



CMSC 858D :

Details of Sequencing Technologies: short & long-read sequencing

NOTE: Illumina sequence slides are taken from
<http://www.slideshare.net/USDBioinformatics/illumina-sequencing>



CMSC 858D :

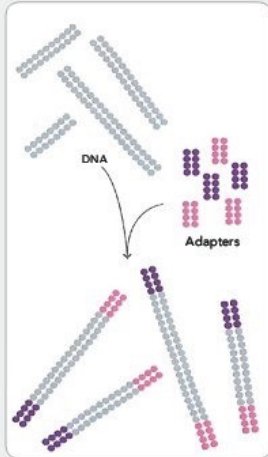
Details of Sequencing Technologies: short & long-read sequencing

NOTE: Illumina sequence slides are taken from
<http://www.slideshare.net/USDBioinformatics/illumina-sequencing>

Illumina Diagram

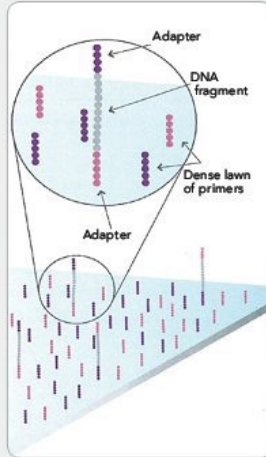


1. PREPARE GENOMIC DNA SAMPLE



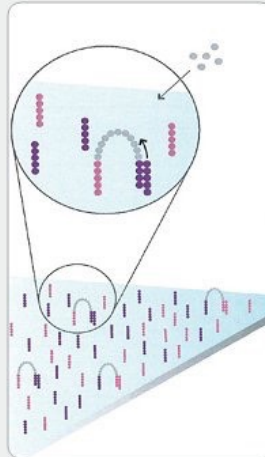
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



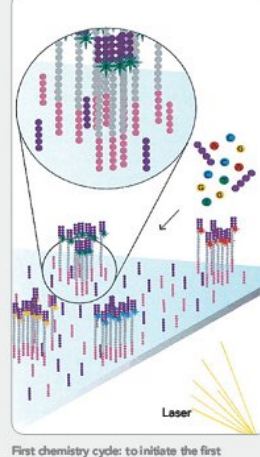
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



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

7. DETERMINE FIRST BASE



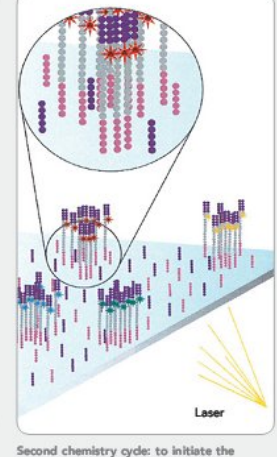
First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

8. IMAGE FIRST BASE



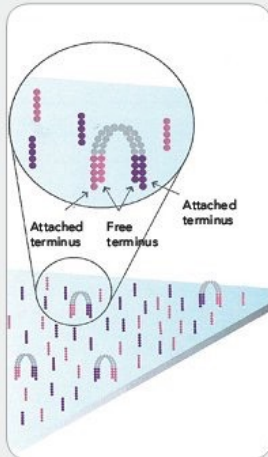
After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

9. DETERMINE SECOND BASE



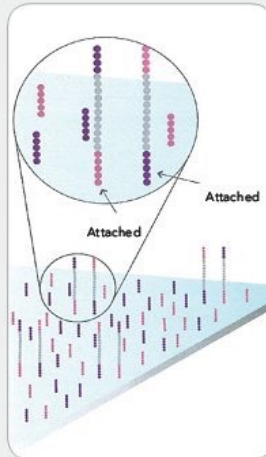
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

4. FRAGMENTS BECOME DOUBLE STRANDED



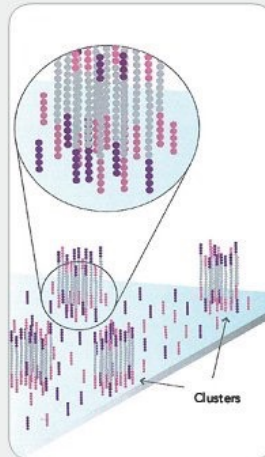
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



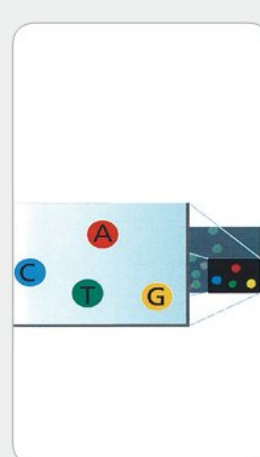
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



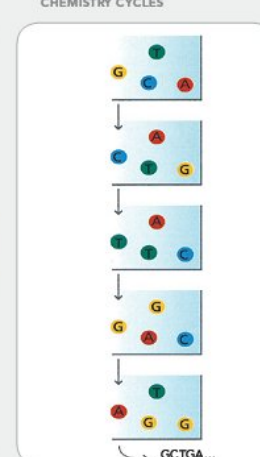
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

10. IMAGE SECOND CHEMISTRY CYCLE



After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA

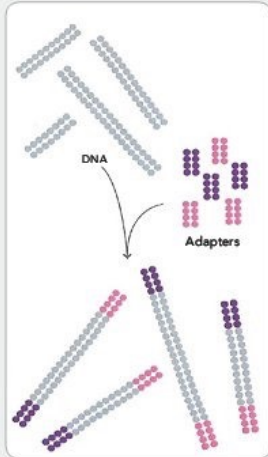


Align data, compare to a reference, and identify sequence differences.

Illumina Diagram

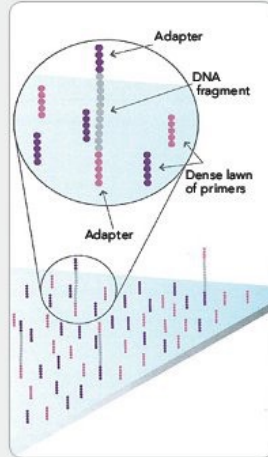


1. PREPARE GENOMIC DNA SAMPLE



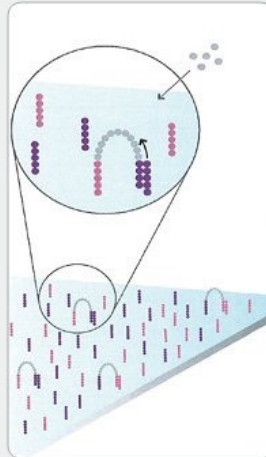
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



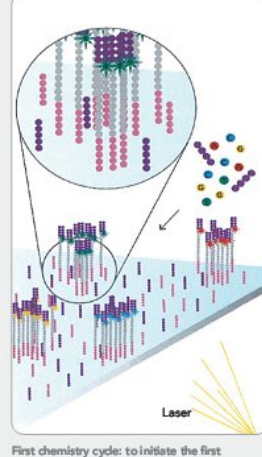
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



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

7. DETERMINE FIRST BASE



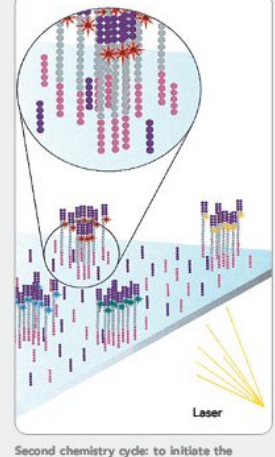
First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

8. IMAGE FIRST BASE



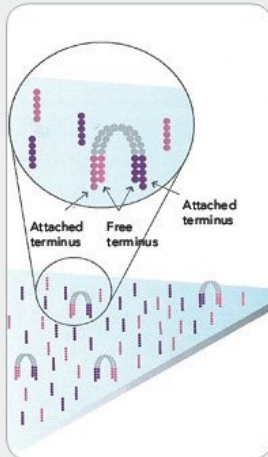
After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

9. DETERMINE SECOND BASE



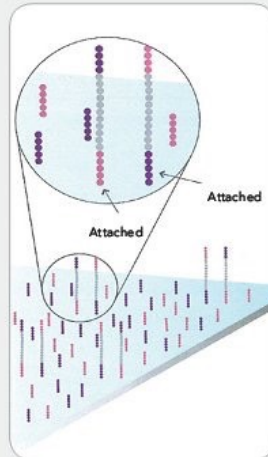
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

4. FRAGMENTS BECOME DOUBLE STRANDED



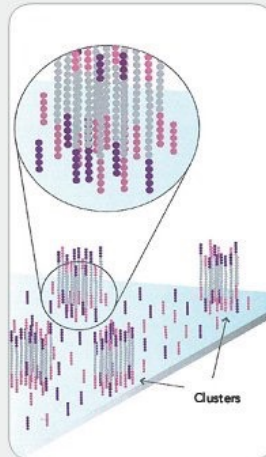
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



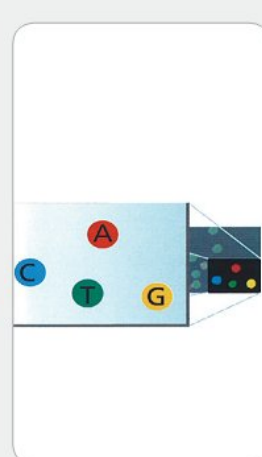
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



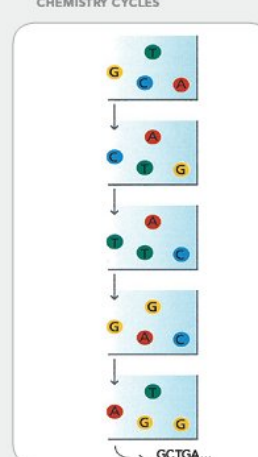
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

10. IMAGE SECOND CHEMISTRY CYCLE



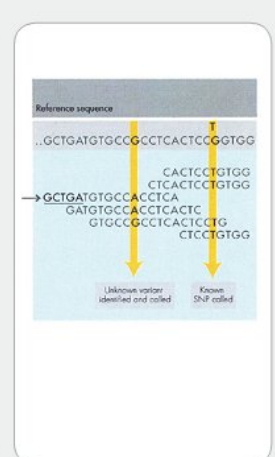
After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



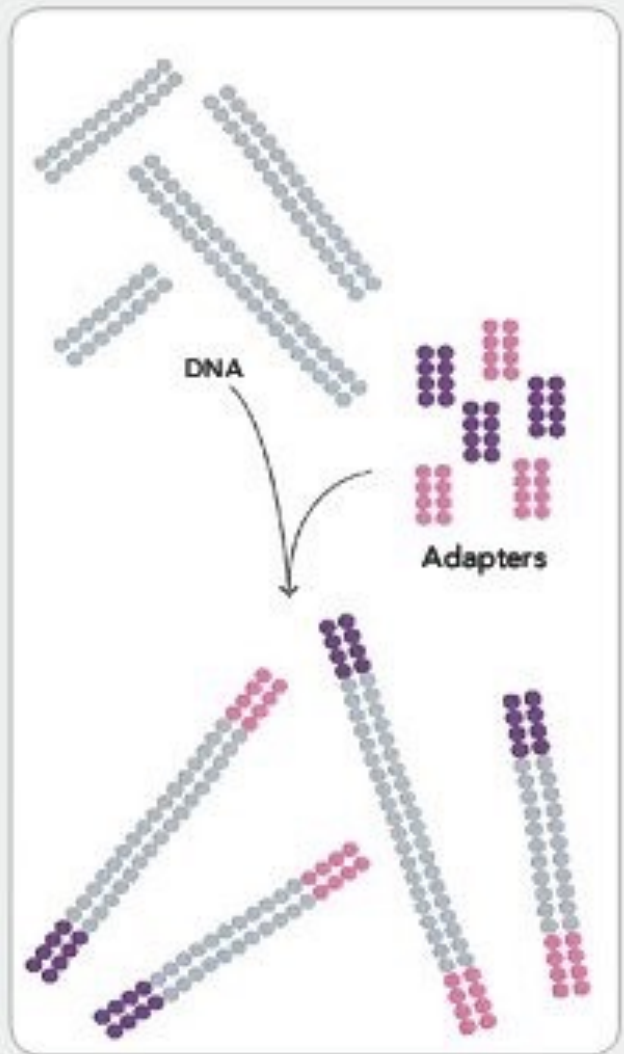
Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA



Align data, compare to a reference, and identify sequence differences.

1. PREPARE GENOMIC DNA SAMPLE



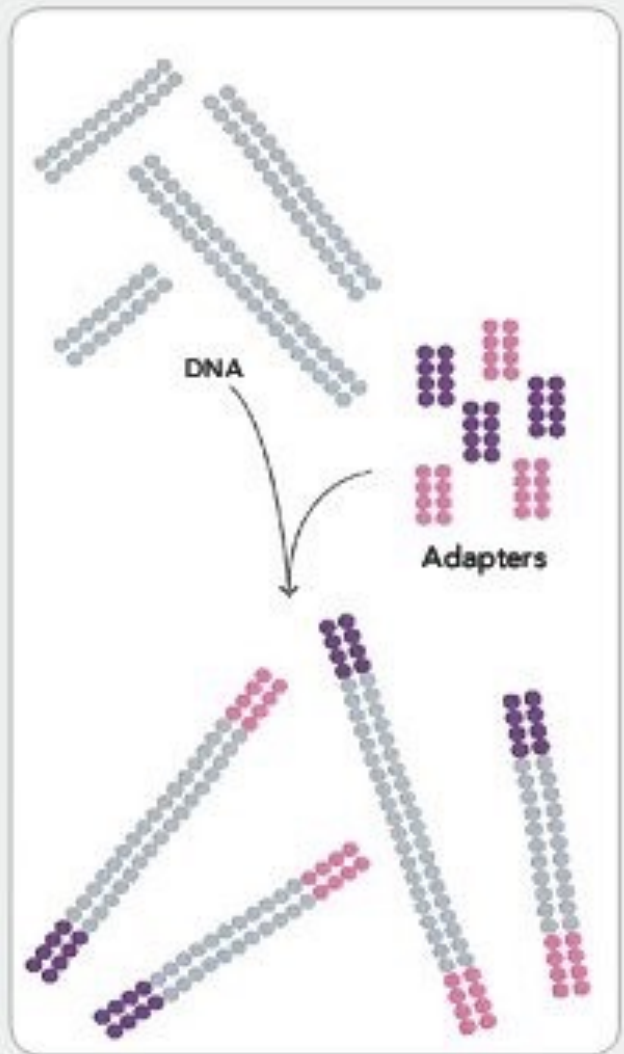
Prepare Genomic DNA Sample



- Fragment DNA of interest into smaller strands that are able to be sequenced
 - Sonication
 - Nebulization
 - Enzyme digestion
- Ligate Adapters
- Denature dsDNA into ssDNA by heating to 95° C

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

1. PREPARE GENOMIC DNA SAMPLE



Prepare Genomic DNA Sample

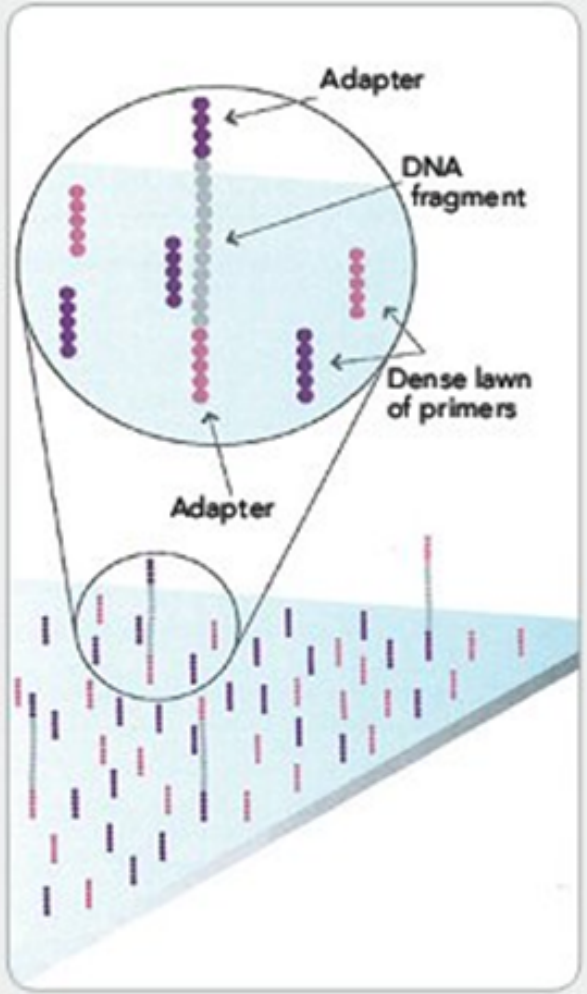


- Fragment DNA of interest into smaller strands that are able to be sequenced
 - Sonication
 - Nebulization
 - Enzyme digestion
- Ligate Adapters
- Denature dsDNA into ssDNA by heating to 95° C

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



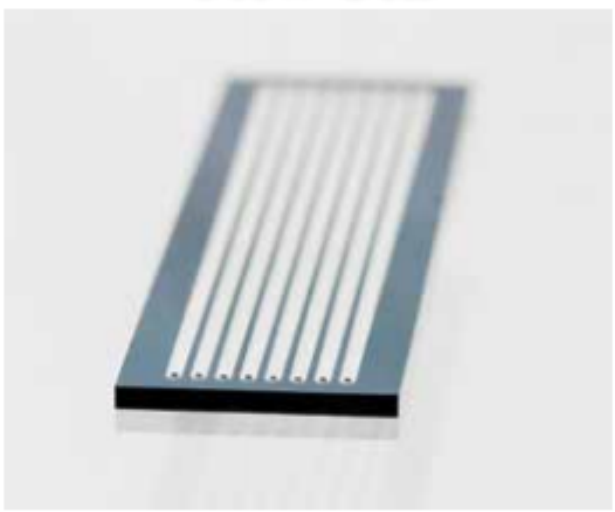
Attach DNA to Surface



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

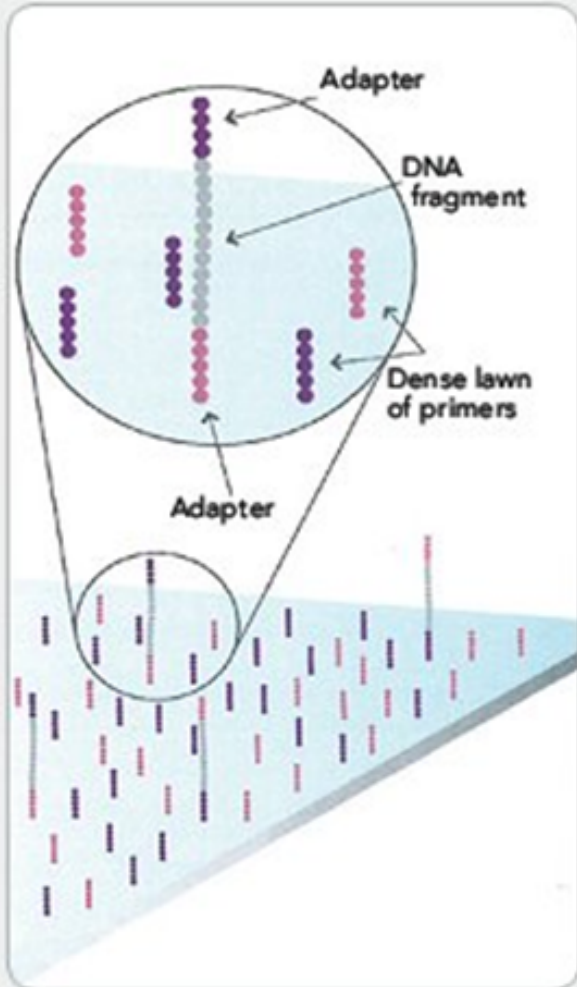
- ssDNA is then bound to inside surface of flow cell channels
- Dense lawn of primer on the surface of the flow cell

