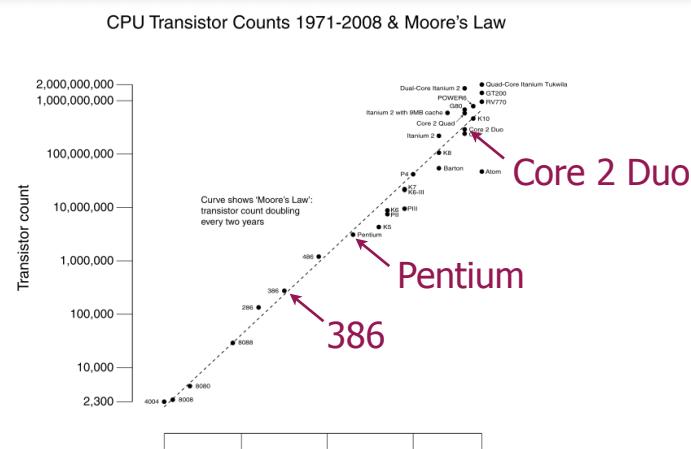
Computational throughput

Moore's Law:

The number of transistors that can be placed inexpensively on an integrated circuit doubles approximately every two years.

18

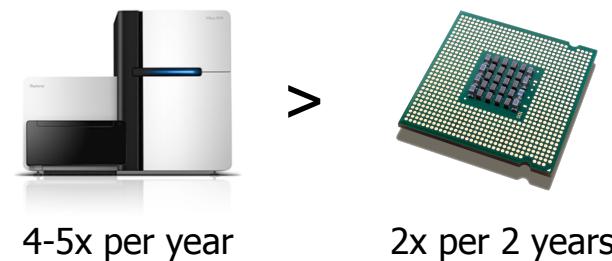
Computational throughput



Source: en.wikipedia.org/wiki/Moore%27s_law

19

Throughput growth gap



20

ionTorrent

Ion Semiconductor Sequencing



Speed, Scalability, Simplicity

Ion semiconductor benchtop sequencers perform real-time measurements of hydrogen ions produced during DNA replication. Ion semiconductor chips employ a massively parallel array of semiconductor sensors, to directly translate genetic information (DNA) to digital information (DNA sequence). This groundbreaking technology enables rapid and scalable sequencing across a range of applications.

Technology



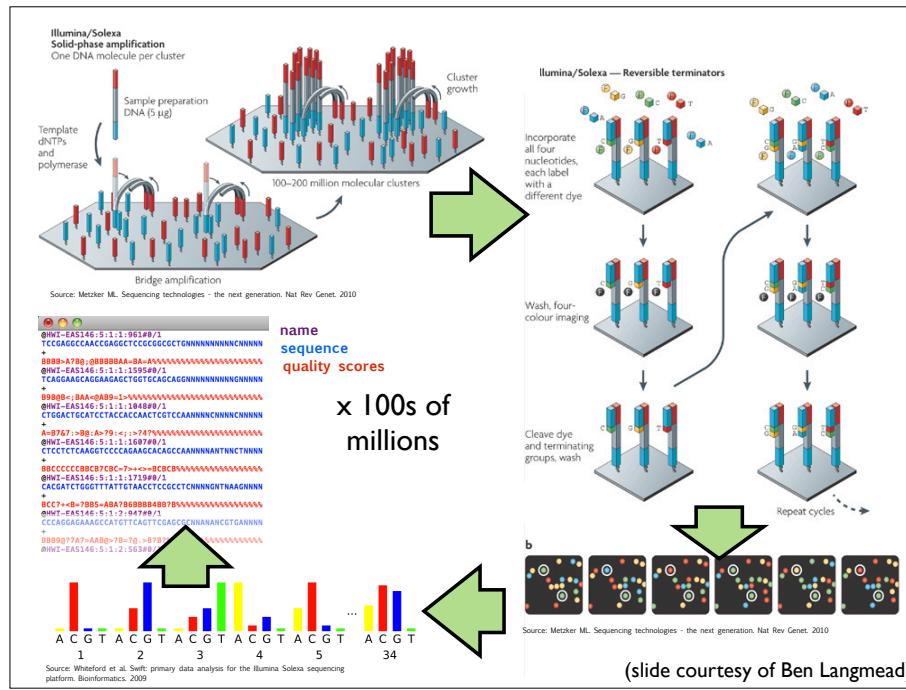
21

Oxford Nanopore

- Nanopore technology
- ultralong reads (48kb genome sequenced as one read)

