# NEXT GENERATION SEQUENCING (NGS)

Andreas Gisel

International Institute of Tropical Agriculture (IITA)
Ibadan, Nigeria

a.gisel@cgiar.org

☐ Short history bio-sequencing
☐ NGS Sequencing technologies
☐ NGS Sequence data
☐ NGS Applications
☐ NGS Data Analysis

# OVERVIEW

First fully sequenced bio-sequence

HISTORY

First fully sequenced bio-sequence
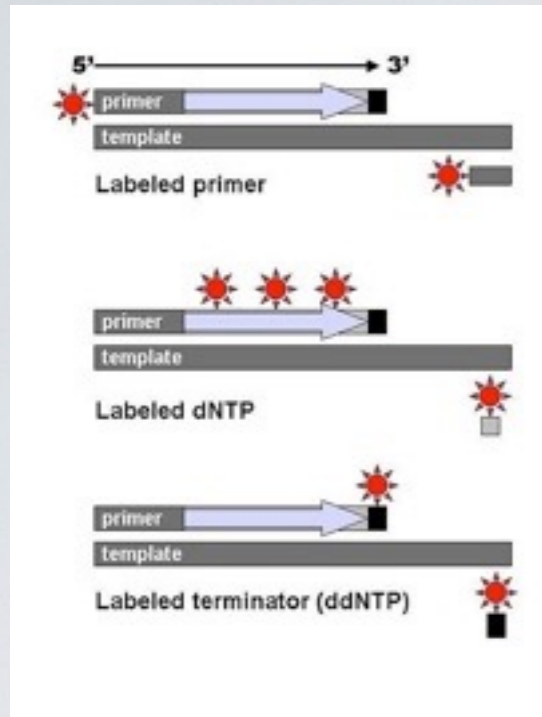- amino acid of insulin (51aa) 1955

# HISTORY

☐ First fully sequenced bio-sequence
   - amini acid of insulin (51aa) 1955
☐ First fully sequence nucleic acid
   - tRNA (75nt) 1965
☐ First DNA
   - Bacteriophage (5375nt) 1977

☐ DNA sequencing
   - Sanger sequencing technology (1975)

# HISTORY

■ sequencing by chain-termination method
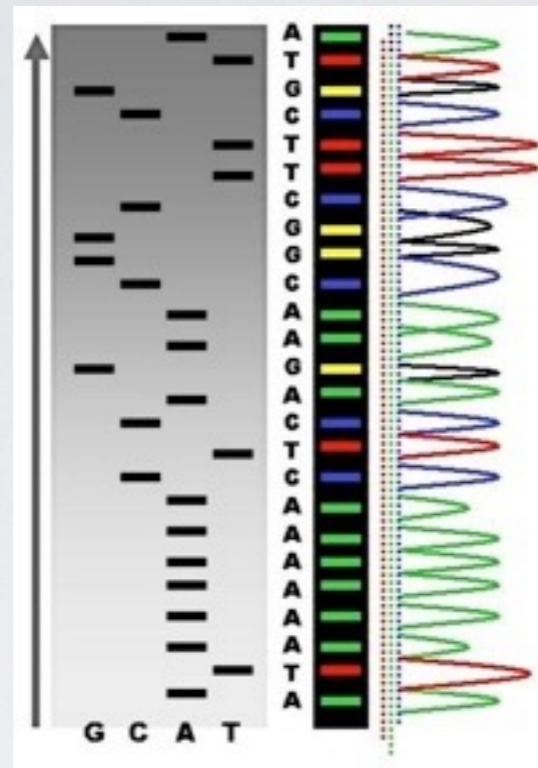
# SANGER TECHNOLOGY