Flow Cell





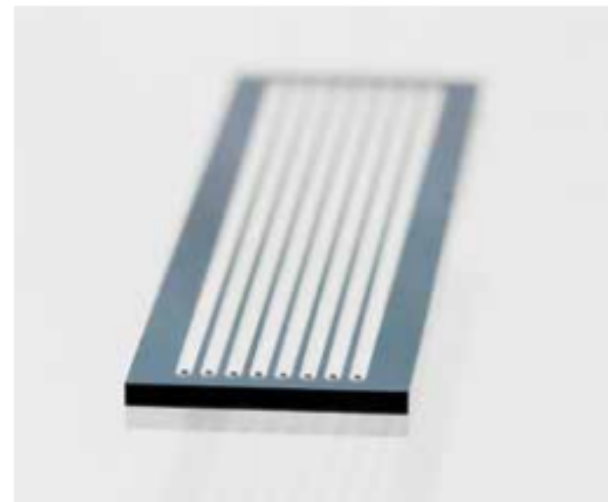
Attach DNA to Surface

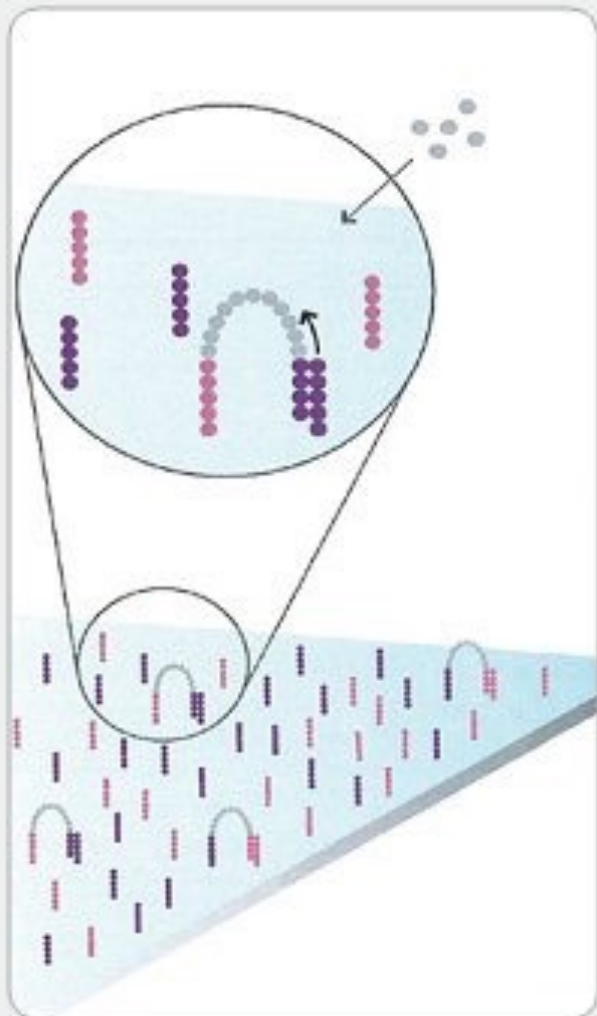


Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

- ssDNA is then bound to inside surface of flow cell channels
- Dense lawn of primer on the surface of the flow cell

Flow Cell



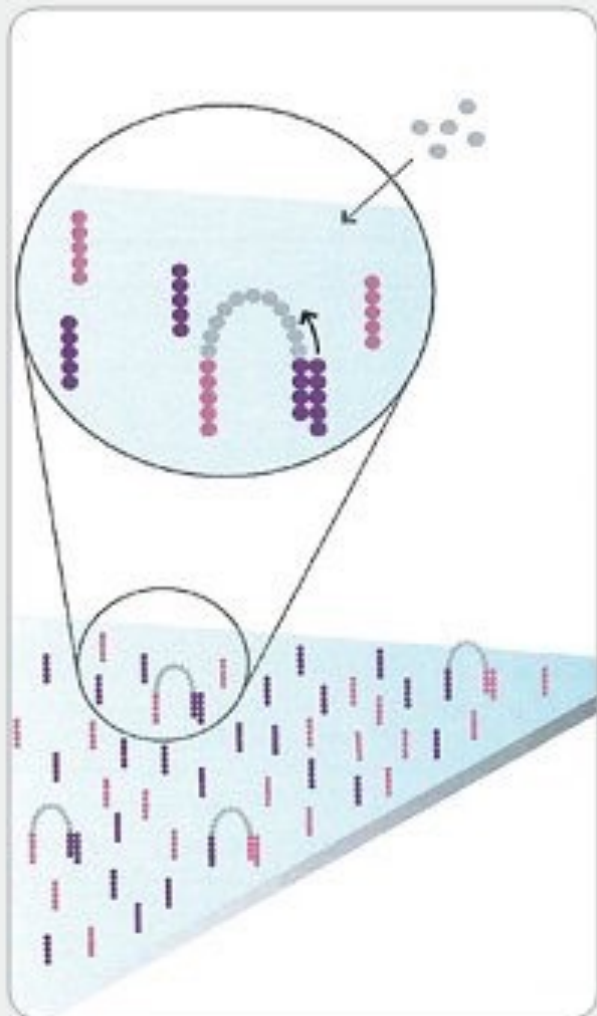


Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

Bridge Amplification



- Unlabeled nucleotides and polymerase enzyme are added to initiate the solid phase bridge amplification



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

Bridge Amplification



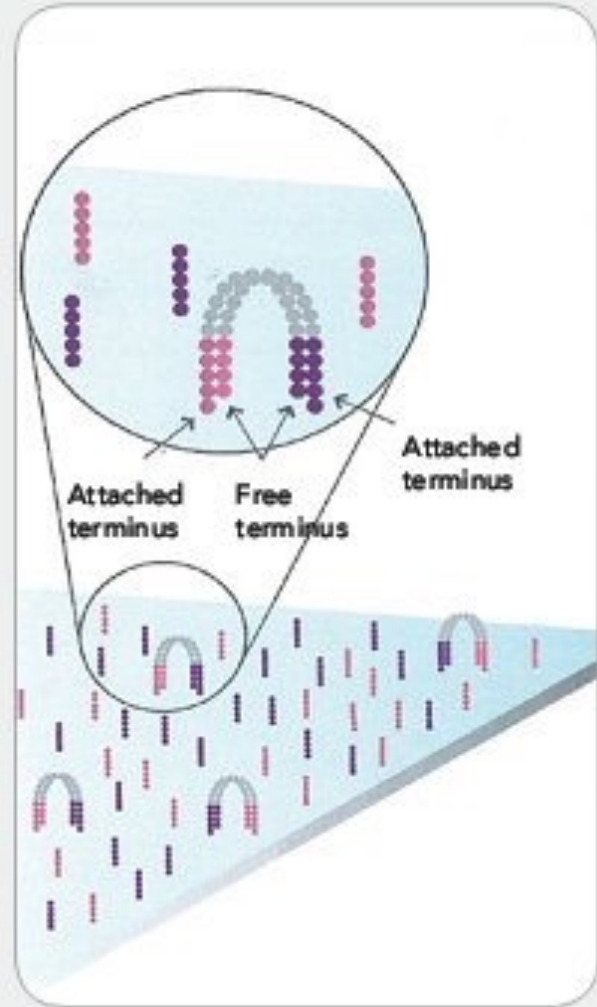
- Unlabeled nucleotides and polymerase enzyme are added to initiate the solid phase bridge amplification

Fragments Become Double Stranded



- In this step it demonstrates the work done by the sequencing reagents
 - Primers
 - Nucleotides
 - Polymerase enzymes
 - Buffer

4. FRAGMENTS BECOME DOUBLE STRANDED



The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

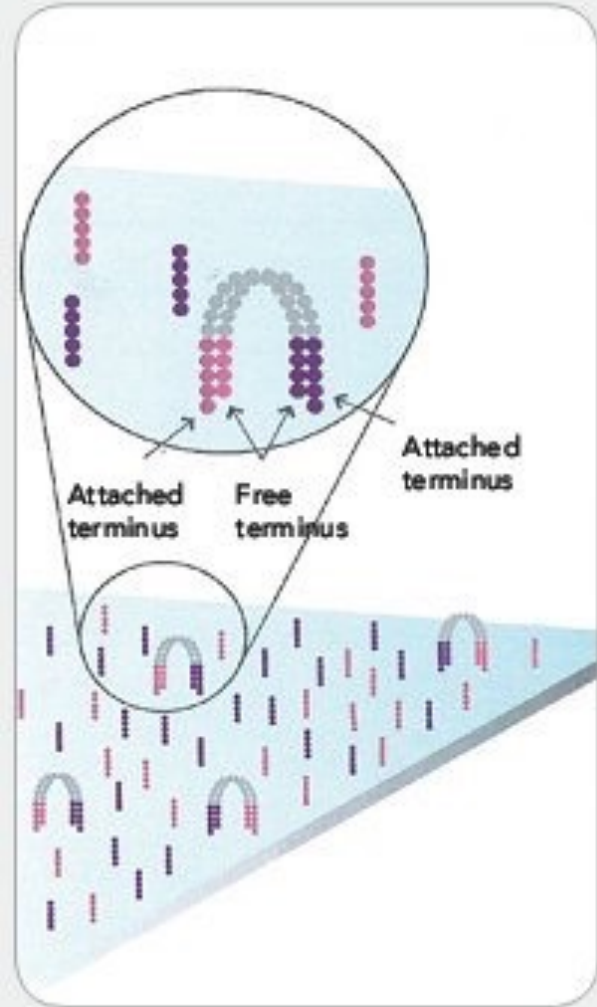
Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Fragments Become Double Stranded



- In this step it demonstrates the work done by the sequencing reagents
 - Primers
 - Nucleotides
 - Polymerase enzymes
 - Buffer

4. FRAGMENTS BECOME DOUBLE STRANDED

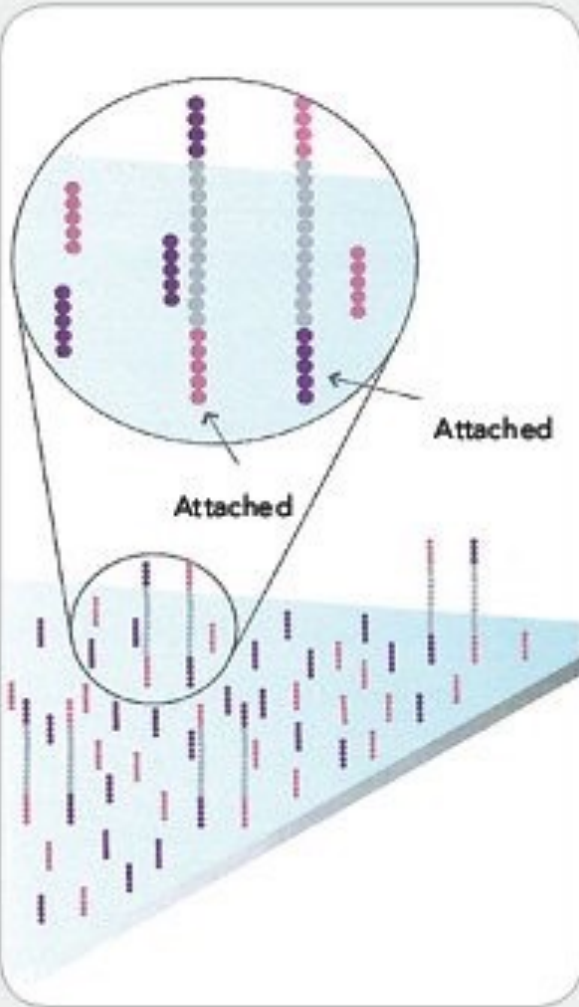


The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

5. DENATURE THE DOUBLE-STRANDED MOLECULES

Denature the Double Stranded Molecules



Attached

Attached

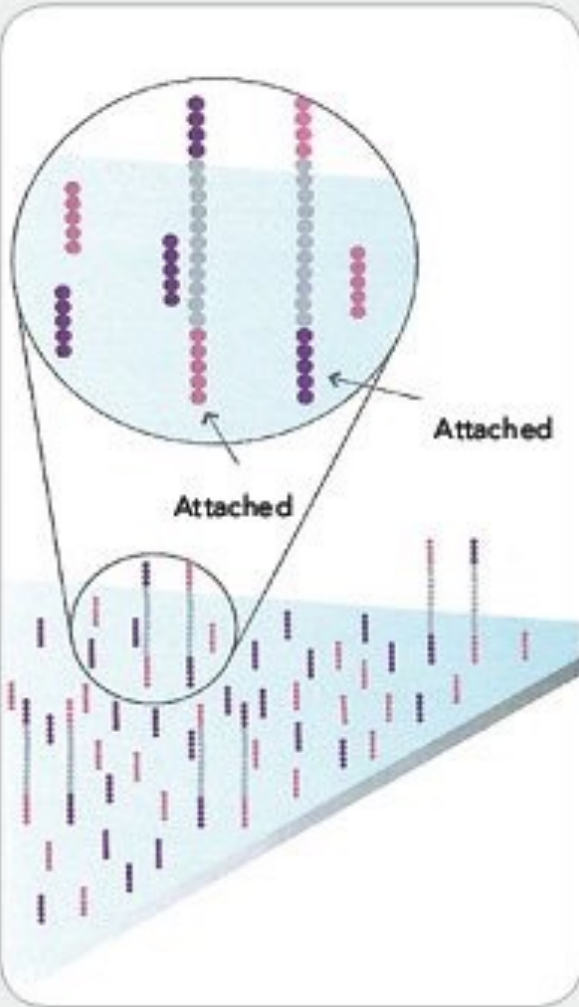
- The original strand is then washed away, leaving only the strands that had been synthesized to the oligos attached to the flow cell

Denaturation leaves single-stranded templates anchored to the substrate.

Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

5. DENATURE THE DOUBLE-STRANDED MOLECULES

Denature the Double Stranded Molecules



Attached

Attached

- The original strand is then washed away, leaving only the strands that had been synthesized to the oligos attached to the flow cell

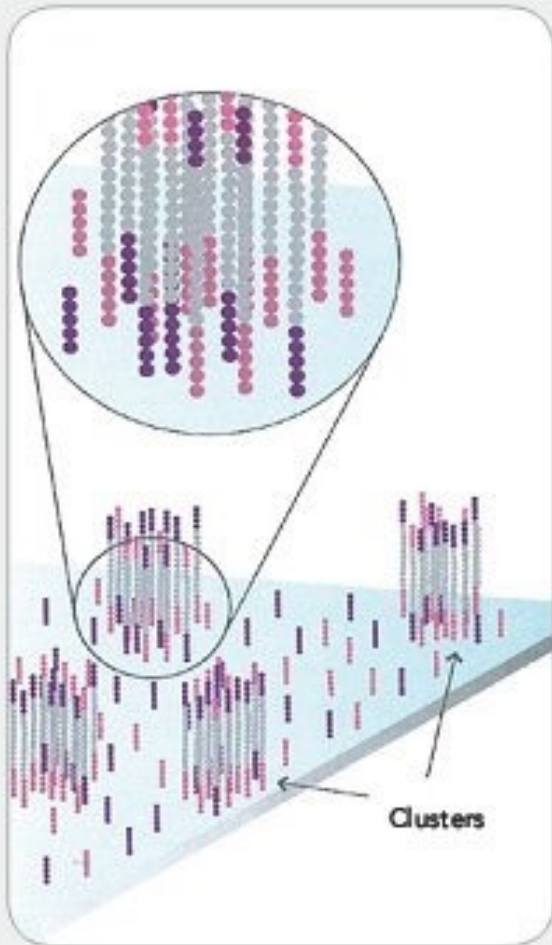
Denaturation leaves single-stranded templates anchored to the substrate.

Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

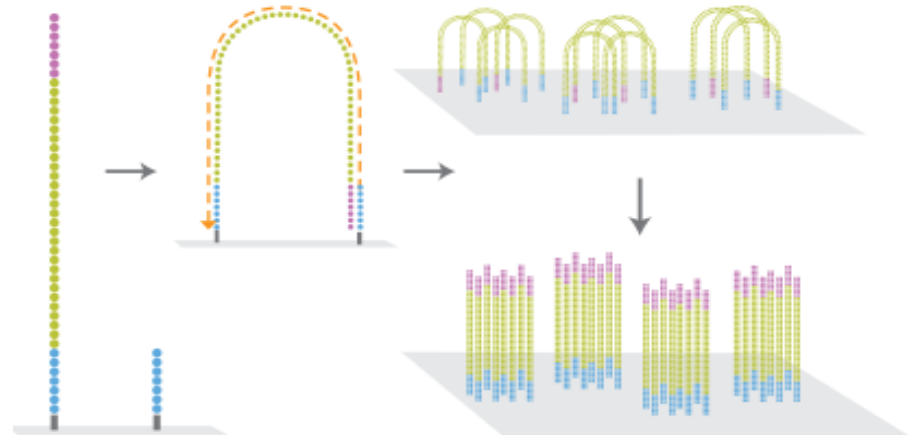


Steps 5-7 Repeats

- Cycle of new strand synthesis and Denaturation to make multiple copies of the same sequence (amplification)
 - Fragments Become Double Stranded
 - Denature the Double Strand Molecules



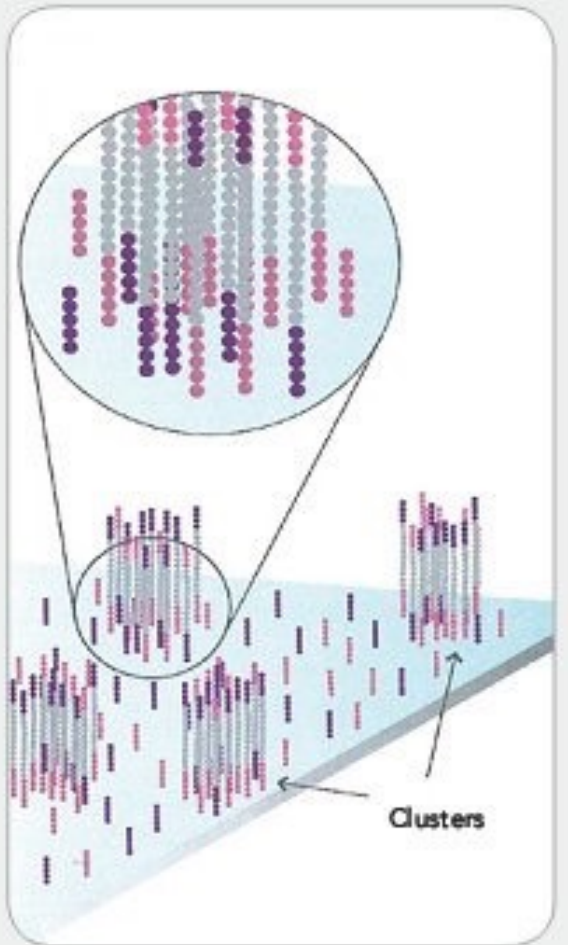
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.