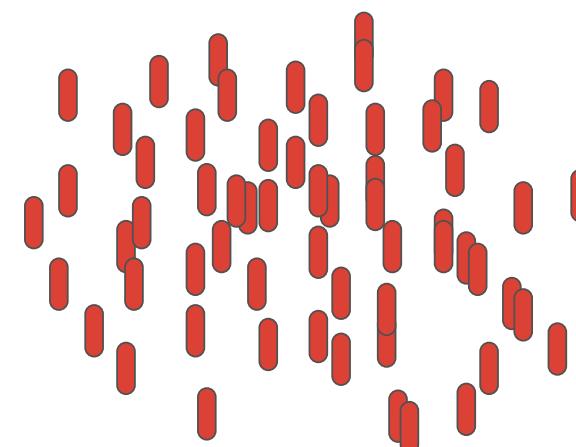


22

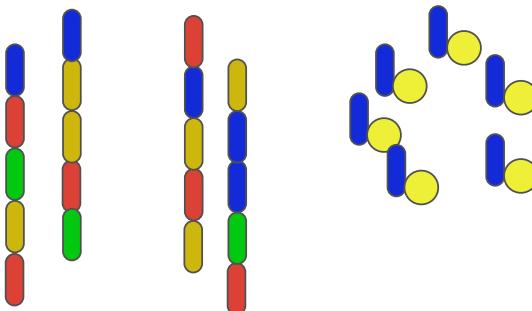


23

SEC-GEN SEQUENCING



SEC-GEN SEQUENCING

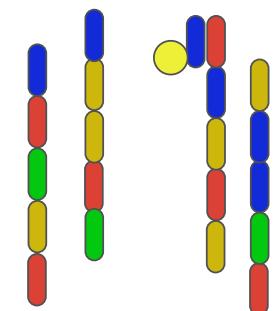


25 Corrada Bravo 10/30/09



25

SEC-GEN SEQUENCING

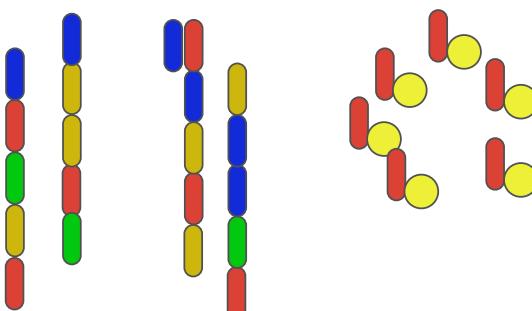


26 Corrada Bravo 10/30/09



26

SEC-GEN SEQUENCING

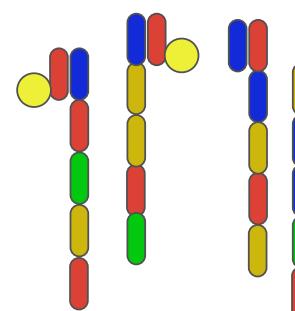


27 Corrada Bravo 10/30/09



27

SEC-GEN SEQUENCING

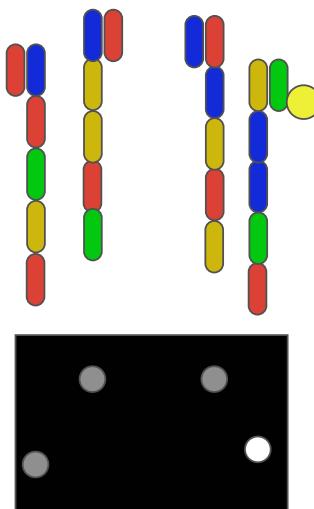


28 Corrada Bravo 10/30/09



28

SEC-GEN SEQUENCING

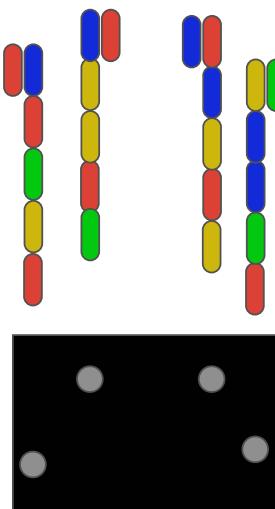