# sequencing by chain-termination method



# SANGER TECHNOLOGY

# sequencing by chain-termination method



# SANGER TECHNOLOGY

# sequencing by chain-termination method
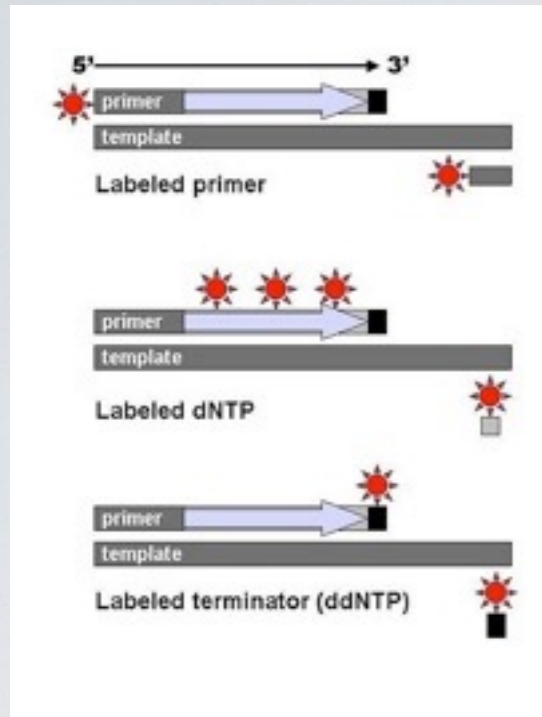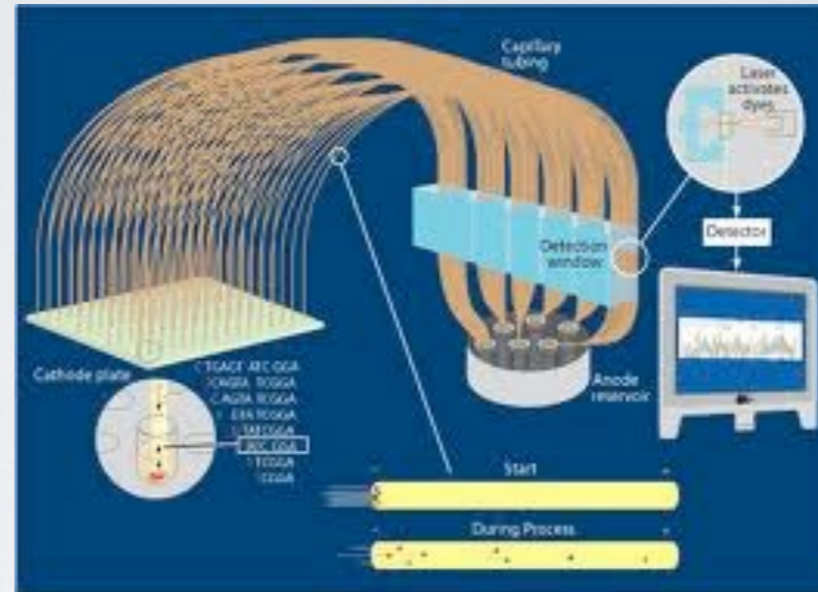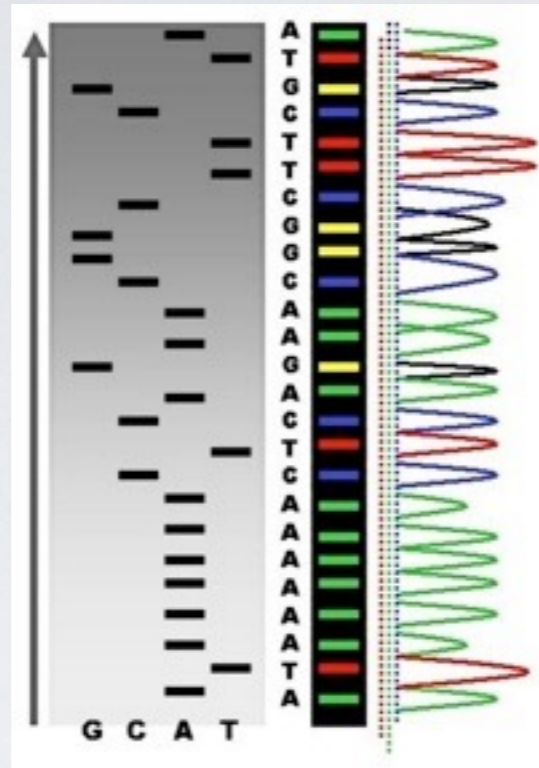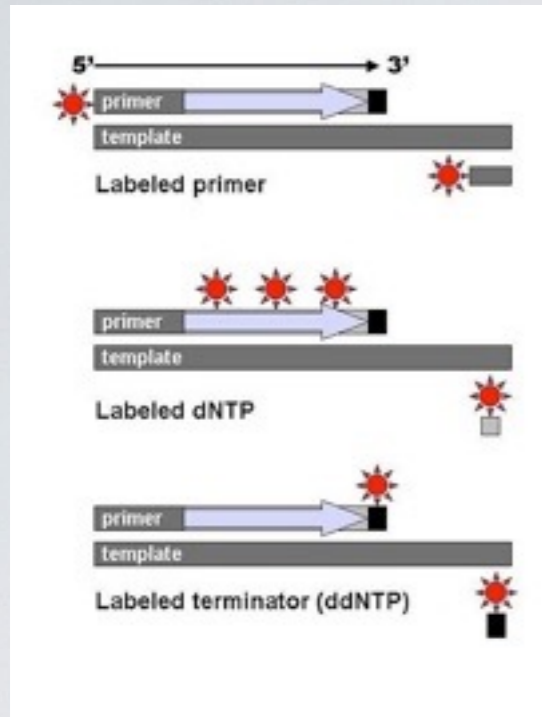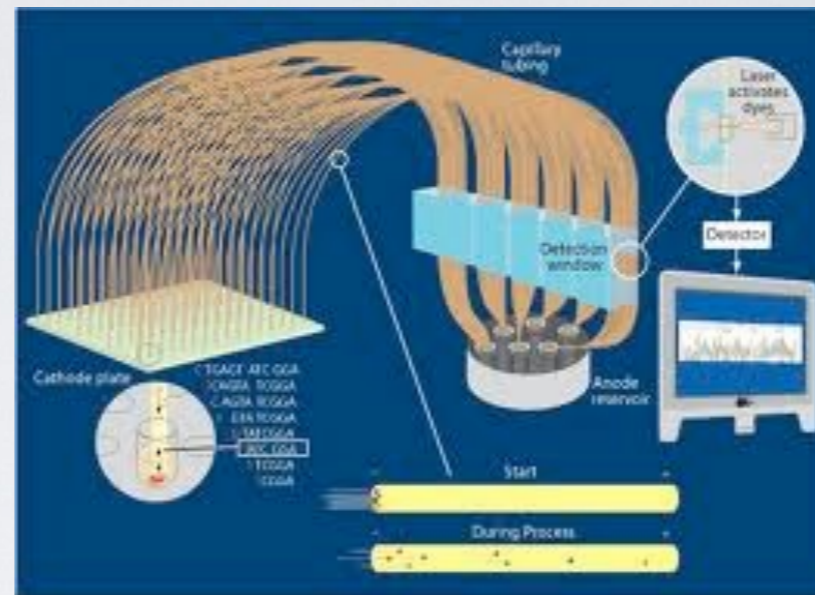