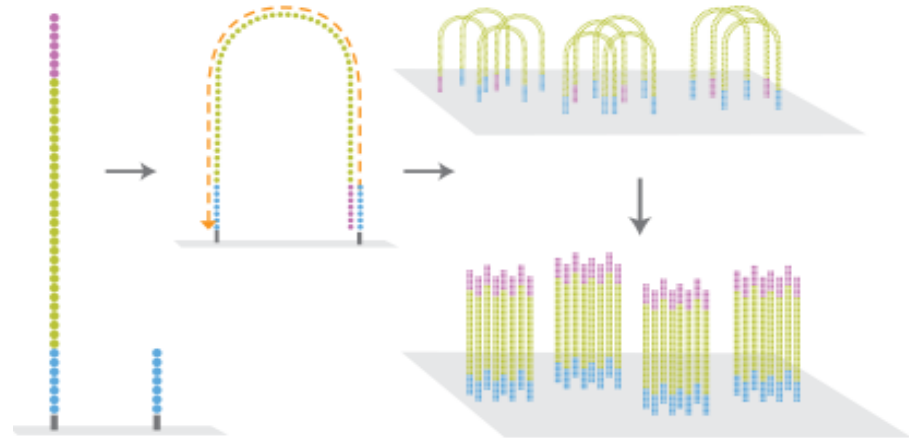


Steps 5-7 Repeats

- Cycle of new strand synthesis and Denaturation to make multiple copies of the same sequence (amplification)
 - Fragments Become Double Stranded
 - Denature the Double Strand Molecules



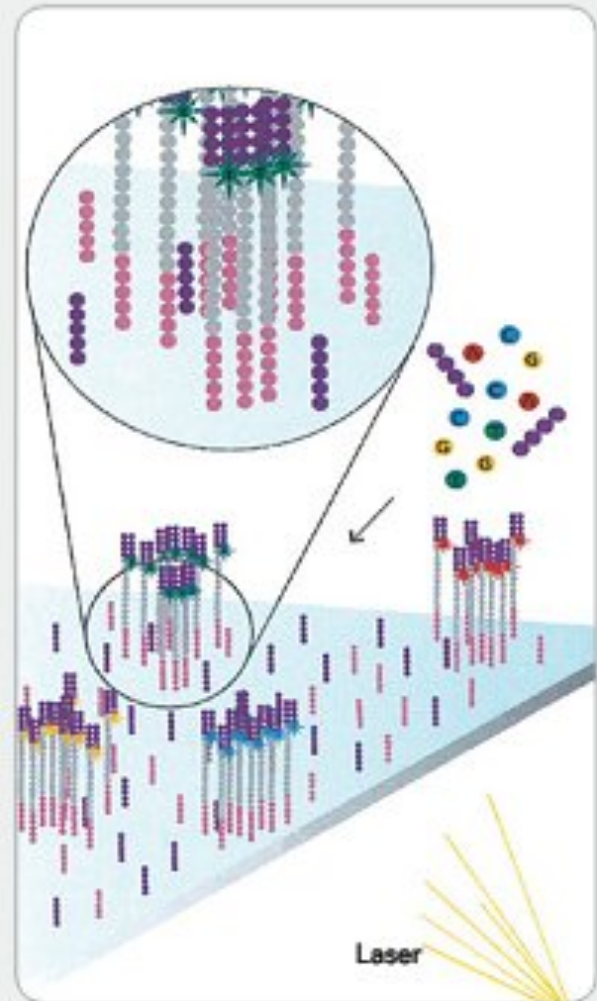
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.





Determine First Base

- The P5 region is cleaved
- Add sequencing reagents
 - Primers
 - Polymerase
 - Fluorescently labeled nucleotides
 - Buffer
- First base incorporated

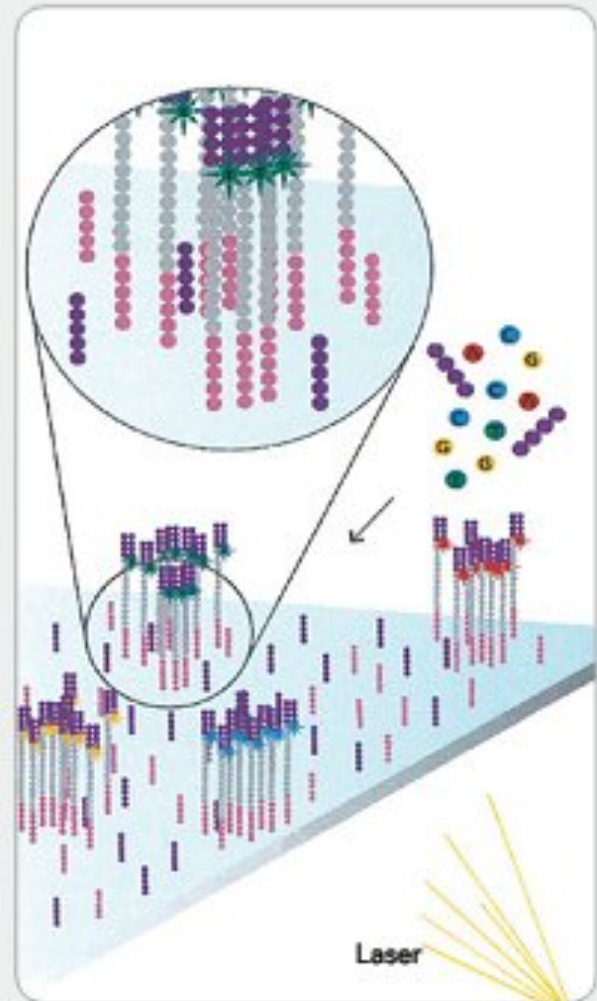


First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.



Determine First Base

- The P5 region is cleaved
- Add sequencing reagents
 - Primers
 - Polymerase
 - Fluorescently labeled nucleotides
 - Buffer
- First base incorporated

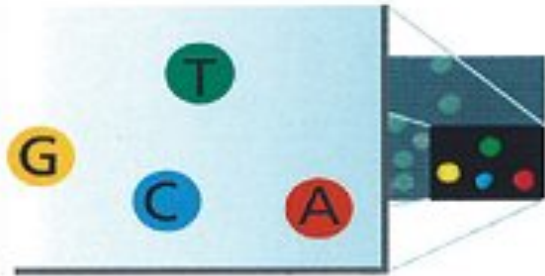


First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.



Image First Base

- Remove unincorporated bases
- Detect Signal
- Deblock and remove the fluorescent signal → new cycle



After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

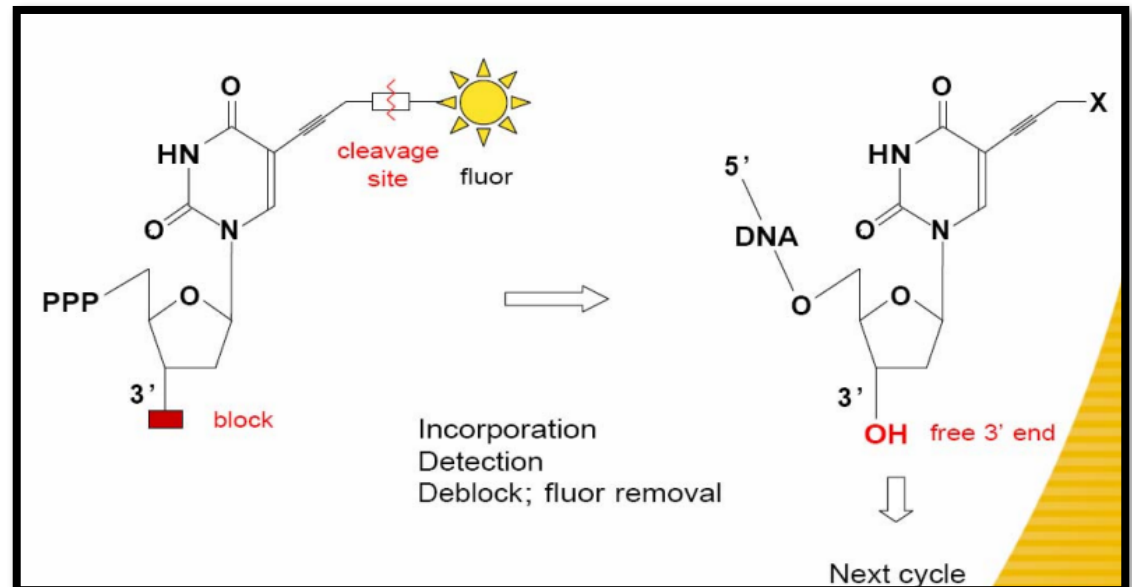


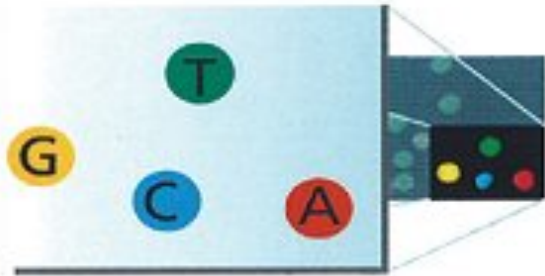
Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Image retrieved from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf



Image First Base

- Remove unincorporated bases
- Detect Signal
- Deblock and remove the fluorescent signal → new cycle



After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

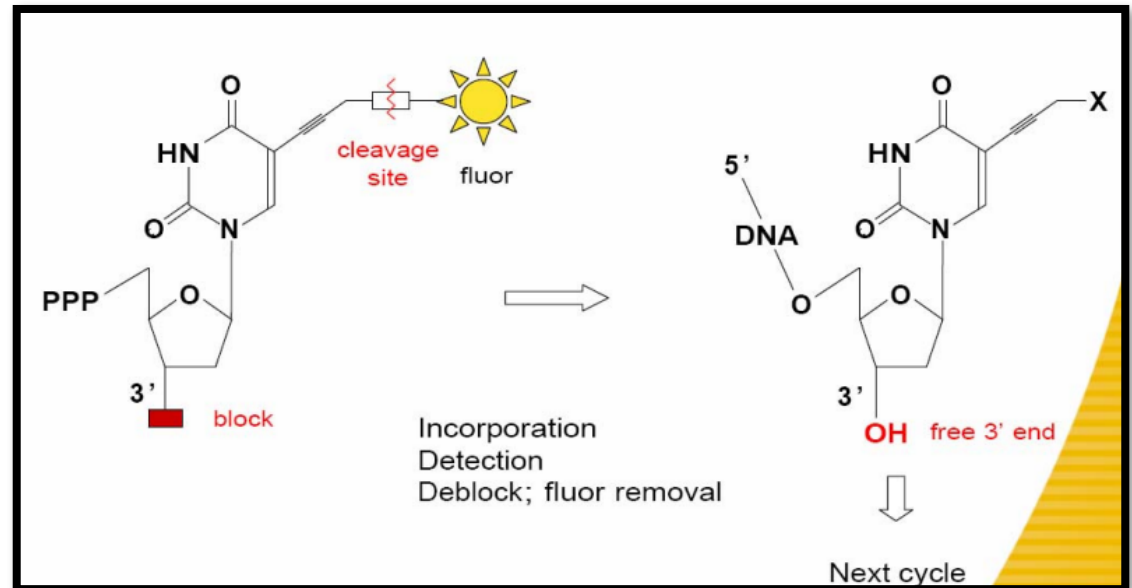
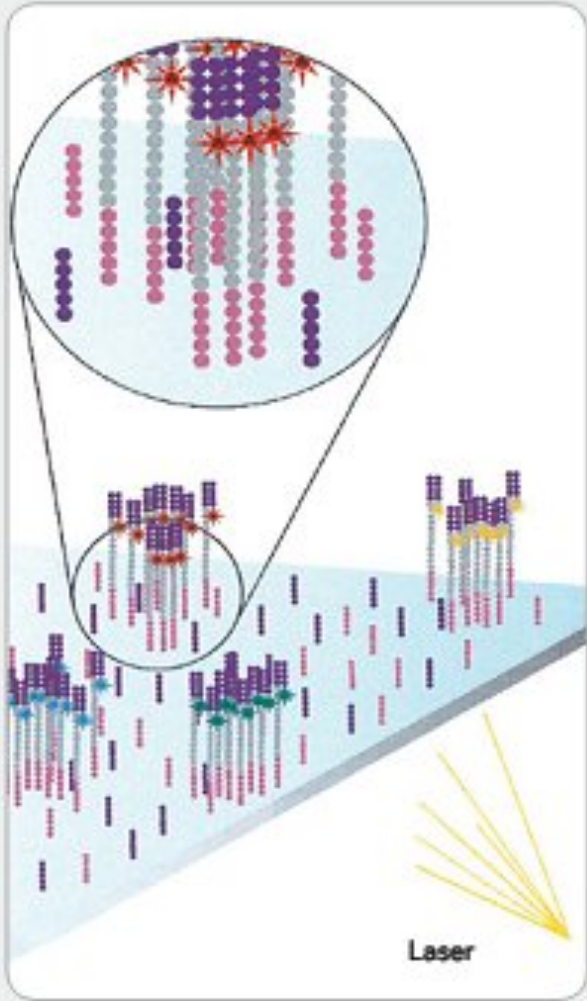


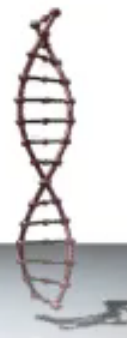
Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Image retrieved from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf

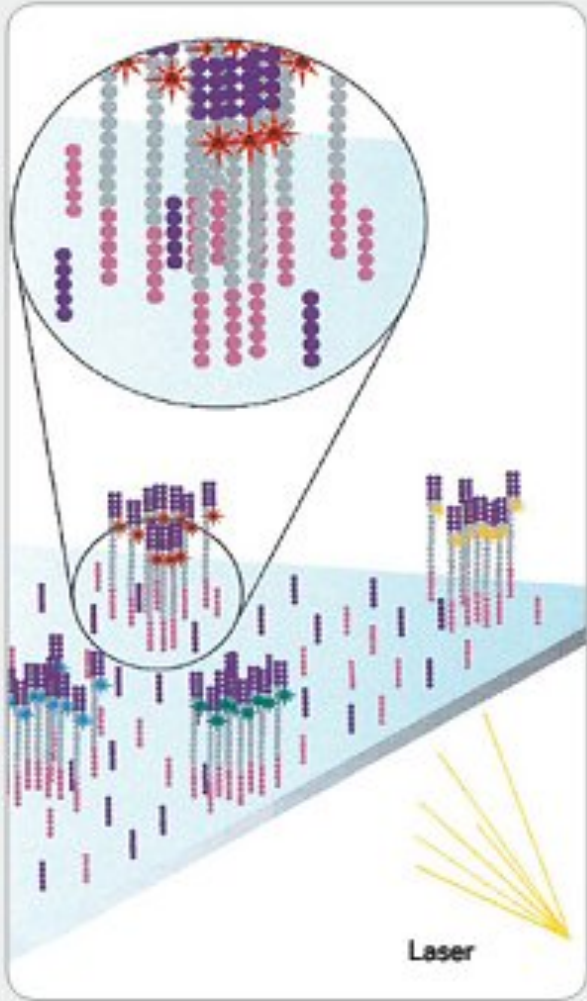


Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

Determine Second Base

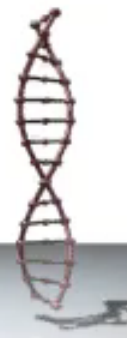


- Add sequencing reagents
 - Primers
 - Polymerase
 - Fluorescently labeled nucleotides
 - Buffer
- Second base incorporated



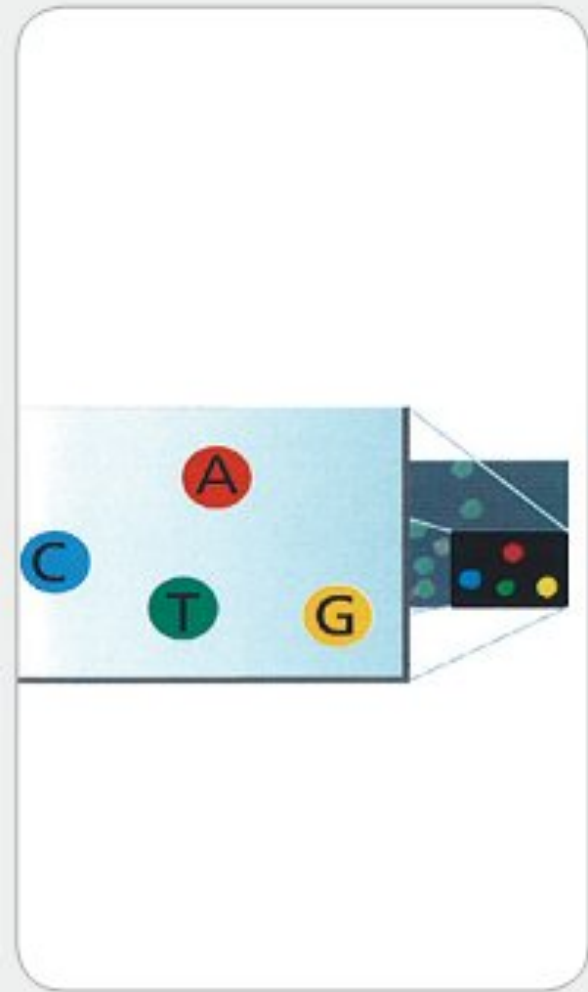
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

Determine Second Base



- Add sequencing reagents
 - Primers
 - Polymerase
 - Fluorescently labeled nucleotides
 - Buffer
- Second base incorporated

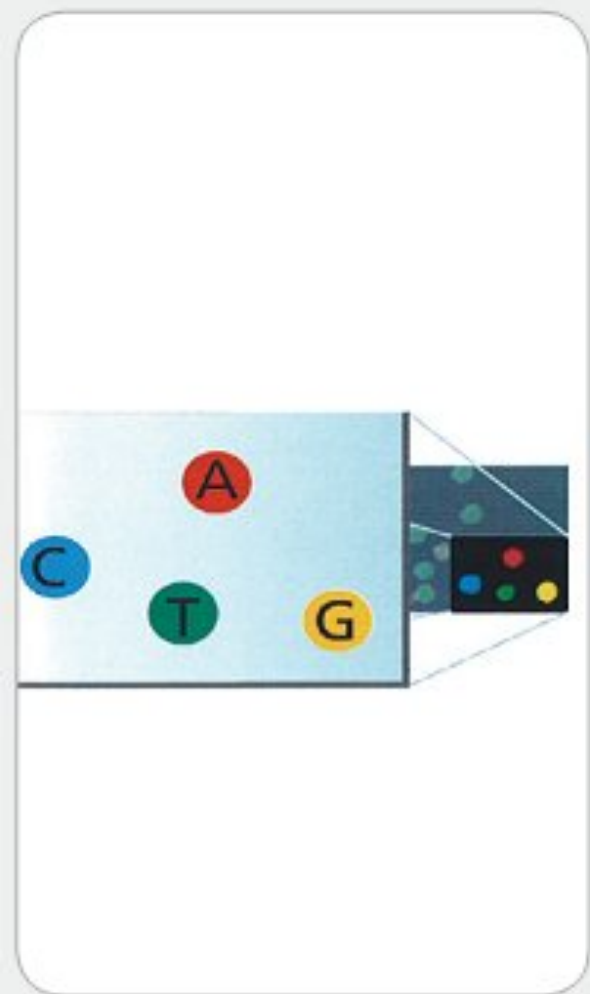
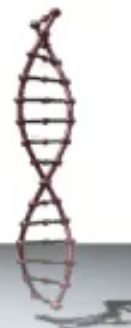
Image Second Chemistry Cycle



After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

- Remove unincorporated bases
- Detect Signal
- Unblock and remove the fluorescent signal → new cycle

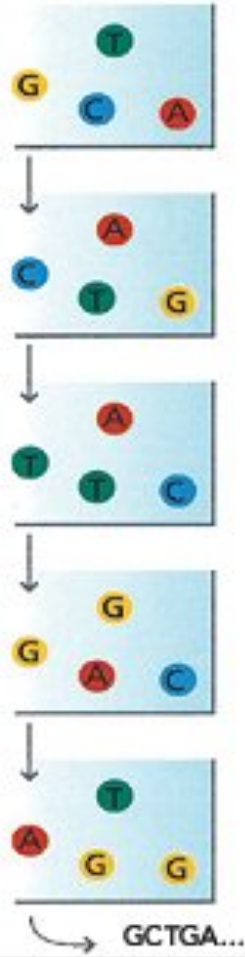
Image Second Chemistry Cycle



After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

- Remove unincorporated bases
- Detect Signal
- Deblock and remove the fluorescent signal → new cycle

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at time.