29 Corrada Bravo 10/30/09

29



SEC-GEN SEQUENCING

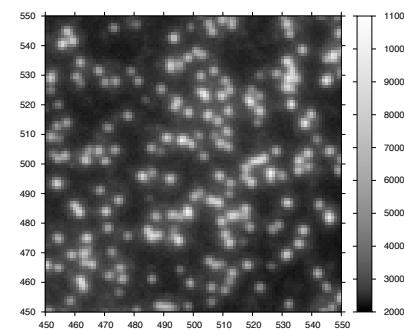
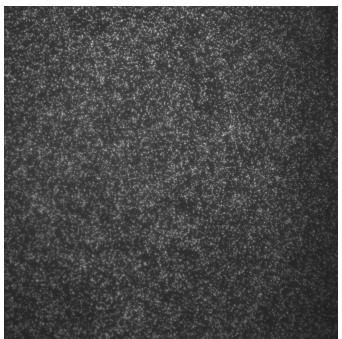


30 Corrada Bravo 10/30/09

30



Image Analysis



An input image and zoomed in section

31

Image Analysis

- 4 images per cycle
- ~100 tiles
- Analysis:
 - Filtering
 - Background subtraction
 - Thresholding
- Each image analysis independent (so can parallelize)

32

Image Analysis

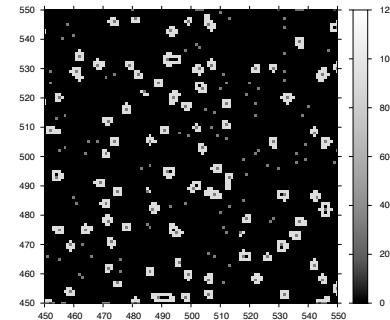
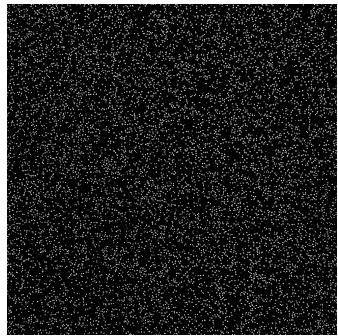
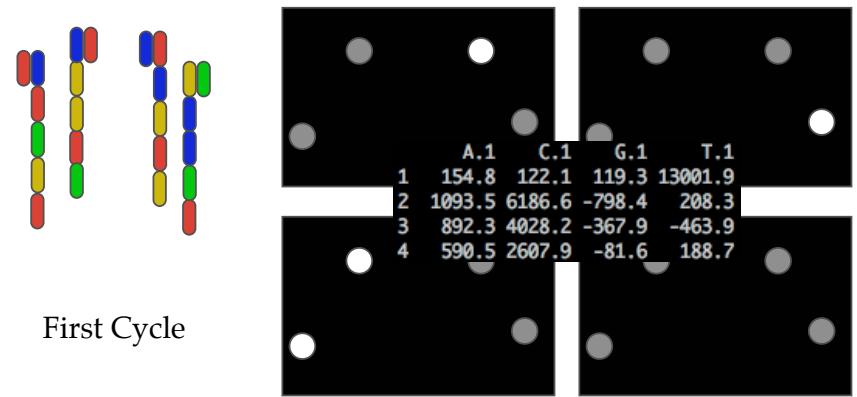