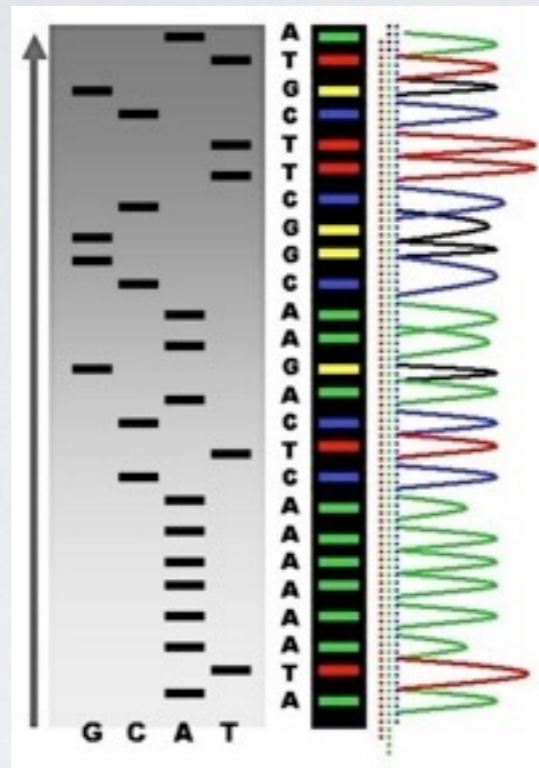# SANGER TECHNOLOGY

# sequencing by chain-termination method



# SANGER TECHNOLOGY

□ sequencing by chain-termination method



□ DNA sequencing by capillary electrophoresis
□ 384 reactions in parallel
□ sequences up to 1000nt

# SANGER TECHNOLOGY

# NEXT GENERATION SEQUENCING

7

☐ Sequencing by synthesis

☐ highly parallelized sequencing

☐ Paired-end sequencing

# NEXT GENERATION SEQUENCING

- Sequencing by synthesis

- highly parallelized sequenc[ing]