Sequence Reads Over Multiple Chemistry Cycles



- The identity of each base of a cluster is read off from sequential images

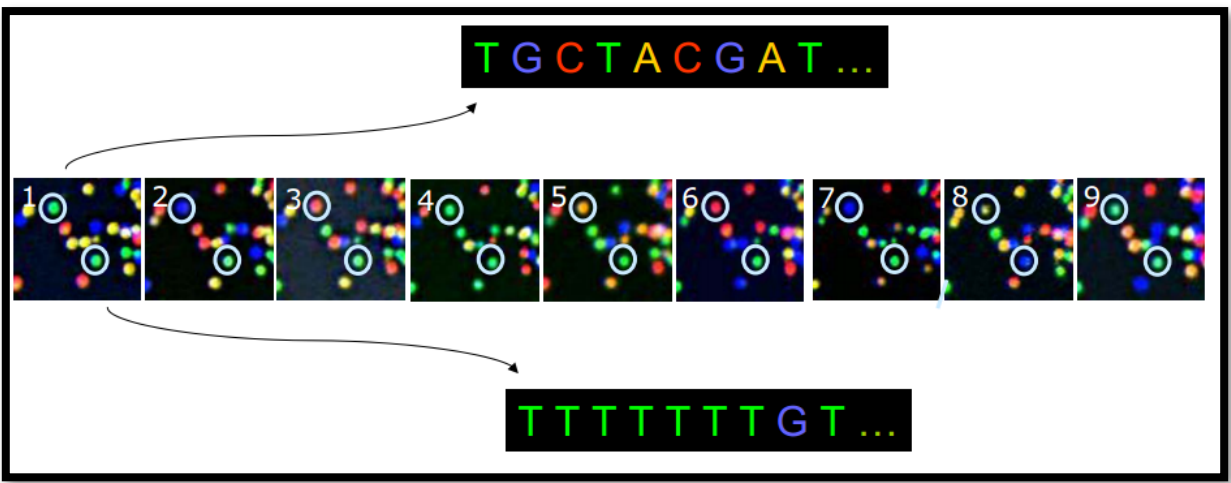
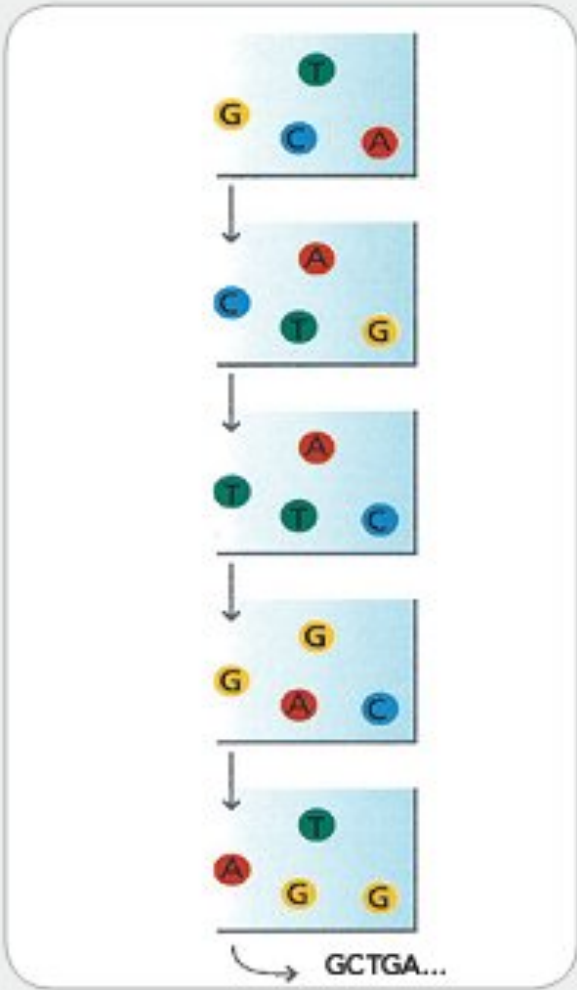


Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

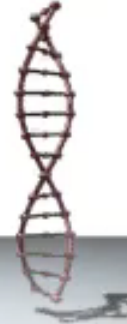
Image retrieved from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at time.

Sequence Reads Over Multiple Chemistry Cycles



- The identity of each base of a cluster is read off from sequential images

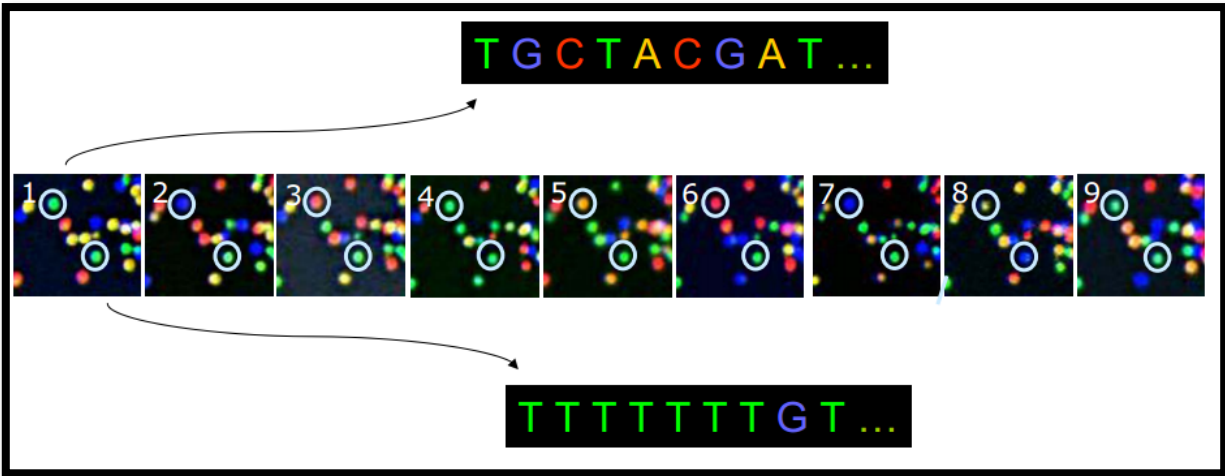
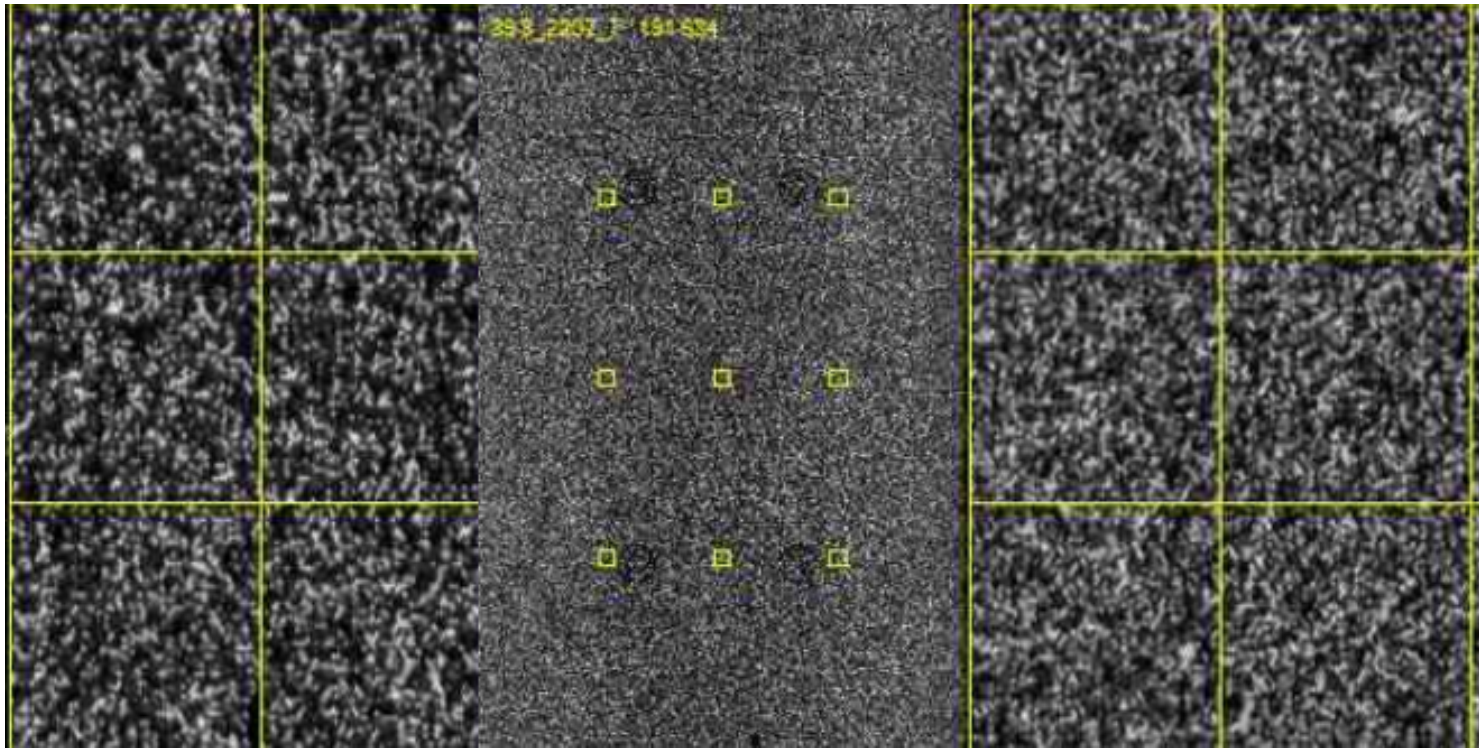


Image retrieved from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf

Image retrieved from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf

Sequencing by synthesis



Actual Illumina HiSeq 3000 image

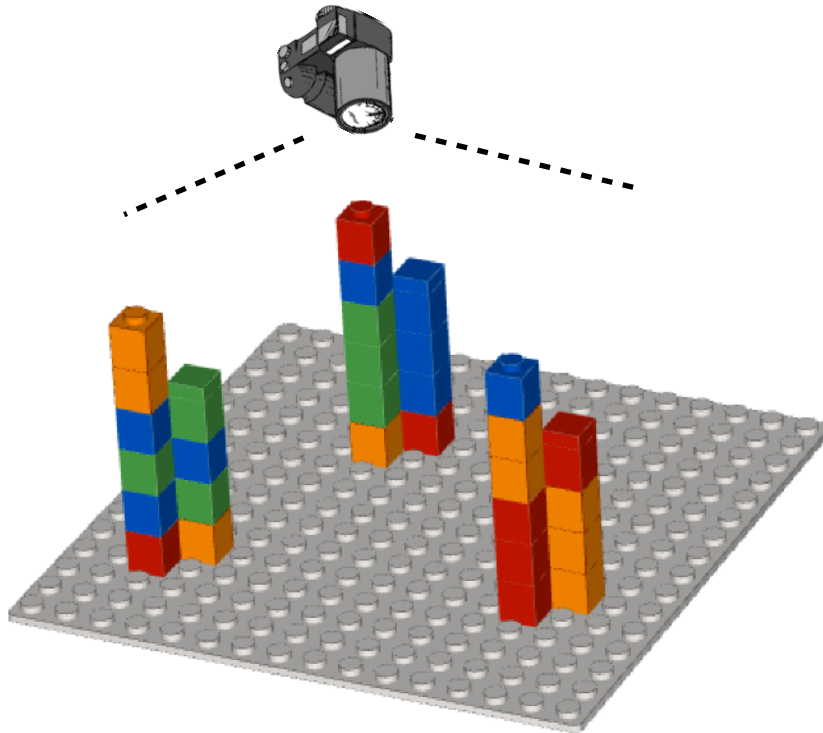
<http://dnatech.genomecenter.ucdavis.edu/2015/05/07/first-hiseq-3000-data-download/>
This and following slides marked with * courtesy of Ben Langmead

Sequencing by synthesis

Billions of templates on a slide

Massively parallel: photograph captures all templates simultaneously

Terminators are “speed bumps,” keeping reactions in sync

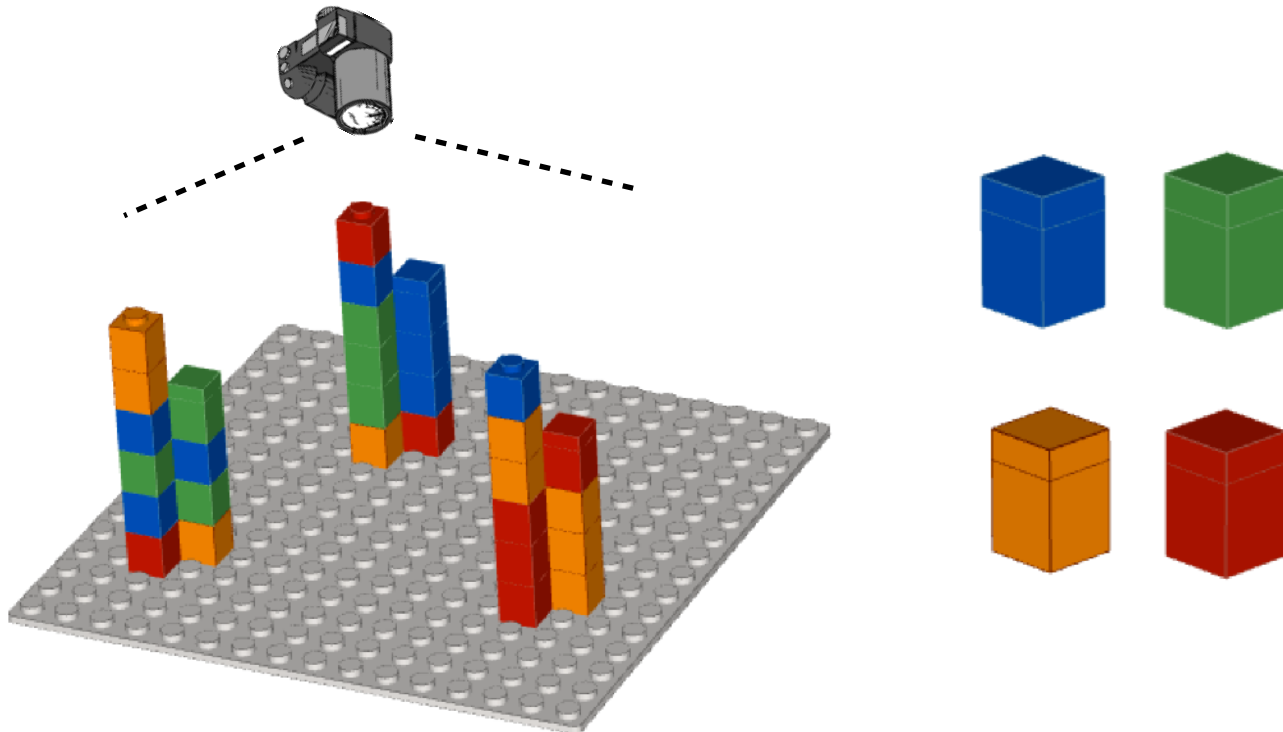


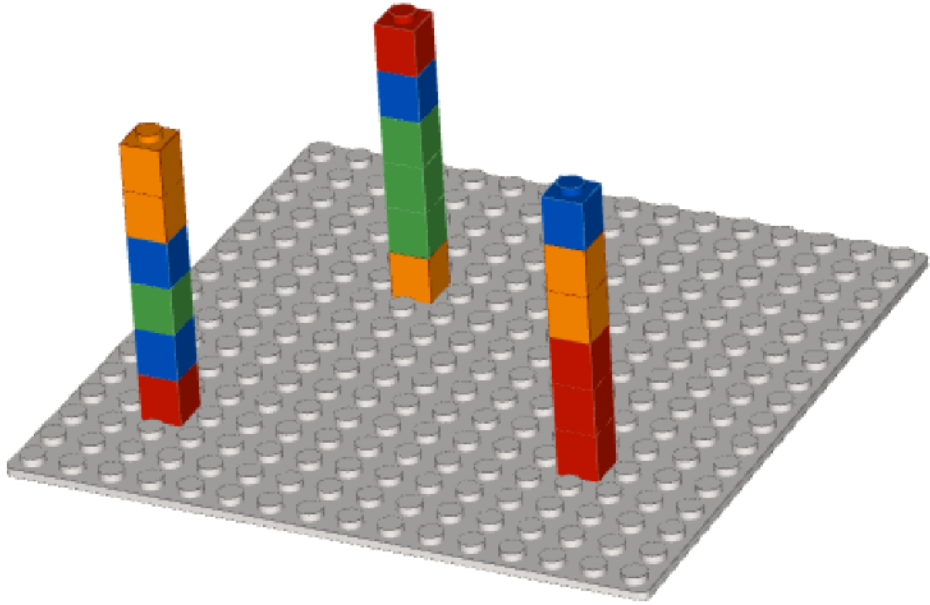
Sequencing by synthesis

Billions of templates on a slide

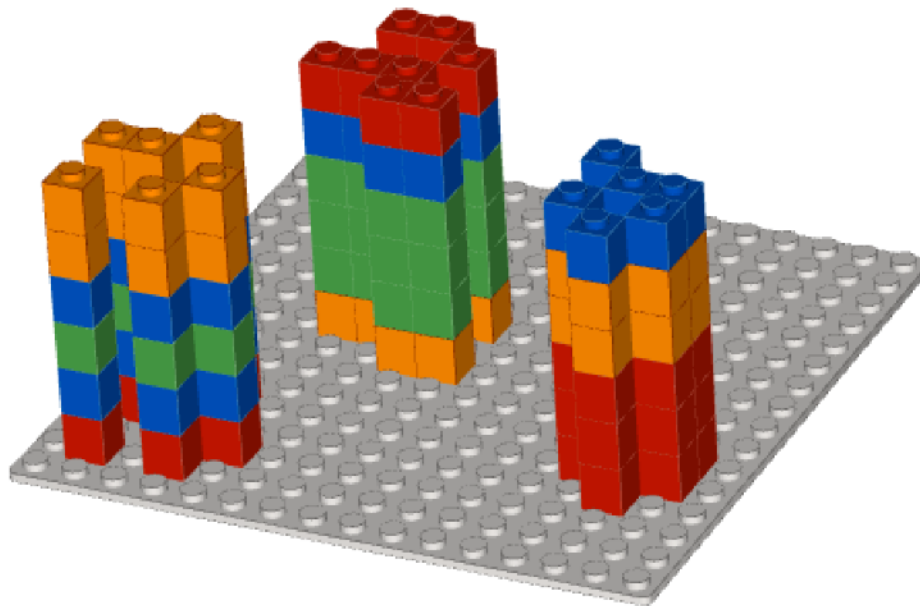
Massively parallel: photograph captures all templates simultaneously