Image after processing. This is old,
cluster density is much higher now

33

SEC-GEN SEQUENCING



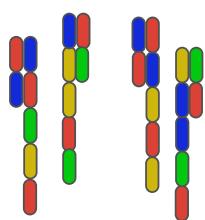
First Cycle

34 Corrada Bravo 10/30/09

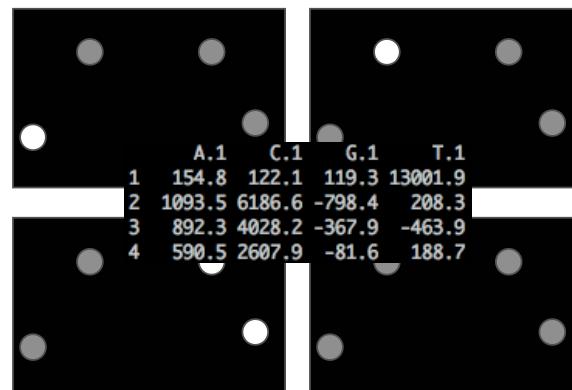


34

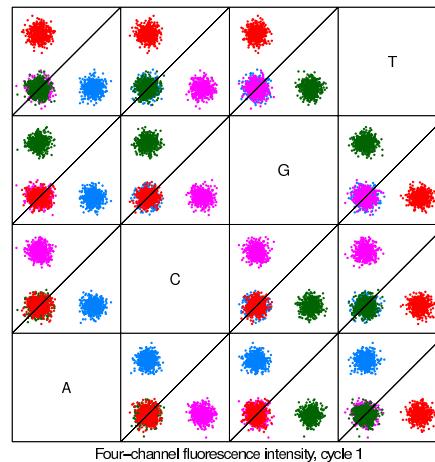
SEC-GEN SEQUENCING



Second Cycle



A THOUGHT EXPERIMENT



Four-channel fluorescence intensity, cycle 1

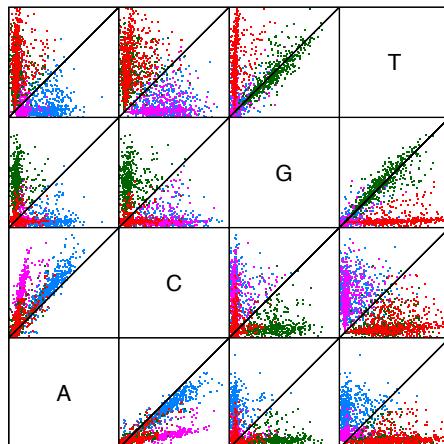
Color coded by call
made: A, C, G, T

35 Corrada Bravo 10/30/09



35

FLUORESCENCE INTENSITY



Color coded by call made: A, C, G, T

Four-channel fluorescence intensity, cycle 1

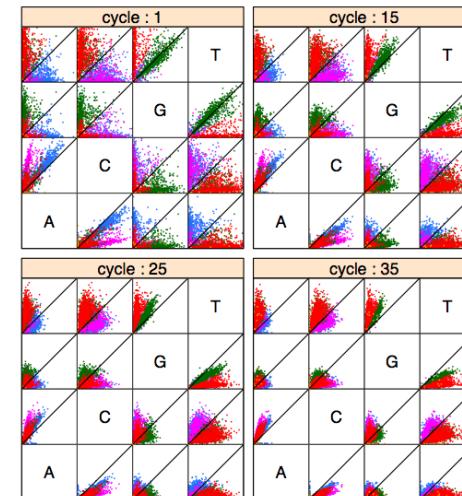
37 Corrada Bravo 10/30/09



37

FLUORESCENCE INTENSITY

MORE ON THIS LATER IN THE COURSE!