- Paired-end sequencing

# NEXT GENERATION SEQUENCING

☐ Sequencing by synthesis

☐ highly parallelized sequencing

☐ Paired-end sequencing

# NEXT GENERATION SEQUENCING

- Sequencing by synthesis

- highly parallelized sequencing

- Paired-end sequencing

# NEXT GENERATION SEQUENCING

☐ Sequencing by synthesis

☐ highly parallelized sequencing

☐ Paired-end sequencing

# NEXT GENERATION SEQUENCING

Amplification steps
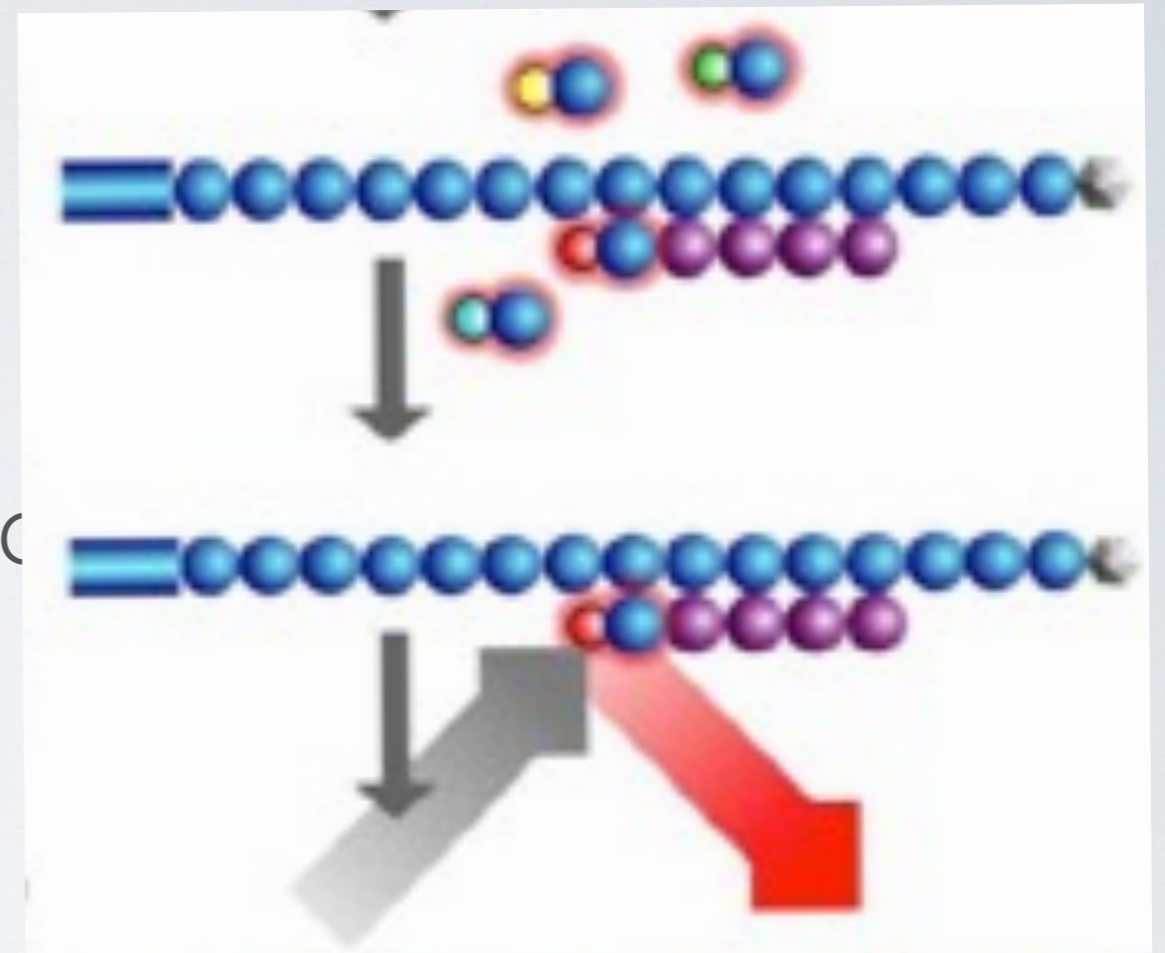☐ 454 Roche
☐ Solexa Illumina
☐ SOLiD Applied Biosystems