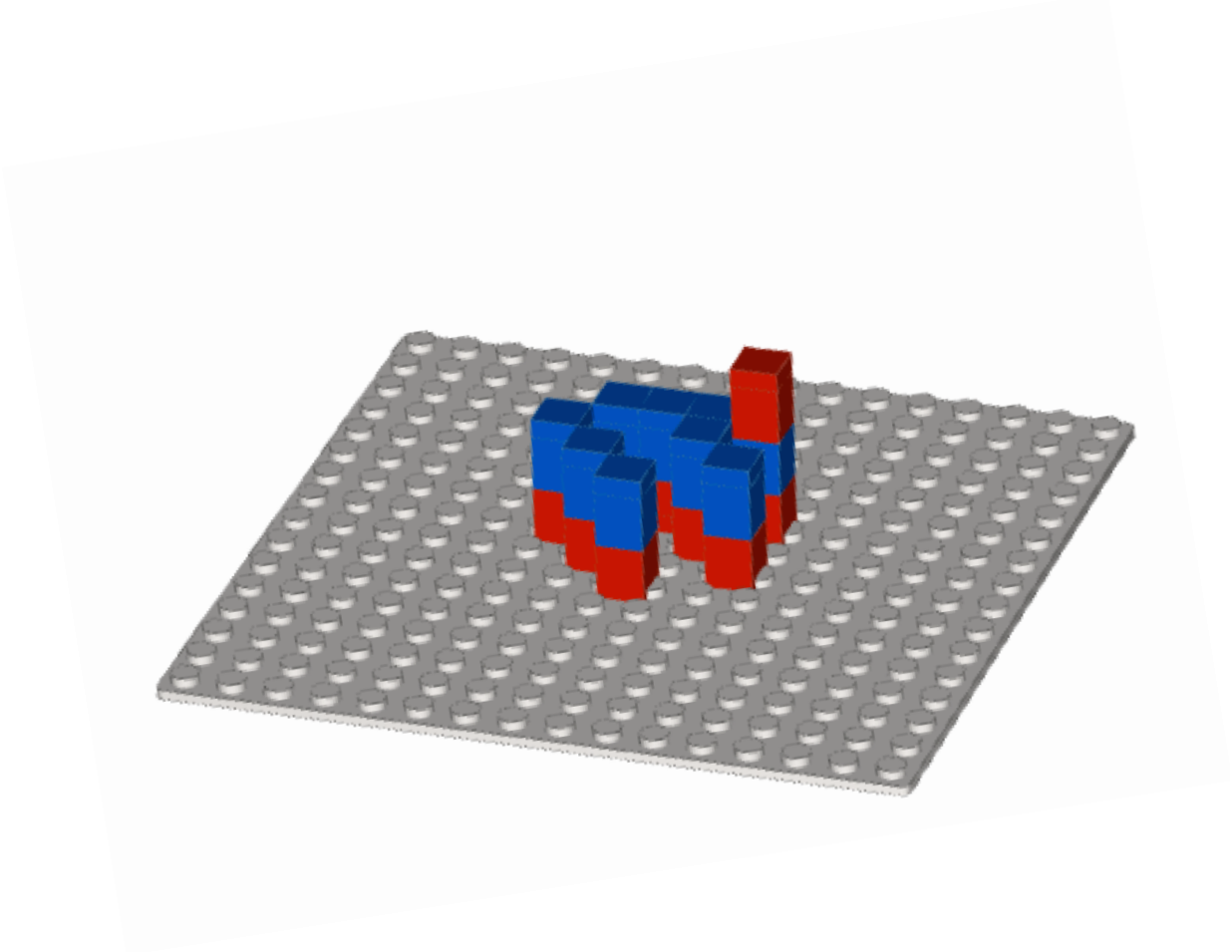
Terminators are “speed bumps,” keeping reactions in sync

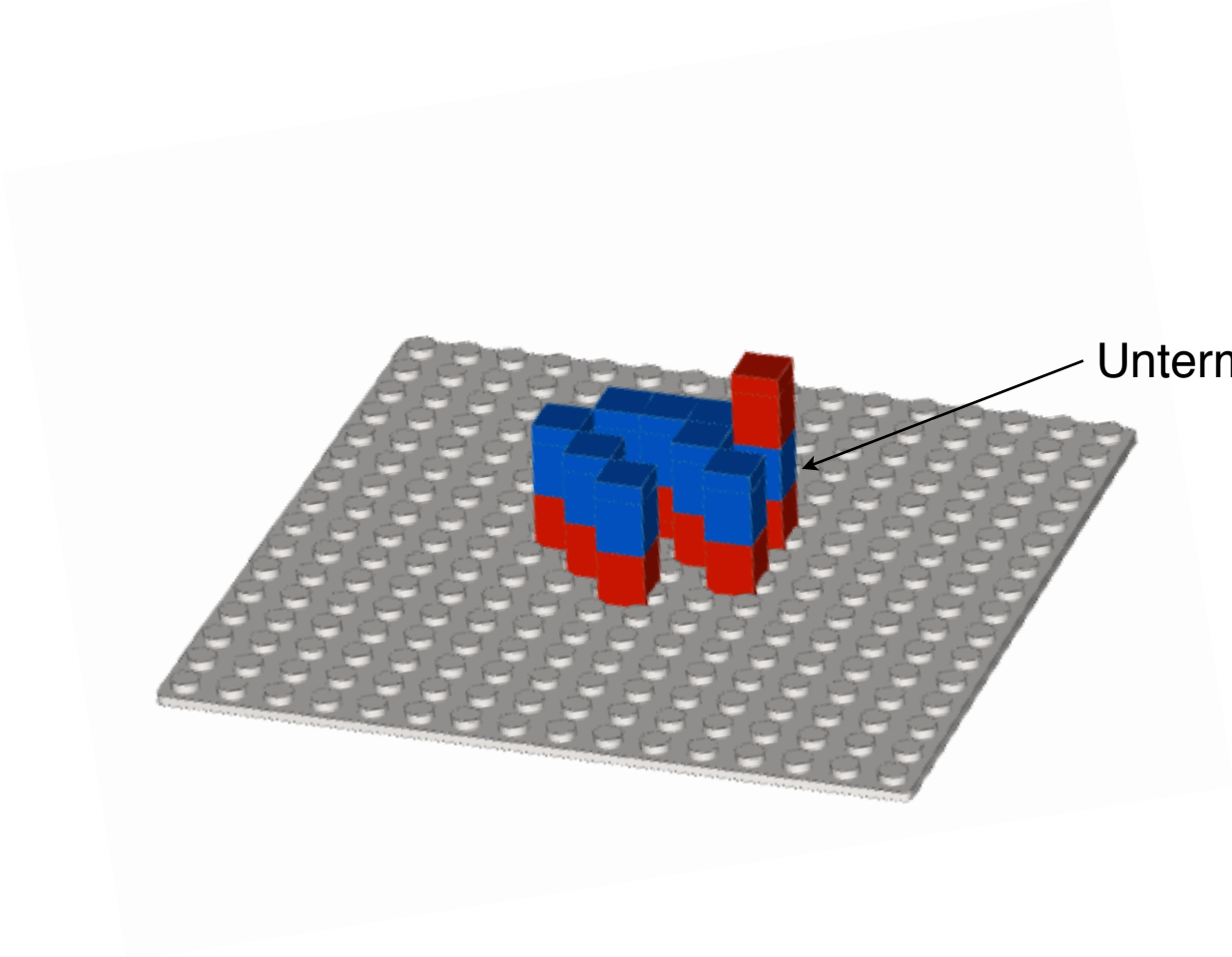




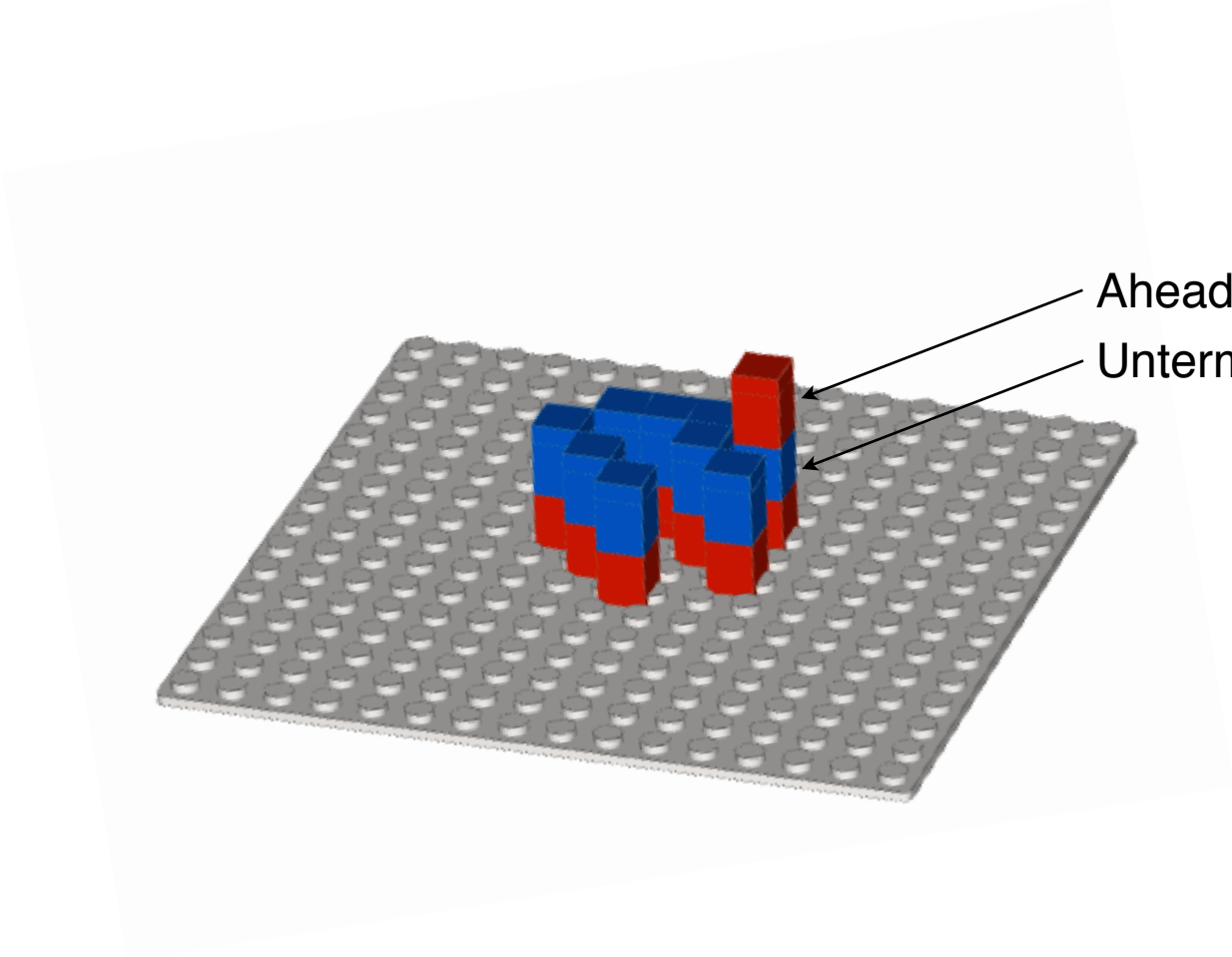
Cluster of clones





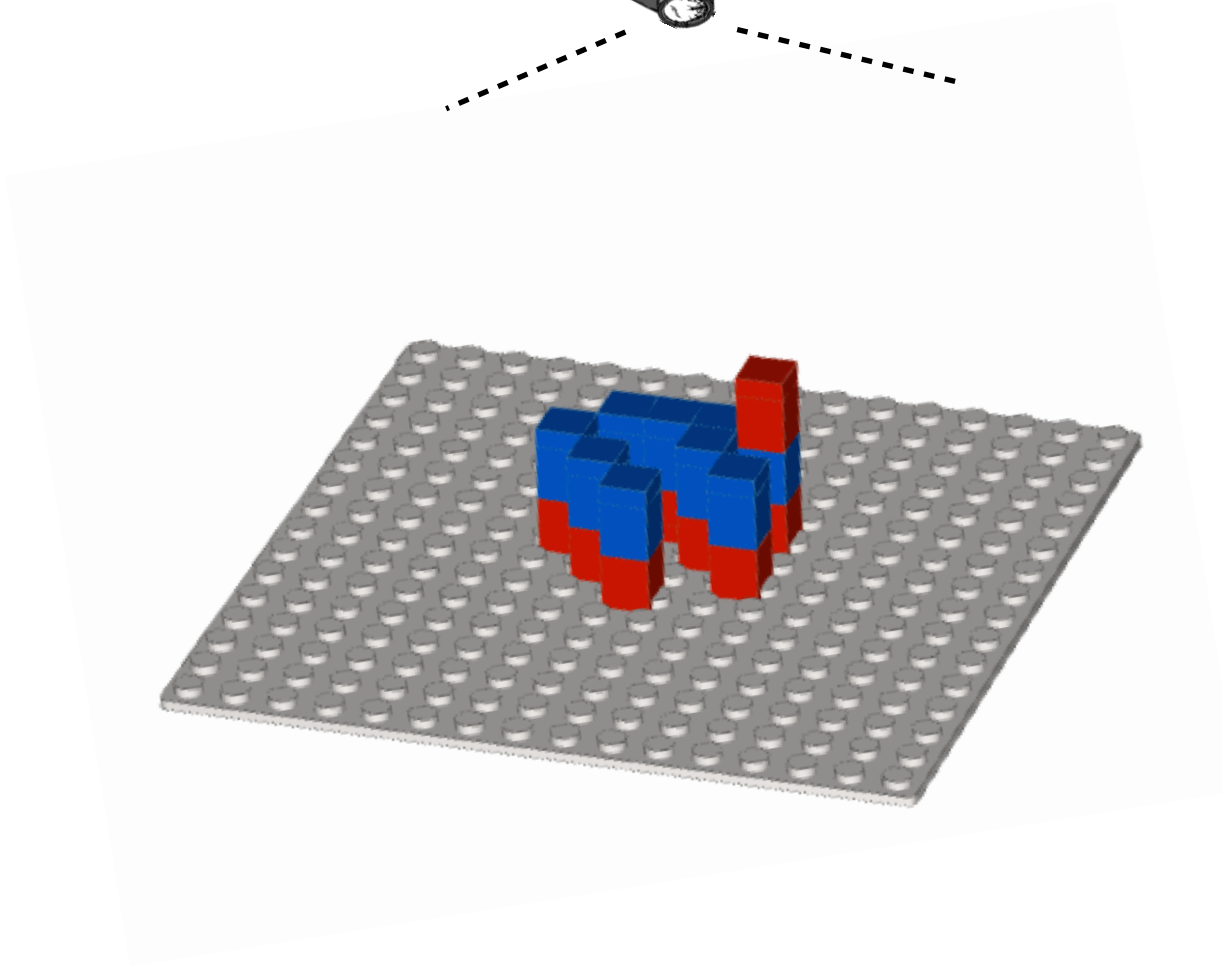
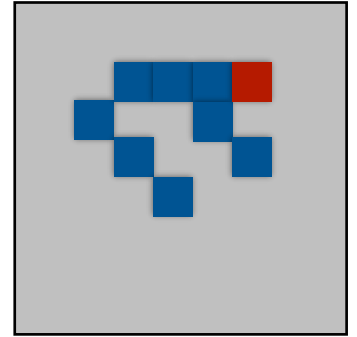
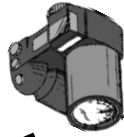


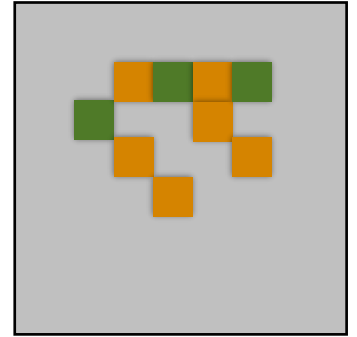
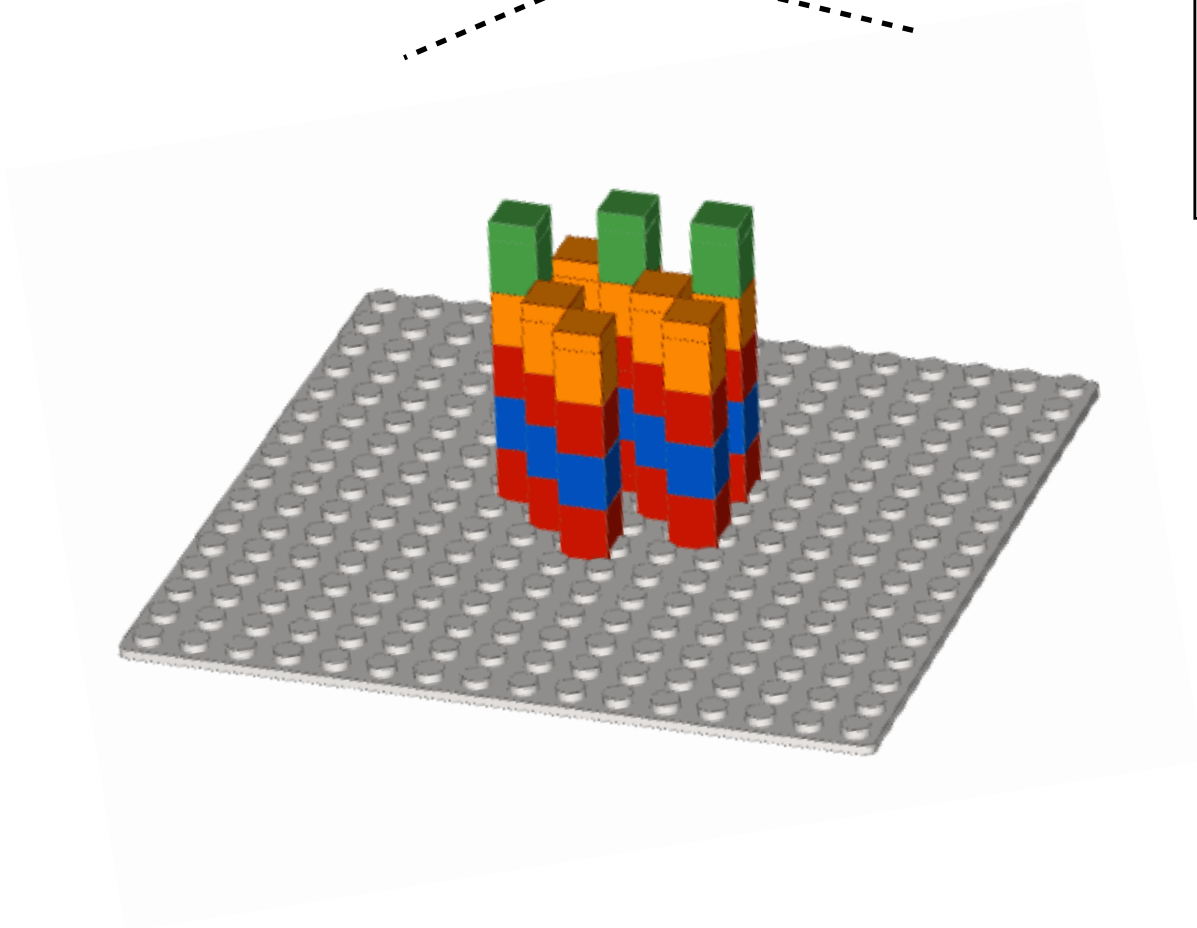
Unterminated



Ahead of schedule

Unterminated





$$Q = -10 \cdot \log_{10} p$$

$$Q = -10 \cdot \log_{10} p$$



Base quality

$$Q = -10 \cdot \log_{10} p$$

Base quality

Probability that
base call is
incorrect

$$Q = -10 \cdot \log_{10} p$$

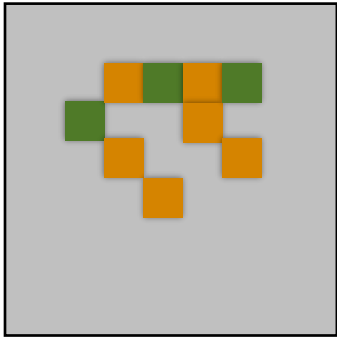
Base quality

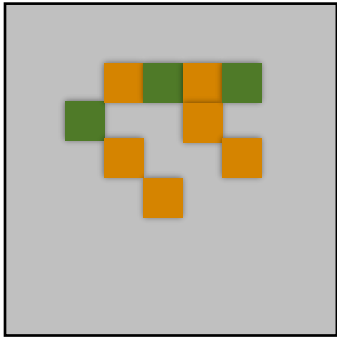
Probability that base call is incorrect

$Q = 10 \rightarrow 1$ in 10 chance call is incorrect

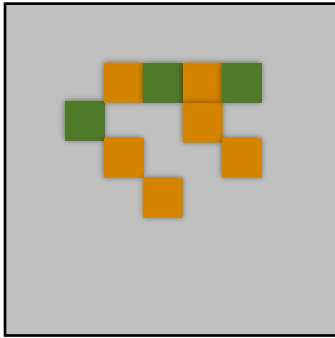
$Q = 20 \rightarrow 1$ in 100

$Q = 30 \rightarrow 1$ in 1,000



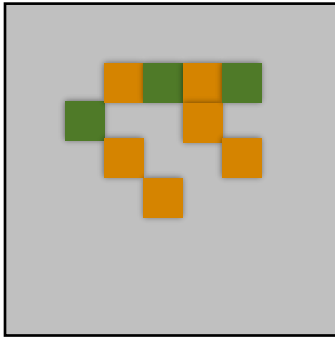


Call: orange (C)