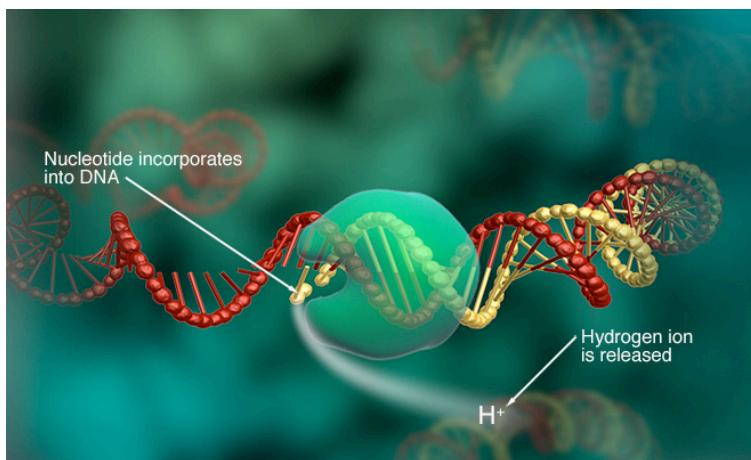


38 Corrada Bravo 10/30



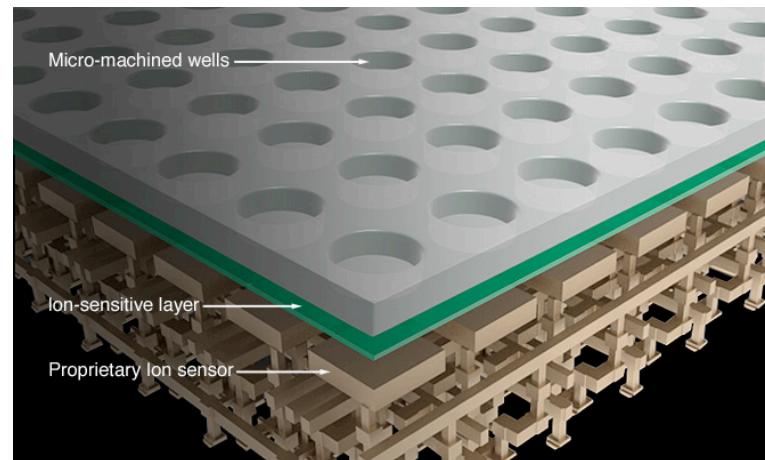
38

ION TORRENT



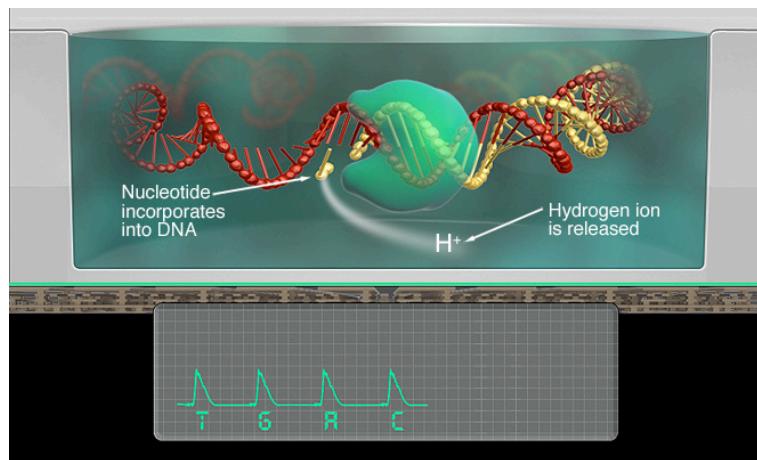
39

ION TORRENT



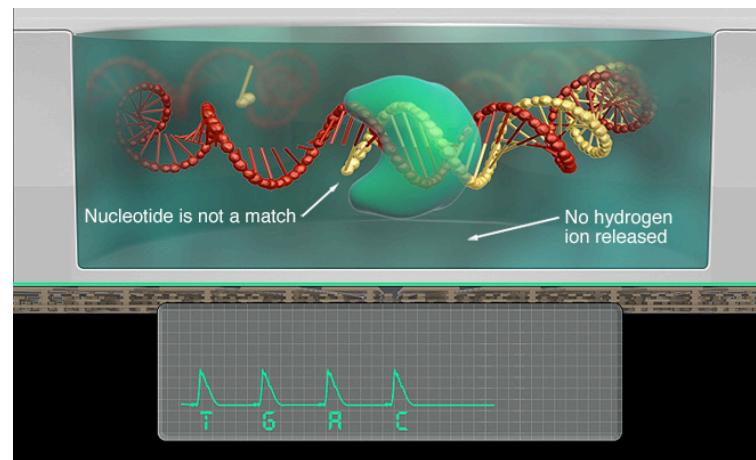
40

ION TORRENT



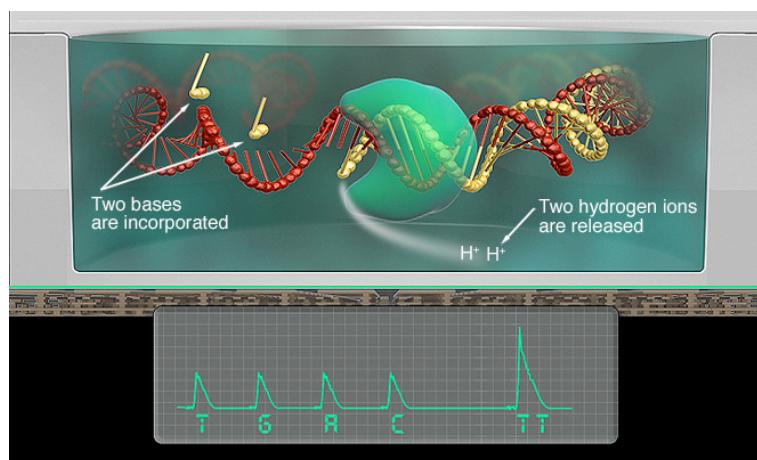
41

ION TORRENT

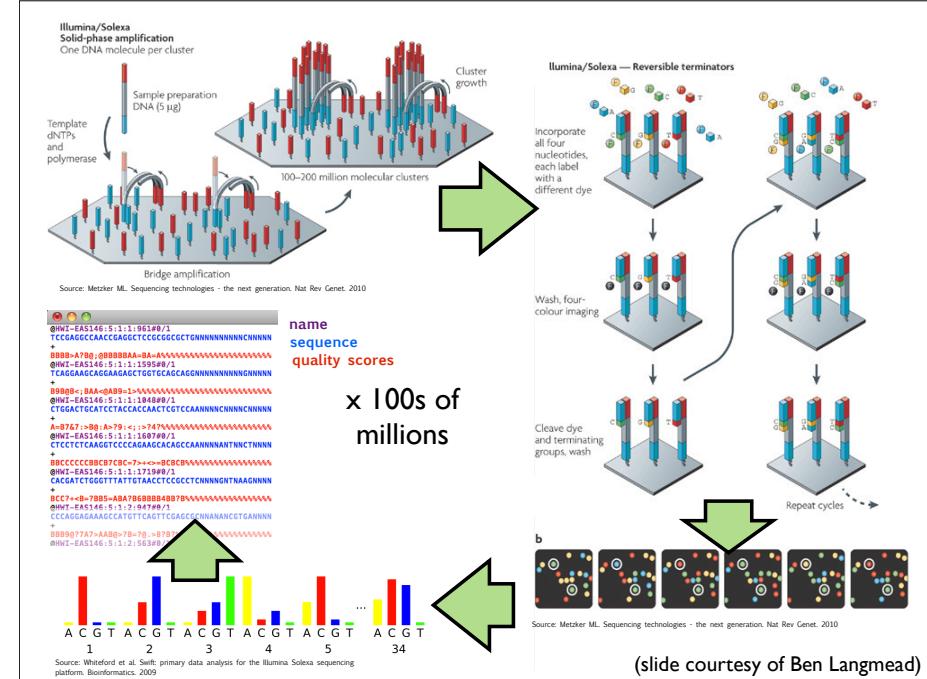


42

ION TORRENT



43



44

From reads to evidence

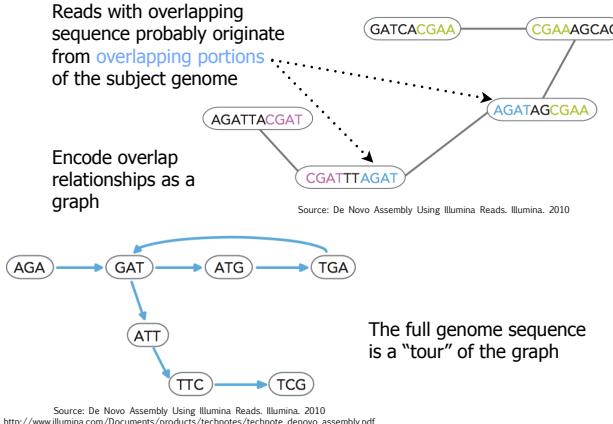
4

From reads to evidence

2. de novo

Assume nothing! - let reads tell us everything

Reads with overlapping sequence probably originate from overlapping portions of the subject genome



The full genome sequence
is a “tour” of the graph

Source: De Novo Assembly Using Illumina Reads. Illumina. 2010
http://www.illumina.com/Documents/products/technotes/technote_denovo_assembly.pdf

4

From reads to evidence

I. Comparative

Sequence-wise, individuals of a species are nearly identical

Well curated, annotated “reference” genomes exist



Idea: "Map" reads to their point of origin with respect to a reference, then study differences