Single molecule
☐ Pacific BioSciences
☐ Ion Torrent
☐ Nanopore

# TECHNOLOGIES

# 454 ROCHE

☐ DNA template immobilized to nano-beads
☐ Emulsion PCR
☐ Pyro-Sequencing in nano wells (1.6M reads)
☐ Sequencing by synthesis

# 454 ROCHE

☐ DNA template immobilized to nano-beads
☐ Emulsion PCR
☐ Pyro-Sequencing in nano wells (1.6M reads)
☐ Sequencing by synthesis



454 ROCHE

- [x] DNA template immobilized to nano-beads
- [x] Emulsion PCR
- [x] Pyro-Sequencing in nano wells (1.6M reads)
- [x] Sequencing by synthesis



# 454 ROCHE

- DNA template immobilized to nano-beads
- Emulsion PCR
- Pyro-Sequencing in nano wells (1.6M reads)
- Sequencing by synthesis



# 454 ROCHE

Sequence fragmentation

Ligation of adaptors

Sequence immobilization

emulsion PCR

distribution in nano wells

454 ROCHE

☐ sequencing length up to 1000nt (800nt)
☐ up to 1.2M reads
☐ 600 - 800Mb per run
☐ problems with homo polymers



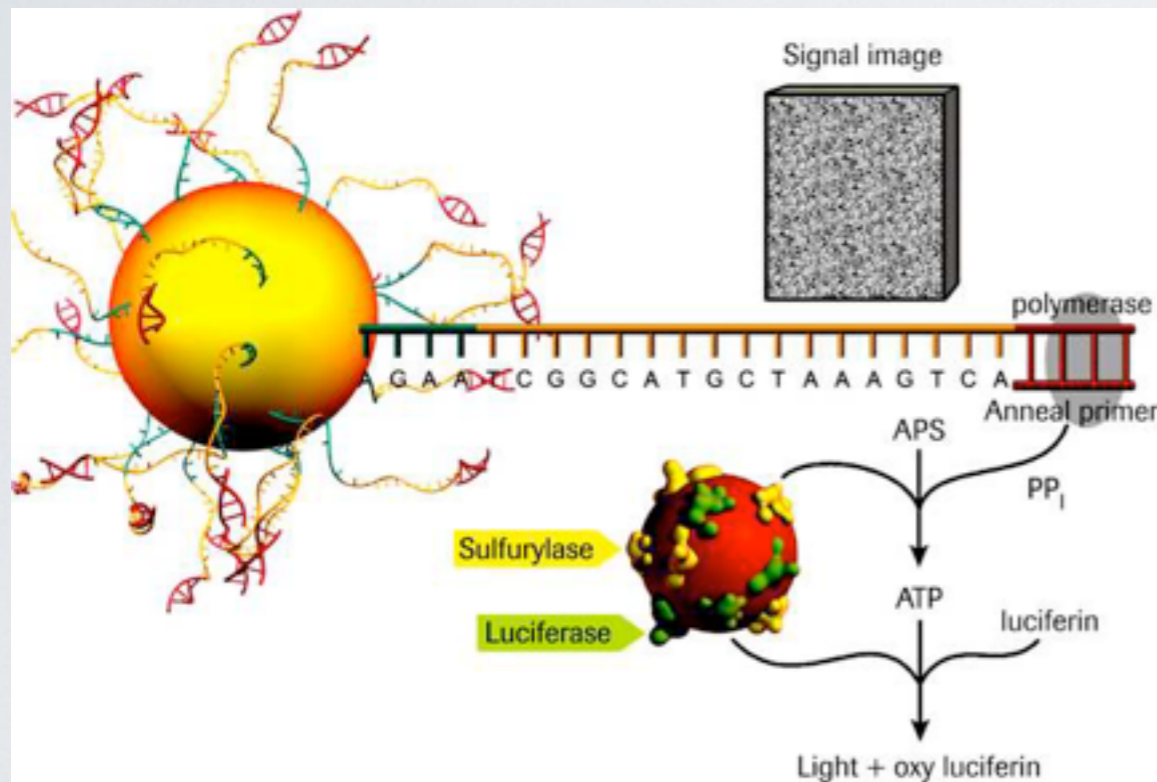| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 | 99,37 | 0,62 | 0,01 | 0 | 0 | 0 | 0 |
| 2 | 0 | 99,74 | 0,26 | 0 | 0 | 0 | 0 |
| 3 | 0 | 2,91 | 95,62 | 1,45 | 0,03 | 0 | 0 |
| 4 | 0 | 0,05 | 9,69 | 89,63 | 0,56 | 0,07 | 0 |
| 5 | 0 | 0 | 0,85 | 29,99 | 68,34 | 0,8 | 0,02 |
| 6* | 0.00 | 0.00 | 0.00 | 5.34 | 67.37 | 26.04 | 1.22 |

Vicarico et al pers comm.

# 454 ROCHE

# ILLUMINA

12

- DNA template immobilized to a flow cell
- Cluster formation by Bridge PCR
- Sequencing on flow cell (3000M reads)
- Sequencing by synthesis (protected nts)

# ILLUMINA

☐ DNA template immobilized to a flow cell
☐ Cluster formation by Bridge PCR
☐ Sequencing on flow cell (3000M reads)
☐ Sequencing by synthesis (protected nts)



# ILLUMINA

☐ DNA template immobilized to a flow cell
☐ Cluster formation by Bridge PCR
☐ Sequencing on flow cell (3000M reads)
☐ Sequencing by synthesis (protected nts)

>100M single
molecules

>100M single
clusters

# ILLUMINA

☐ DNA template immobilized to a flow cell
☐ Cluster formation by Bridge PCR
☐ Sequencing on flow cell (3000M reads)
☐ Sequencing by synthesis (protected nts)



>100M single molecules → >100M single clusters



Sequencing

Add 4 Fl-NTP's + Polymerase — Incorporated Fl-NTP is imaged — Terminator and fluorescent dye are cleaved from the Fl-NTP