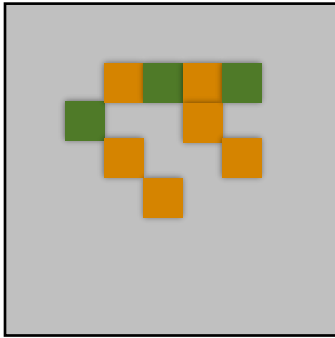
Call: orange (C)

Estimate p , probability incorrect:



Call: orange (C)

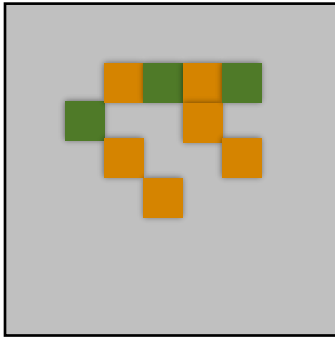
Estimate p , probability incorrect:
non-orange light / total light



Call: orange (C)

Estimate p , probability incorrect:
non-orange light / total light

$$p = 3 \text{ green} / 9 \text{ total} = 1/3$$

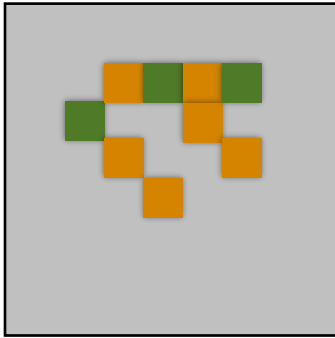


Call: orange (C)

Estimate p , probability incorrect:
non-orange light / total light

$$p = 3 \text{ green} / 9 \text{ total} = 1/3$$

$$Q = -10 \log_{10} 1/3$$



Call: orange (C)

Estimate p , probability incorrect:
non-orange light / total light

$$p = 3 \text{ green} / 9 \text{ total} = 1/3$$

$$Q = -10 \log_{10} 1/3 = 4.77$$

A read in FASTQ format

```
@ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1  
ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT  
+  
?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G
```

A read in FASTQ format

Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
+
?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

A read in FASTQ format

Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
+
?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

A read in FASTQ format

Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
(ignore) +
?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

A read in FASTQ format

```
Name @ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
(ignore) +
Base qualities ?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G
```


Base qualities

Bases and qualities line up:

```
AGCTCTGGTGACCCATGGGCAGCTGCTAGGGA
|||||
HHHHHHHHHHHHHHGCGC5FEFFFGHHHHHH
```

Base quality is ASCII-encoded version of $Q = -10 \log_{10} p$

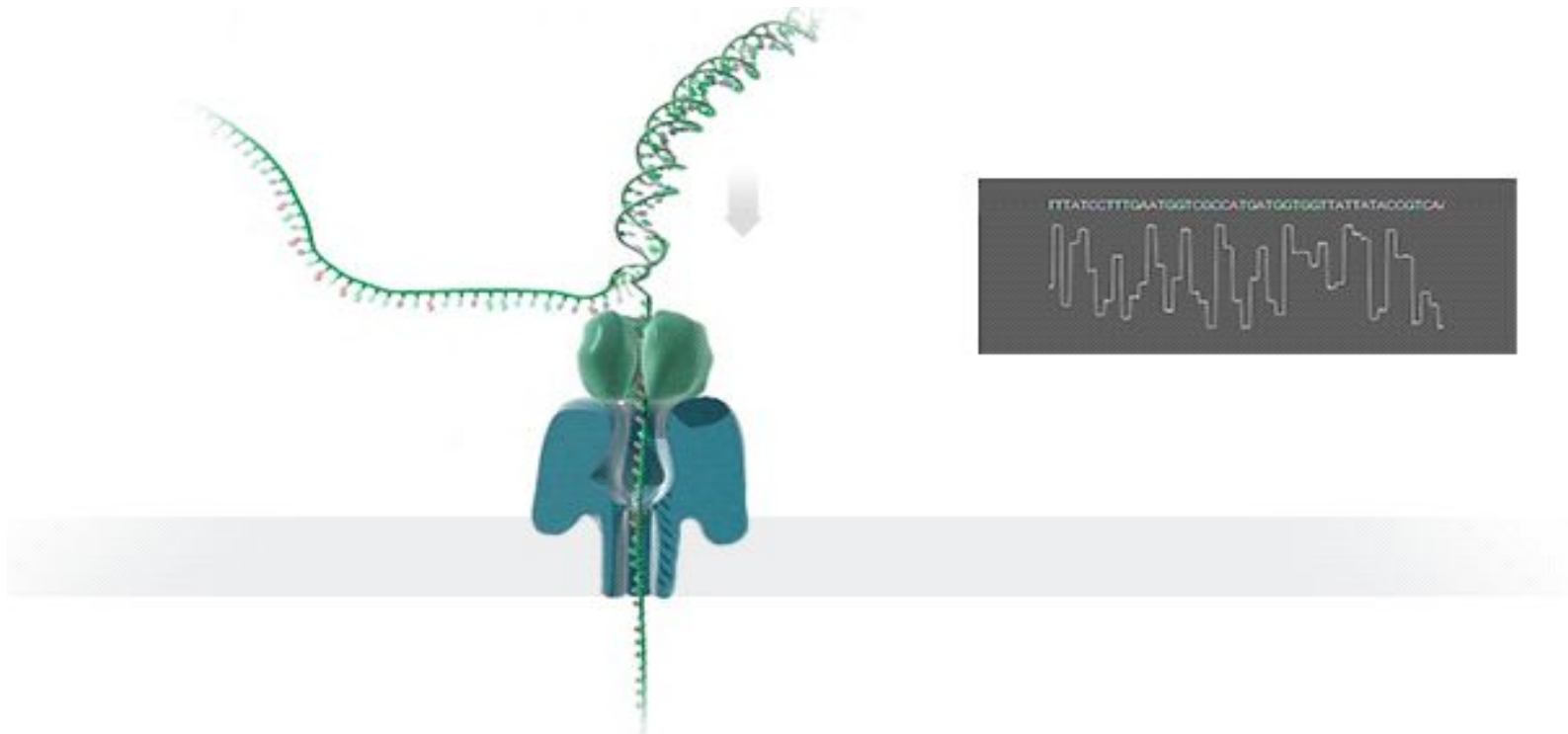
Long-read sequencing via nanopores



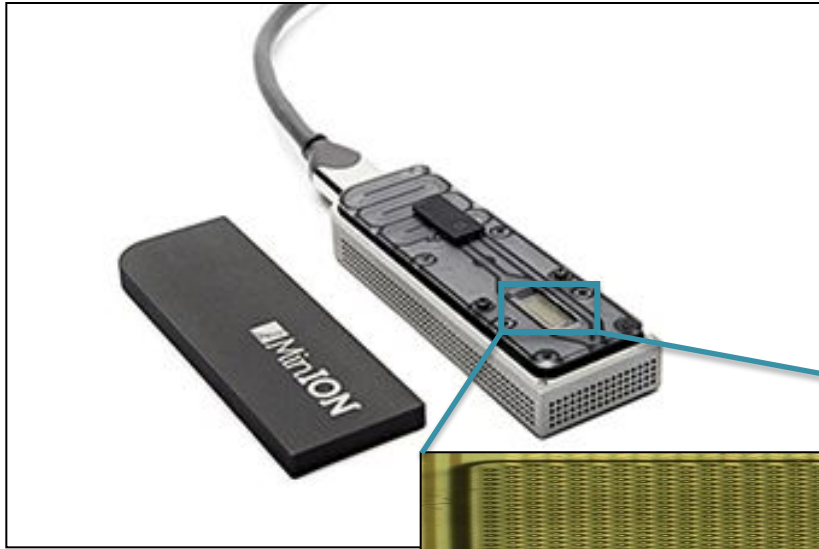
Following slides taken or adapted from Mike Schatz:
<http://schatz-lab.org/appliedgenomics2019/lectures05.LinkedReaderAndLongReads.pdf>

Nanopore Sequencing

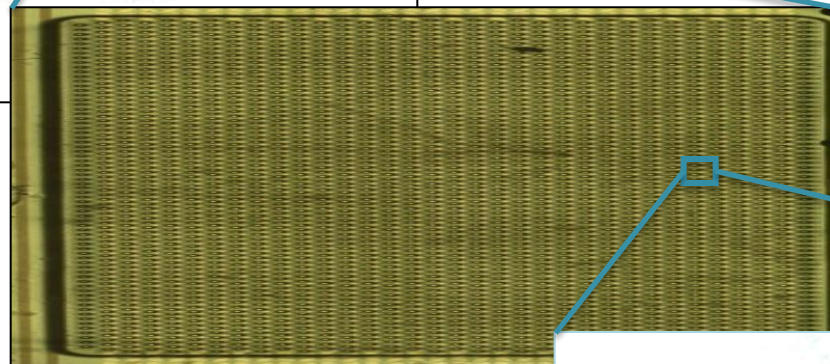
Sequences DNA/RNA by measuring changes in ionic current as nucleotide strand passes through a pore



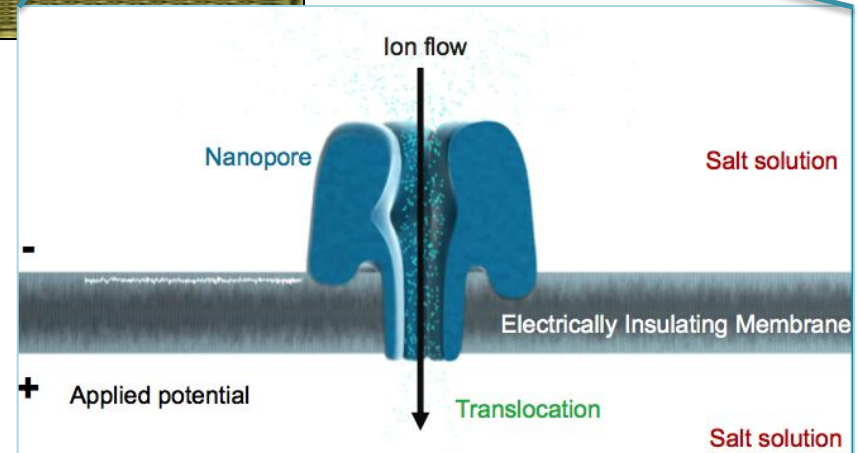
Oxford Nanopore MinION



- Thumb drive sized sequencer powered over USB
- Contains 512 channels
- Four pores per channel, only one pore active at a time



- Early access began in 2014
- Officially released in 2015



“Ultra-Long Read” Assembly

nature
biotechnology

OPEN

Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain^{1,13}, Sergey Koren^{2,13}, Karen H Miga^{1,13}, Josh Quick^{3,13}, Arthur C Rand^{1,13}, Thomas A Sasaki^{4,5,13}, John R Tyson^{6,13}, Andrew D Beggs⁷, Alexander T Dilthey², Ian T Fiddes¹, Sunir Malla⁸, Hannah Marriott⁸, Tom Nieto⁷, Justin O’Grady⁹, Hugh E Olsen¹, Brent S Pedersen^{4,5}, Arang Rhie², Hollian Richardson⁹, Aaron R Quinlan^{4,5,10}, Terrance P Snutch⁶, Louise Tee⁷, Benedict Paten¹, Adam M Phillippy², Jared T Simpson^{11,12}, Nicholas J Loman³ & Matthew Loose⁸

We report the sequencing and assembly of a reference genome for the human GM12878 Utah/Ceph cell line using the MinION (Oxford Nanopore Technologies) nanopore sequencer. 91.2 Gb of sequence data, representing ~30x theoretical coverage, were produced. Reference-based alignment enabled detection of large structural variants and epigenetic modifications. *De novo* assembly of nanopore reads alone yielded a contiguous assembly (NG50 ~3 Mb). We developed a protocol to generate ultra-long reads (N50 > 100 kb, read lengths up to 882 kb). Incorporating an additional 5x coverage of these ultra-long reads more than doubled the assembly contiguity (NG50 ~6.4 Mb). The final assembled genome was 2,867 million bases in size, covering 85.8% of the reference. Assembly accuracy, after incorporating complementary short-read sequencing data, exceeded 99.8%. Ultra-long reads enabled assembly and phasing of the 4-Mb major histocompatibility complex (MHC) locus in its entirety, measurement of telomere repeat length, and closure of gaps in the reference human genome assembly GRCh38.

The human genome is used as a yardstick to assess performance of DNA sequencing instruments^{1–5}. Despite improvements in sequencing technology, assembling human genomes with high accuracy and completeness remains challenging. This is due to size (~3.1 Gb), heterozygosity, regions of GC% bias, diverse repeat families, and segmental duplications (up to 1.7 Mbp in size) that make up at least 50% of the genome⁶. Even more challenging are the pericentromeric, centromeric, and acrocentric short arms of chromosomes, which contain satellite DNA and tandem repeats of 3–10 Mb in length^{7,8}. Repetitive structures pose challenges for *de novo* assembly using “short read” sequencing technologies, such as Illumina. Such data, while enabling highly accurate genotyping in non-repetitive regions, do not provide contiguous *de novo* assemblies. This limits the ability to reconstruct repetitive sequences, detect complex structural variation, and fully characterize the human genome.

Single-molecule sequencers, such as Pacific Biosciences’ (PacBio), can produce read lengths of 10 kb or more, which makes *de novo* human genome assembly more tractable⁹. However, single-molecule sequencing reads have significantly higher error rates compared with Illumina sequencing. This has necessitated development of *de novo* assembly

algorithms and the use of long noisy data in conjunction with accurate short reads to produce high-quality reference genomes¹⁰. In May 2014, the MinION nanopore sequencer was made available to early-access users¹¹. Initially, the MinION nanopore sequencer was used to sequence and assemble microbial genomes or PCR products^{12–14} because the output was limited to 500 Mb to 2 Gb of sequenced bases. More recently, assemblies of eukaryotic genomes including yeasts, fungi, and *Caenorhabditis elegans* have been reported^{15–17}.

Recent improvements to the protein pore (a laboratory-evolved *Escherichia coli* CsgG mutant named R9.4), library preparation techniques (1D ligation and 1D rapid), sequencing speed (450 bases/s), and control software have increased throughput, so we hypothesized that whole-genome sequencing (WGS) of a human genome might be feasible using only a MinION nanopore sequencer^{17–19}.

We report sequencing and assembly of a reference human genome for GM12878 from the Utah/CEPH pedigree, using MinION R9.4 1D chemistry, including ultra-long reads up to 882 kb in length. GM12878 has been sequenced on a wide variety of platforms, and has well-validated variation call sets, which enabled us to benchmark our results²⁰.

¹UC Santa Cruz Genomics Institute, University of California, Santa Cruz, California, USA. ²Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, Bethesda, Maryland, USA. ³Institute of Microbiology and Infection, University of Birmingham, Birmingham, UK. ⁴Department of Human Genetics, University of Utah, Salt Lake City, Utah, USA. ⁵USTAR Center for Genetic Discovery, University of Utah, Salt Lake City, Utah, USA. ⁶Michael Smith Laboratories and Djavad Mowafaghian Centre for Brain Health, University of British Columbia, Vancouver, Canada. ⁷Surgical Research Laboratory, Institute of Cancer & Genomic Science, University of Birmingham, UK. ⁸DeepSec, School of Life Sciences, University of Nottingham, UK. ⁹Nottingham Medical School, University of East Anglia, Norwich, UK. ¹⁰Department of Biomedical Informatics, University of Utah, Salt Lake City, Utah, USA. ¹¹Ontario Institute for Cancer Research, Toronto, Canada. ¹²Department of Computer Science, University of Toronto, Toronto, Canada. ¹³These authors contributed equally to this work. Correspondence should be addressed to N.J.L. (n.j.loman@bham.ac.uk) or M.L. (matt.loose@nottingham.ac.uk).

Received 20 April 2017; accepted 11 December 2017; published online 29 January 2018; doi:10.1038/nbt.4060

Current Nanopore Assembly

nature
biotechnology

OPEN

Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain^{1,13}, Sergey Koren^{2,13}, Karen H Miga^{1,13}, Josh Quick^{3,13}, Arthur C Rand^{1,13}, Thomas A Sasaki^{4,5,13}, John R Tyson^{6,13}, Andrew D Beggs⁷, Alexander T Dilthey², Ian T Fiddes¹, Sunir Malla⁸, Hannah Marriott⁸, Tom Nieto⁷, Justin O'Grady⁹, Hugh E Olsen¹, Brent S Pedersen^{4,5}, Arang Rhie², Hollian Richardson⁹, Aaron R Quinlan^{4,5,10}, Terrance P Snutch⁶, Louise Tee⁷, Benedict Paten¹, Adam M Phillippy², Jared T Simpson^{11,12}, Nicholas J Loman³ & Matthew Loose⁸

We report the sequencing and assembly of a reference genome for the human GM12878 Utah/Ceph cell line using the MinION (Oxford Nanopore Technologies) nanopore sequencer. 91.2 Gb of sequence data, representing ~30x theoretical coverage, were produced. Reference modifications. *De novo* assembly protocol to generate ultra-long reads of these ultra-long reads more than 2,867 million bases in size, cover short-read sequencing data, exceed histocompatibility complex (MHC) reference human genome assembly

The human genome is used as a yardstick for DNA sequencing instruments¹⁻⁵. Despite technology, assembling human genomes remains challenging. The heterozygosity, regions of GC% bias, divergent duplications (up to 1.7 Mbp in size of the genome⁶). Even more challenging are centromeric, and acrocentric short arms of satellite DNA and tandem repeats of 3-structures pose challenges for *de novo* sequencing technologies, such as Illumina highly accurate genotyping in non-repetitive contiguous *de novo* assemblies. This is due to repetitive sequences, detect complex structures to characterize the human genome.

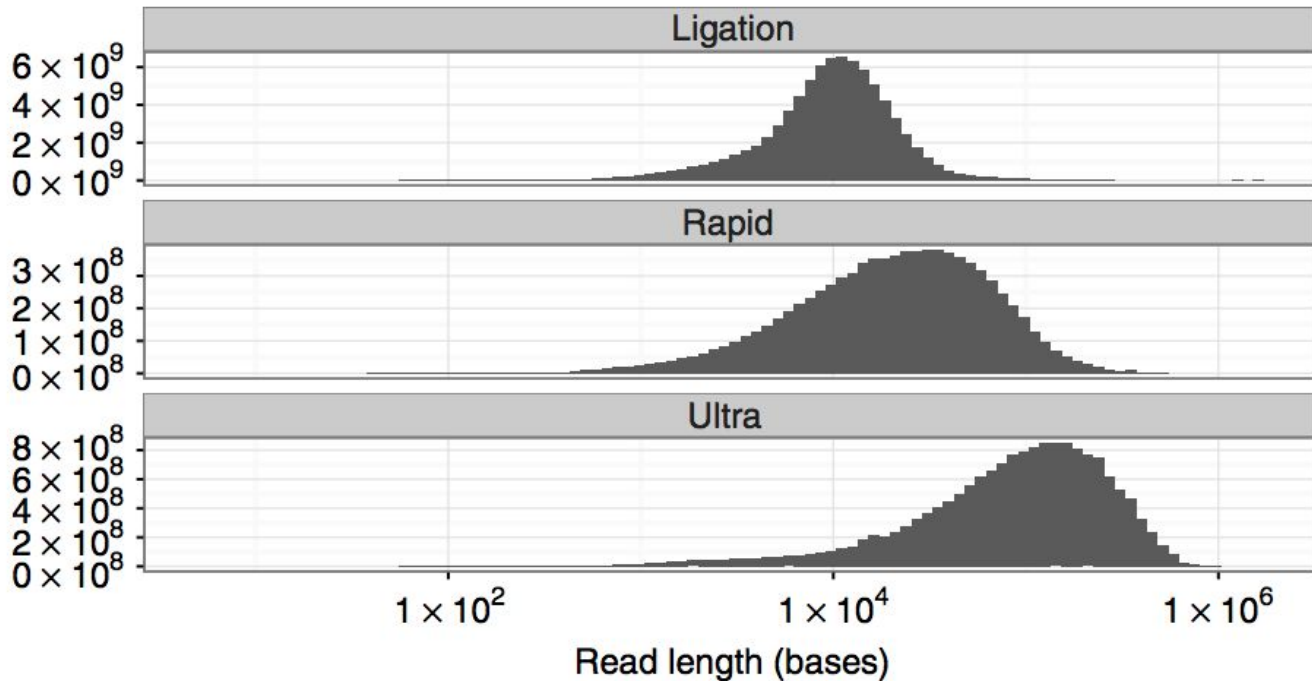
Single-molecule sequencers, such as PacBio, can produce read lengths of 10 kb or more, making genome assembly more tractable⁹. However, long reads have significantly higher error rates than short reads. This has necessitated development of new assembly algorithms.