4

Mapping

Take a read:

CTCAAACCTCCTGACCTTGGTATCCACCCGCCTNGGCCTT

How do we determine the read's point of origin with respect to the reference?

Answer: sequence similarity

Hypothesis 1

```

          Read
CTCAAAGACCTGACCTTTGGTATGCCACCC-----GCCTNGGCCCTC
||||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |
CTAAACTCTGGATTGG---GATCCACCCAGCTGGCCTGGCCATAA
          Reference

```

Hypothesis 2

Which hypothesis is better?

Say hypothesis 2 is correct. Why are there still mismatches and gaps?

4

Mapping

Take a read:

CTCAAACCTCTGACCTTGGTATCCA

And a reference sequence:

Which hypothesis is better?

```
> NT_dna:chr000508141>00050814:GRCh37:MT-1-15689-1
GATCACAGCTTATACCTCTTAACTCACTTACAGCGATCTCCCTTGATTTGGTATTTC
CTCTGGGGGTATGACGGGGATAGGCTTGGGAGGCCCTGGGAGGCCCTGGT
GCAGATCTGTCTTTGTTCTGGCTCATCTTATTTATTCGACACTAGGTCATAATT
ACAGGGACACATACTAATCAGGTGTAAATTAACTGCTTAGGACATAAATAA
ACAATTGATGTGCAAGCGCACCTTCACACAGACATCATACAAAAAACTTCCCA
AAACCCCCCTGGGCAACAGCTAACACATCTGGCAAAACCCCCAAA
ACAAGAACCTTACACAGCGCTAACAGATTTCATTTGGCGGTATCAC
TTTAAACAGTACCCCCCTAACACATTTTCCCTAACACTCCATCTACTAAAT