¹UC Santa Cruz Genomics Institute, University Branch, National Human Genome Research Institute, University of California, Santa Cruz, CA, USA. ²Department of Human Genetics, University of Utah, USA. ³Michael Smith Laboratories and Centre for Genome Sciences and Policy, Genome Sciences Centre, University of British Columbia, Vancouver, BC, Canada. ⁴Institute of Cancer & Genomic Sciences, University of East Anglia, Norwich, UK. ⁵Department of Human Genetics, University of Utah, USA. ⁶Department of Human Genetics, University of Utah, USA. ⁷Department of Human Genetics, University of Utah, USA. ⁸Department of Human Genetics, University of Utah, USA. ⁹Department of Human Genetics, University of Utah, USA. ¹⁰Department of Human Genetics, University of Utah, USA. ¹¹Department of Human Genetics, University of Utah, USA. ¹²Department of Human Genetics, University of Utah, USA. ¹³Department of Human Genetics, University of Utah, USA.

Received 20 April 2017; accepted 11 December 2017

b

Cumulative length (bases)



Current Nanopore Assembly

nature
biotechnology

OPEN

Nanopore sequencing and assembly of a human genome with ultra-long reads

Miten Jain^{1,13}, Sergey Koren^{2,13}, Karen H Miga^{1,13}, Josh Quick^{3,13}, Arthur C Rand^{1,13}, Thomas A Sasaki^{4,5,13}, John R Tyson^{6,13}, Andrew D Beggs⁷, Alexander T Dilthey², Ian T Fiddes¹, Sunir Malla⁸, Hannah Marriott⁸, Tom Nieto⁷, Justin O'Grady⁹, Hugh E Olsen¹, Brent S Pedersen^{4,5}, Arang Rhie², Hollian Richardson⁹, Aaron R Quinlan^{4,5,10}, Terrance P Snutch⁶, Louise Tee⁷, Benedict Paten¹, Adam M Phillippy², Jared T Simpson^{11,12}, Nicholas J Loman³ & Matthew Loose⁸

We report the sequencing and assembly of a reference genome for the human GM12878 Utah/Ceph cell line using the MinION (Oxford Nanopore Technologies) nanopore sequencer. 91.2 Gb of sequence data, representing ~30x theoretical coverage, were produced. Reference modifications. *De novo* assembly protocol to generate ultra-long reads of these ultra-long reads more than 2,867 million bases in size, cover short-read sequencing data, exceed histocompatibility complex (MHC) reference human genome assembly

The human genome is used as a yardstick for DNA sequencing instruments¹⁻⁵. Despite the technology, assembling human genomes remains challenging. The heterozygosity, regions of GC% bias, divergent duplications (up to 1.7 Mbp in size of the genome⁶). Even more challenging are centromeric, and acrocentric short arms of satellite DNA and tandem repeats of 3-structures pose challenges for *de novo* sequencing technologies, such as Illumina highly accurate genotyping in non-repetitive sequences, detect complex structures to characterize the human genome.

Single-molecule sequencers, such as PacBio, can produce read lengths of 10 kb or more, making genome assembly more tractable⁹. However, long reads have significantly higher error rates than short reads. This has necessitated development of new assembly algorithms.

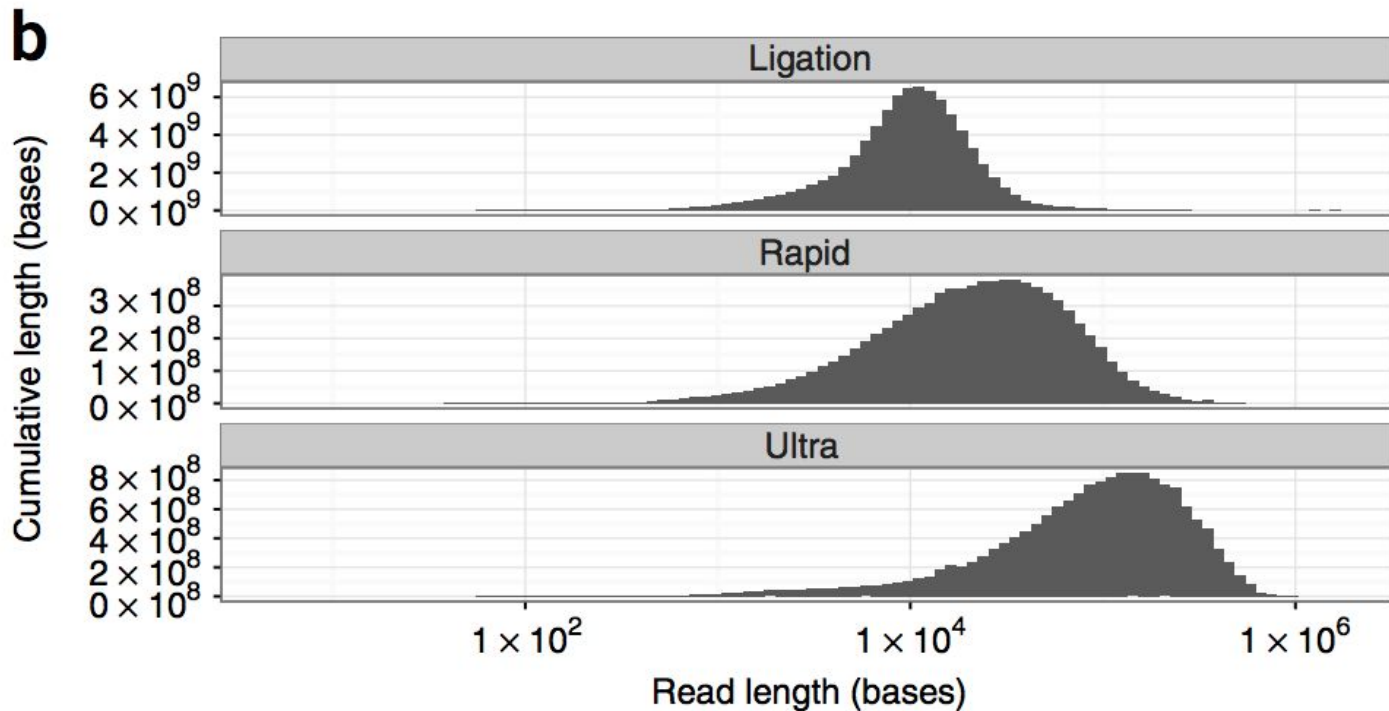
¹UC Santa Cruz Genomics Institute, University of California, Santa Cruz, CA, USA. ²Department of Human Genetics, University of Cambridge, Cambridge, UK. ³Department of Human Genetics, University of Cambridge, Cambridge, UK. ⁴Michael Smith Laboratories and Centre for Genome Sciences and Policy, Genome Sciences Centre, University of British Columbia, Vancouver, BC, Canada. ⁵Institute of Cancer & Genomic Sciences, University of East Anglia, Norwich, UK. ⁶Department of Cancer Research, Toronto, Canada. ⁷Department of Cancer Research, Toronto, Canada. ⁸Department of Cancer Research, Toronto, Canada. ⁹Department of Cancer Research, Toronto, Canada. Correspondence should be addressed to N.J.L.

Received 20 April 2017; accepted 11 December 2017

Same group recently reported a read 2.3 million bases long!

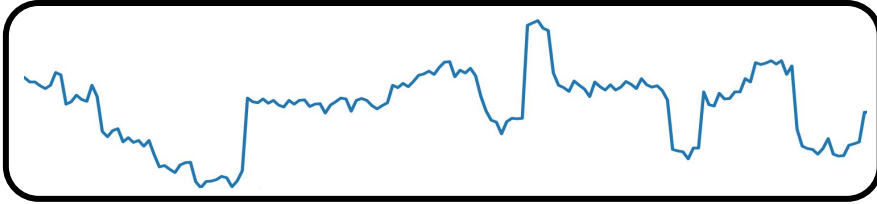
2,272,580 nt, to be exact.

No theoretical upper limit



Nanopore Basecalling

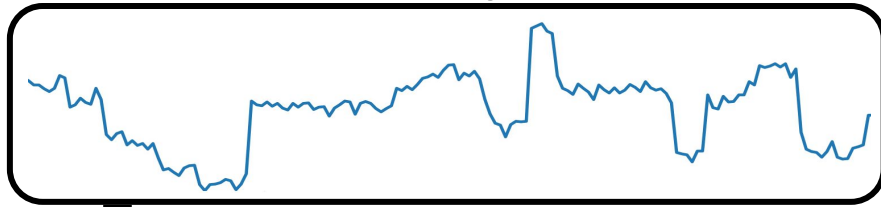
Raw Signal



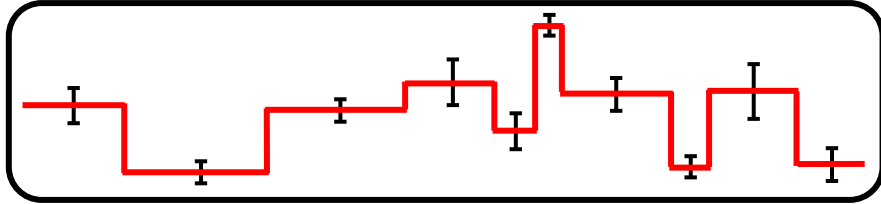
Translation of raw signal
into basepairs

Nanopore Basecalling

Raw Signal



Events



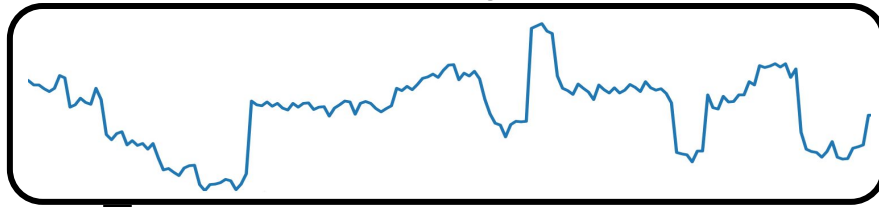
Translation of raw signal into basepairs

Early basecallers began by estimating k-mer boundaries using “events”, which were then input to an HMM

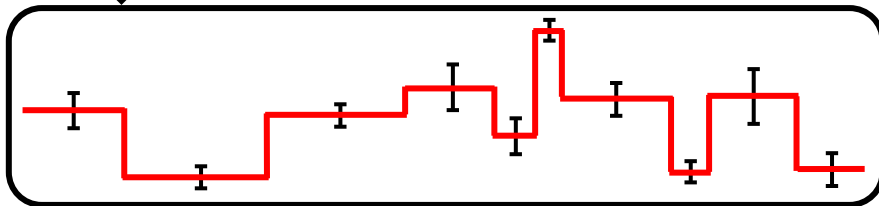
Modern basecallers use neural networks directly on raw signal

Nanopore Basecalling

Raw Signal



Events

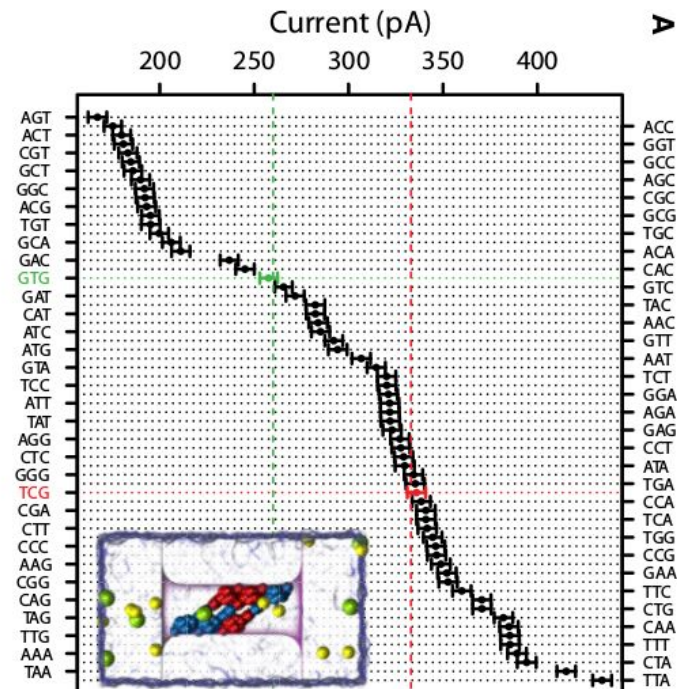


Possible k-mers

0	1	2	3
TCCA	CCAT	CATG	TACA
AGCA	TGGC	TTAC	TCCA
GTCT	ATTA	ACGT	GACG
GATT	ATTG	GTCT	ACGG

(Based on probability of event matches)

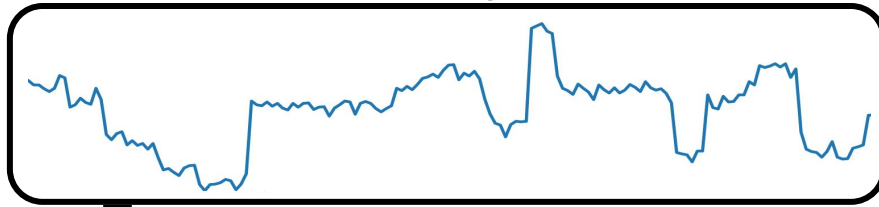
ONT releases k-mer models with expected current distribution of every k-mer



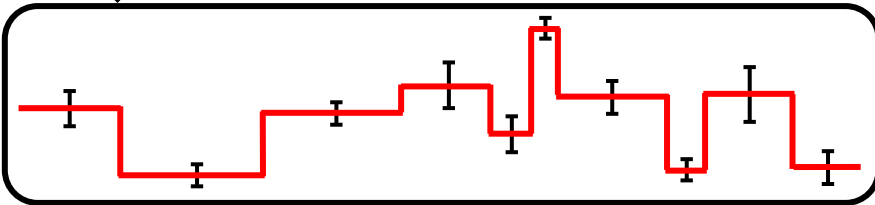
DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling

Raw Signal



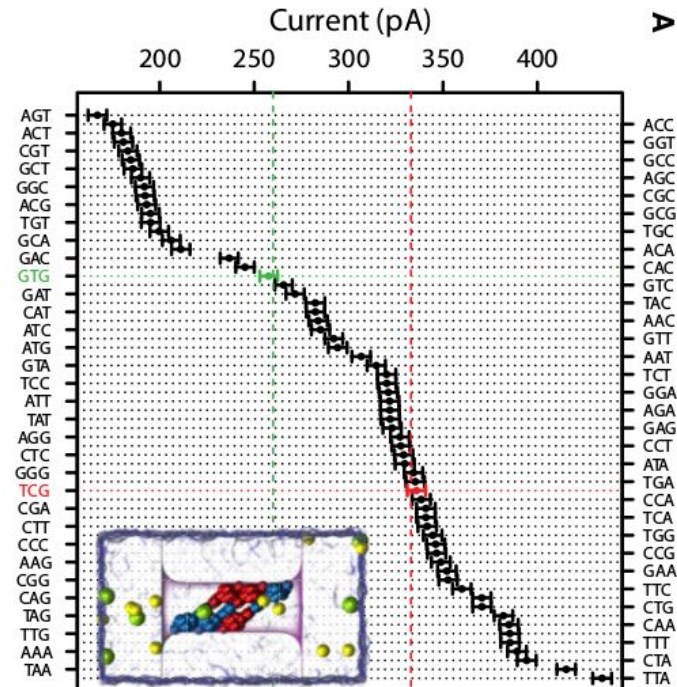
Events



Possible k-mers

0	1	2	3
TCCA →	CCAT	CATG	TACA
AGCA	AGGG	TTAC	TCCA
GTCT	ATTA	ACGT	GACG
GATT →	ATTG	GTCT	ACGG

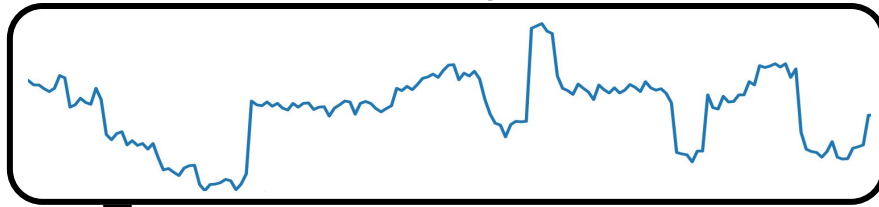
Certain k-mers can be eliminated based on possible transitions



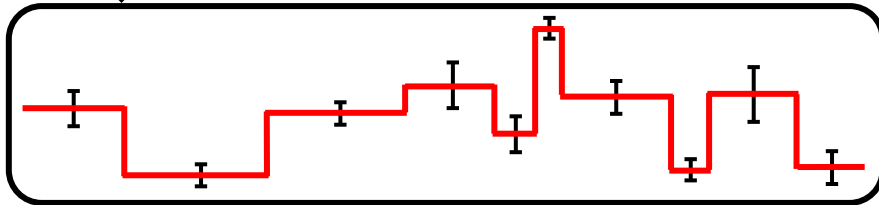
DNA Base-Calling from a Nanopore Using a Viterbi Algorithm
Timp et al. (2012) *Biophysical Journal*

Nanopore Basecalling

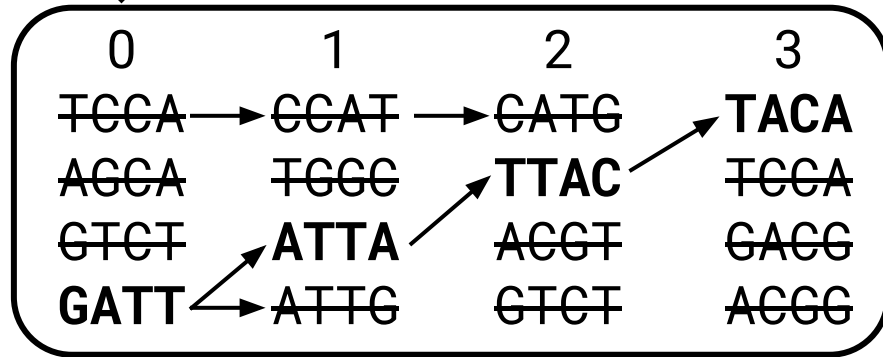
Raw Signal



Events

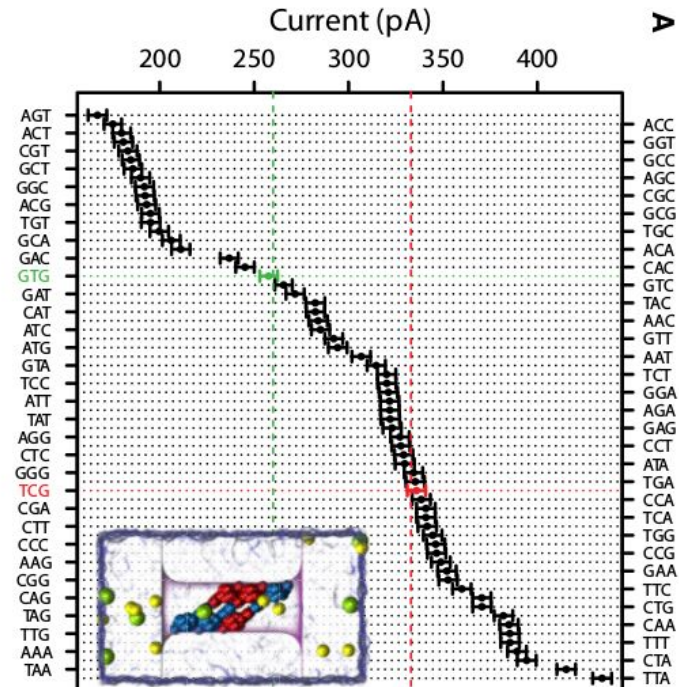


Possible k-mers



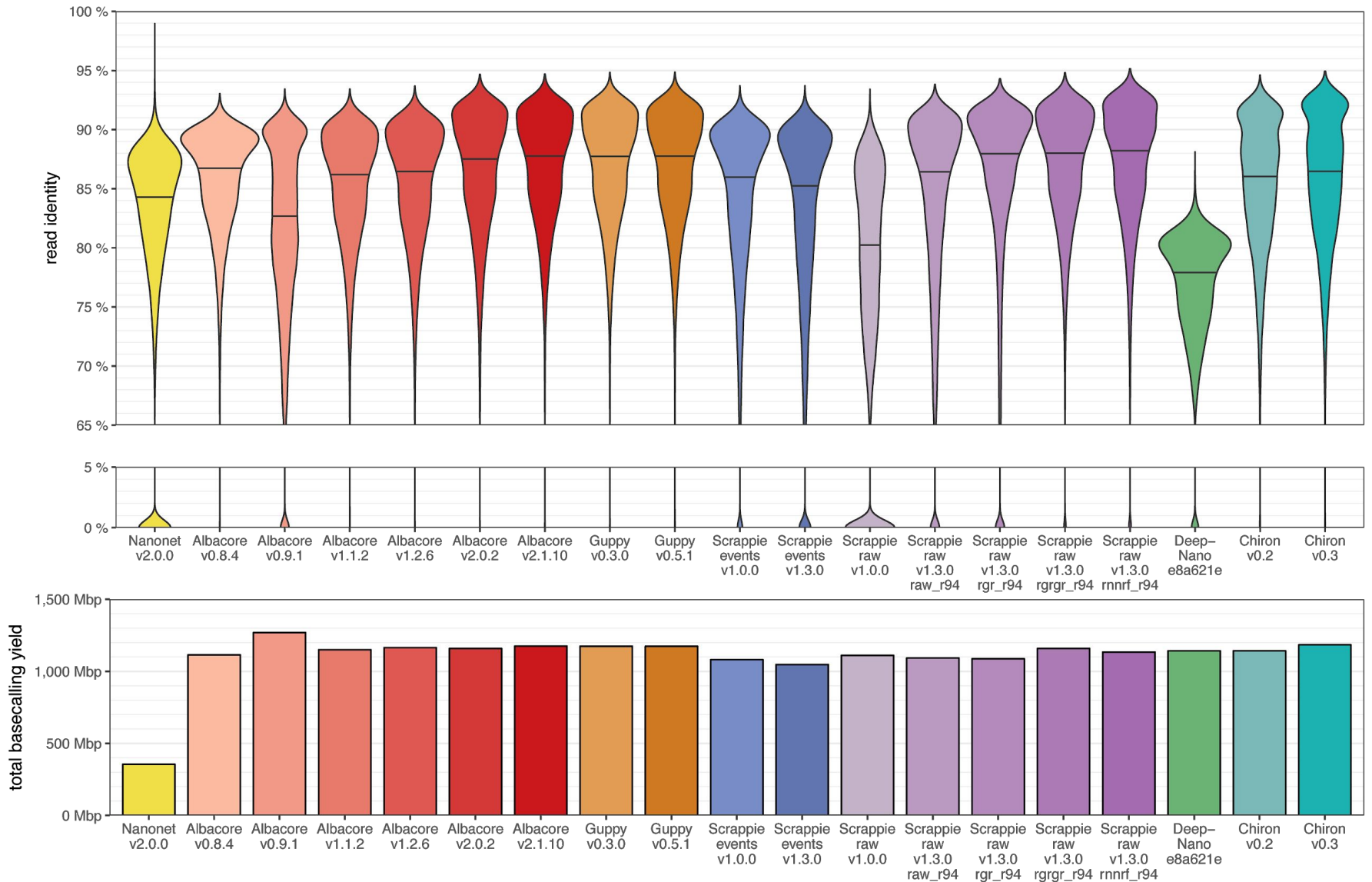
GATTACA

Final sequence determined by most probable k-mers



“DNA Base-Calling from a Nanopore Using a Viterbi Algorithm”
Timp et al. (2012) *Biophysical Journal*

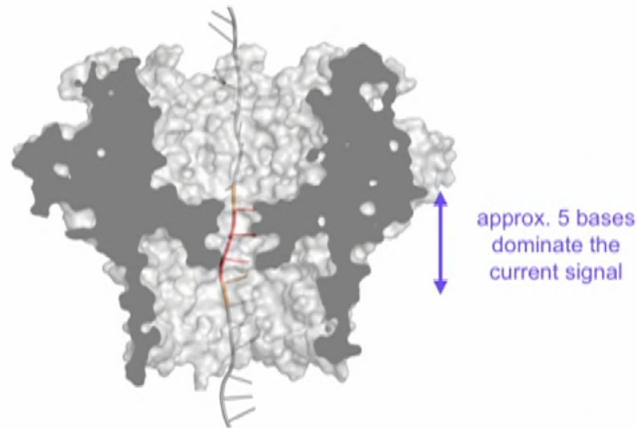
Basecaller Comparison



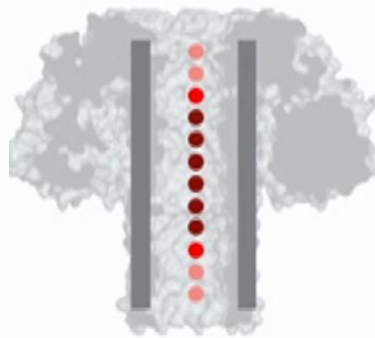
<https://github.com/rrwick/Basecalling-comparison>

New Pore Chemistries

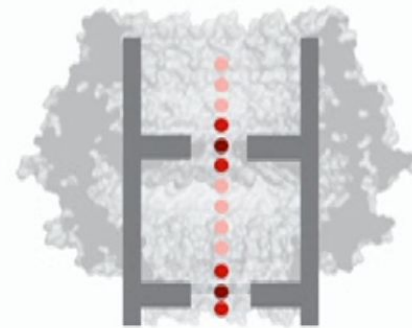
ONT is developing alternate pore chemistries to improve accuracy, particularly for homopolymers



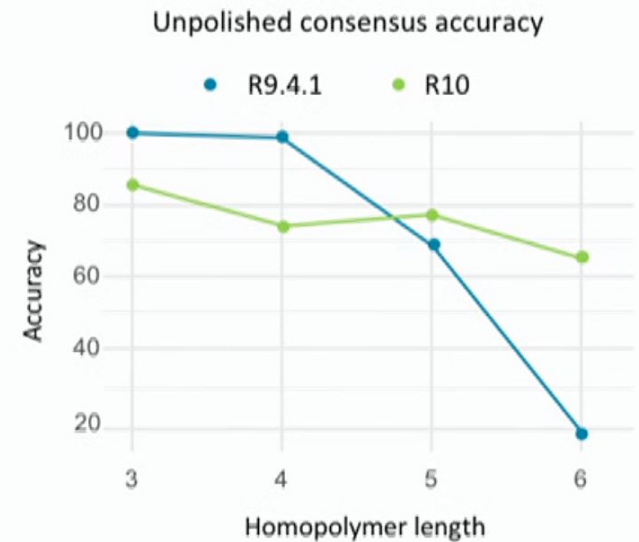
Standard pore
chemistry
“R9”



Pore with long
reader head
Lysenin –
“R8”



Multiple points of
contribution
“R10”



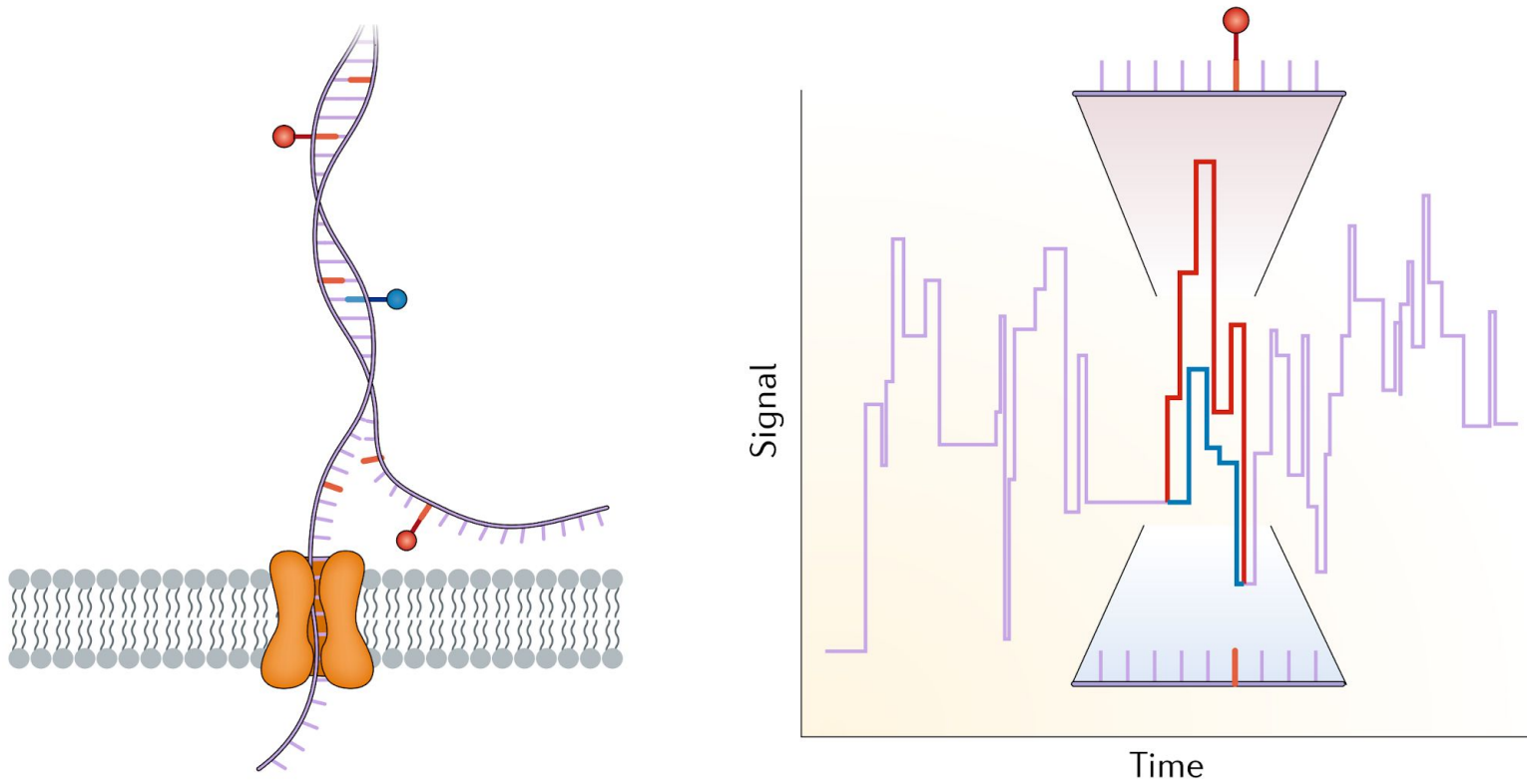
From 2018 London Calling Keynote

<https://vimeo.com/272526835>

DNA Modification Detection

Like PacBio, ONT can detect methylation from raw signal

- Or any other modification that changes ionic current



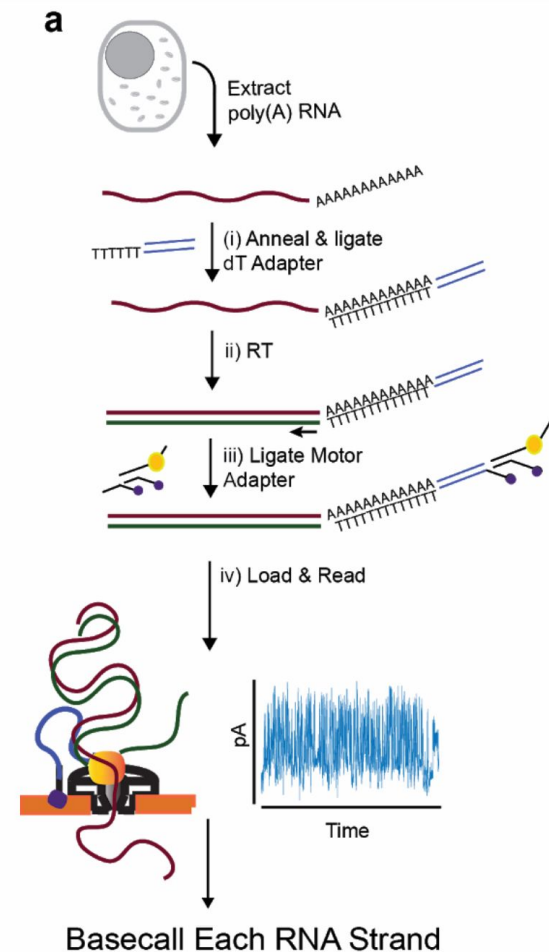
Piercing the dark matter: bioinformatics of long-range sequencing and mapping
Sedlazeck et al. (2018) *Nature Reviews Genetics*. 19:329

Direct RNA-seq

Standard RNA sequencing (RNA-seq) requires creation of complementary DNA (cDNA)

ONT recently introduced direct RNA sequencing

Allows detection of RNA modifications, and potentially secondary structure



Nanopore native RNA sequencing of a human poly(A) transcriptome

Workman et al. *BioRxiv* (<https://www.biorxiv.org/content/10.1101/459529v1>)

ReadUntil Sequencing

ONT machines can stop sequencing a read and immediately start on another in real-time

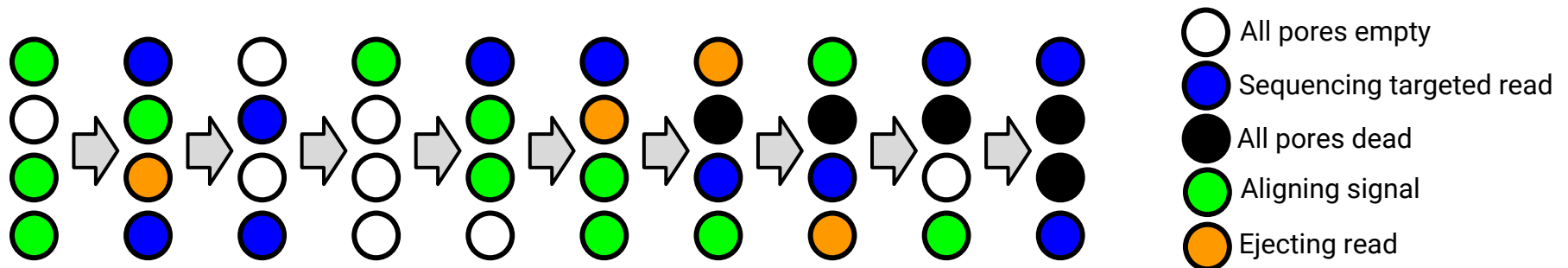
- Each channel has four pores, non-active pores have reads docked

Can potentially avoid sequencing unwanted reads

- For example: reads that align to the human genome, reads that *do not* align to a database of pathogens, reads that align to a region already sequenced to a desired depth

MinION has up to 512 active channels, each reading 450 bp/sec

- Actual number of active channels is variable



TRADEOFFS OF DIFFERENT TECHS.

- Illumina sequencing is cheap & ubiquitous
 - Can sequence *very* deeply, good for measuring abundance
 - Error rate is *very* low, can be good for detection of small variants
 - Can be combined with other technologies (e.g. linked-reads) to provide many different types of information
 - Reads are short (≤ 350 bp each, fragments ≤ 1000 bp), making assembly difficult
 - Library prep (prior to sequencing) can introduce many biases
- ONT sequencing is getting cheaper
 - Can sequence *very* long reads, transformative for assembly
 - Can be good for detection of large (e.g. structural) variants
 - Error rate is *much* higher than with short reads (getting better)
 - Fewer individual reads raise challenges in quantification-related tasks

Work Cited



- Illumina Sequencing Technology. (2010). In Techonlogy Spotlight: Illumina Sequencing. Retrieved July 30, 2014, from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf
- Illumina Solexa Sequencing. (Jan 22, 2010). On YouTube uploaded by Draven1983101. Retrieved July 30, 2014, from <https://www.youtube.com/watch?v=77r5p8lBwJk>
- Overview of Illumina Chemistry. In Massachusetts General Hospital. Retrieved July 30, 2014, from <http://nextgen.mgh.harvard.edu/IlluminaChemistry.html>
- Introduction to Next Generation Sequencing Using the Illumina 1G Genome Analyzer (Solexa). (Jan 31, 2008). Retrieved July 30, 2014, from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf
- Sequencing technology – Past, Present and Future. (2013). Wei Chen. Berlin Institute for Medical Systems Biology. Max-Delbrueck-Center for Molecular Medicine Retrieved July 30, 2014, from http://www.molgen.mpg.de/899148/OWS2013_NGS.pdf
- DNA barcoding. (2014, July 30). In *Wikipedia, The Free Encyclopedia*. Retrieved July 30, 2014, from http://en.wikipedia.org/w/index.php?title=DNA_barcoding&oldid=619163634
- Reference genome. (2014, June 22). In *Wikipedia, The Free Encyclopedia*. Retrieved 19:54, August 4, 2014, from http://en.wikipedia.org/w/index.php?title=Reference_genome&oldid=613984719

Work Cited



- Illumina Sequencing Technology. (2010). In Techonlogy Spotlight: Illumina Sequencing. Retreived July 30, 2014, from http://res.illumina.com/documents/products/techspotlights/techspotlight_sequencing.pdf
- Illumina Solexa Sequencing. (Jan 22, 2010). On YouTube uploaded by Draven1983101. Retrieved July 30, 2014, from <https://www.youtube.com/watch?v=77r5p8lBwJk>
- Overview of Illumina Chemistry. In Massachusetts General Hospital. Retreived July 30, 2014, from <http://nextgen.mgh.harvard.edu/IlluminaChemistry.html>
- Introduction to Next Generation Sequencing Using the Illumina 1G Genome Analyzer (Solexa). (Jan 31, 2008). Retrieved July 30, 2014, from http://research.stowers-institute.org/microscopy/external/PowerpointPresentations/ppt/Methods_Technology/KSH_Tech&Methods_012808Final.pdf
- Sequencing technology – Past, Present and Future. (2013). Wei Chen. Berlin Institute for Medical Systems Biology. Max-Delbrueck-Center for Molecular Medicine Retrieved July 30, 2014, from http://www.molgen.mpg.de/899148/OWS2013_NGS.pdf
- DNA barcoding. (2014, July 30). In *Wikipedia, The Free Encyclopedia*. Retrieved July 30, 2014, from http://en.wikipedia.org/w/index.php?title=DNA_barcoding&oldid=619163634
- Reference genome. (2014, June 22). In *Wikipedia, The Free Encyclopedia*. Retrieved 19:54, August 4, 2014, from http://en.wikipedia.org/w/index.php?title=Reference_genome&oldid=613984719