CTCTCATACACCCCCGGCCATCCCTGGGACACACACACACCGCTGCTAACCCATA
CCCTACACACACACACACACACACACACACACACACACACACACACACACAC
GAGATACAGTACCCGGTAAACTCTGGGATTTCTGGCTCATCTGGCTAAC
CTAGCTCTCTTCTTATGCTTGTAAATTACATCGACAGCATCCGGTCTCAGTGACT
TCACCCCTCTAAATACCCAGATCAAAAGAACAGCATCAAGCAGCAGAAATGAGCTC
AAAAGCTTAGCTGACACCCCCGGAAACAGCTGGGAGGAACTTACCTTAAATAA
ACGAAAGTTTAACTAGCTTAACTAACCCAGGGTTGTAATTCTGGCAGGACCCG
GGTCACAGATTACCAAGCTAATGAAGGGGGTAAAGAGTGTGTTAGATCACC
TCCCCTAAAGCTAAACATCACCTGAGTTGTAATAAAACTCAGTGCACAAAATAGC
TACCTGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGGCTGG
TACCCACACCTTGGCCAAACACTAACAGCAACATACACACACACACACACAC
CACTACAGCCACAGCTAAACTAACAGGACCTGGGGTCTGGCTGGCTGGCTGG
AGCTGTGTTGTAATTCTGGTAAACTCTGGGAGGAACTTACCTGGCTGGCTGGCTGG
CCGGCATCTTCAAGAACCTGTGATGAGGCTAACAGTAAAGGAGCTACCCAGTAAG
AGCTAGGGTCAAGGTGAGGGCATTAAGGGTCAATTCTGGGAGGAACTTACCTGG
AAAATCTAGACGACCTTAACTTGAAGGCTGAAAGTGTGAGGAGCTAACCTACAA
AGTATAGGGTGAAGAATTTGAAACTGGGCAATATAGATATGACGGAAAGGAGATG
AAAATTATAACCAAGCTAATATAGCAAGGACTAACCCCTATACCTTCTGCTATAA
GTTACTGAAATAACCTTGGCAAGGAGGACCAAGTGAACCCCCAACAGACGAGCT
```

Hypothesis 1:

Read
CTCAAACCTCTG**A**CTTTGGTATCCA
Reference
CTCAAACCTCTG-C**CT**TTGGTATCCA

Hypothesis 2:

Read
CTCAAACCTCTGAC**C**CTTTGGTATCCA
Reference
CTCAAACCTCTGAC**T**CTTGGTATCCA

Is there any way to break the tie?

Hint: In Illumina sequencing, sequencing errors almost never manifest as gaps

53

Mapping

Aligners can employ **penalties** to account for the relative probability of seeing different dissimilarities

Estimates vary, but small gaps ("indels") occur in humans at 1 in ~10-100K positions.

$$\begin{aligned} Pen_{\text{gap}} &\equiv -10 \log_{10}(P_{\text{gap}}) \\ &= -10 \log_{10}(0.00005) \\ &\approx 45 \end{aligned}$$

SNPs occur in humans at 1 in ~1K positions, but depending on Q, sequencing error may be more likely

$$\begin{aligned} Pen_{\text{mm}} &\equiv \arg \min(-10 \log_{10}(P_{\text{miscal}}), -10 \log_{10}(P_{\text{SNP}})) \\ &= \arg \min(Q, -10 \log_{10}(0.001)) \\ &= \arg \min(Q, 30) \end{aligned}$$

<p>Read CTCAAACCTCTGACTTTGGTATCCA Reference CTCAAACCTCTG-CCTTTGGTATCCA</p>	Penalty = 45
<p>Read CTCAAACTCTGACCCTTTGGTATCCA Reference CTCAA-CTCCTGACTCTTGGGTATCCA</p>	Q=10
<p>Read CTCAAACTCTGACCCTTTGGTATCCA Reference CTCAA-CTCCTGACTCTTGGGTATCCA</p>	Penalty = 55
<p>Read CTCAAACTCTGACCCTTTGGTATCCA Reference CTCAA-CTCCTGACTCTTGGGTATCCA</p>	Q=40

54

Resources

- Bowtie: ultra-fast mapping of short reads to reference genome

Bowtie
An ultrafast memory-efficient short read aligner



Bowtie is an ultrafast, memory-efficient short read aligner. It aligns short DNA sequences (reads) to the human genome at a million 35-bp reads per hour. Bowtie indexes the genome with a Burrows-Wheeler index to keep its memory footprint small: 1 GB for the human genome (2.9 GB for paired-end).

• <http://bowtie-bio.sourceforge.net>