X 36 - 75

T G C G A A T T

illumina

# ILLUMINA

- DNA template immobilized to a flow cell
- Cluster formation by Bridge PCR
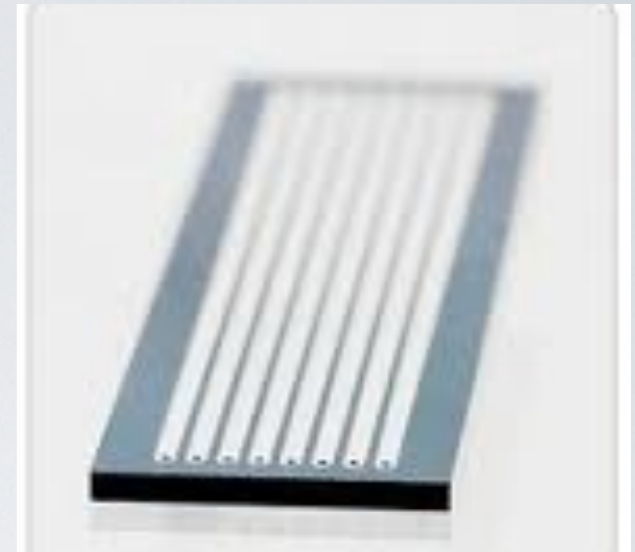- Sequencing on flow cell (3000M reads)
- Sequencing by synthesis (protected nts)



>100M single molecules

>100M single clusters



Sequencing

Add 4 Fl-NTP's + Polymerase

Incorporated Fl-NTP is imaged

Terminator and fluorescent dye are cleaved from the Fl-NTP

X 36 - 75

T G C G A A T T

illumina



# ILLUMINA

Fragmentation and adaptor ligation

Immobilization and strand synthesis

Bridge PCR to form clusters

# ILLUMINA

- sequencing length up to 250nt
- up to 3000M sequences - high coverage
- 400 - 600Gb per run
- sequence size limitation



ILLUMINA

# SOLID

☐ DNA template immobilized to a nano bead
☐ Cluster formation by emulsion PCR
☐ Sequencing on flow cell (4800M reads)
☐ Sequencing by ligation

# SOLID

- DNA template immobilized to a nano bead
- Cluster formation by emulsion PCR
- Sequencing on flow cell (4800M reads)
- Sequencing by ligation



# SOLID

- DNA template immobilized to a nano bead
- Cluster formation by emulsion PCR
- Sequencing on flow cell (4800M reads)
- Sequencing by ligation

# SOLID

- DNA template immobilized to a nano bead
- Cluster formation by emulsion PCR
- Sequencing on flow cell (4800M reads)
- Sequencing by ligation



SOLID
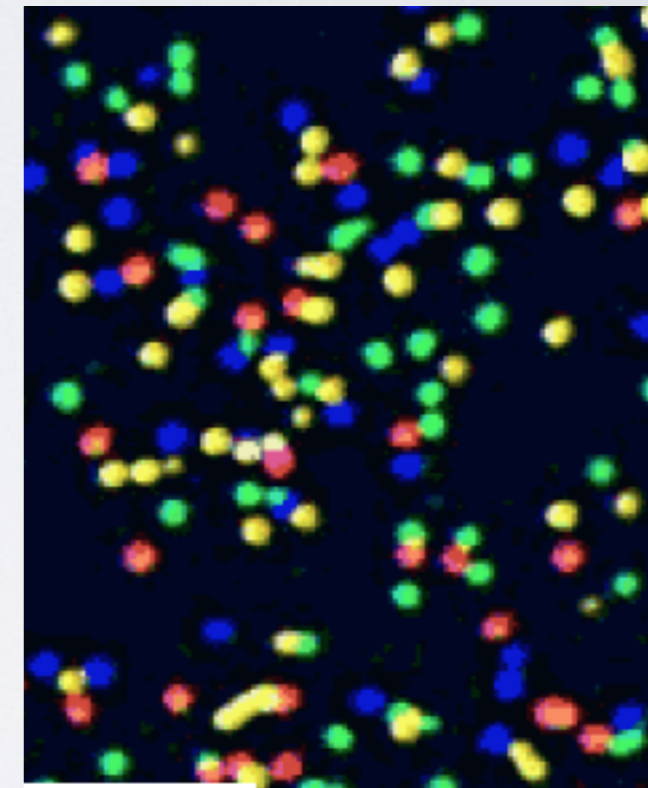
- DNA template immobilized to a nano bead
- Cluster formation by emulsion PCR
- Sequencing on flow cell (4800M reads)
- Sequencing by ligation



SOLID

- DNA template immobilized to a nano bead
- Cluster formation by emulsion PCR
- Sequencing on flow cell (4800M reads)
- Sequencing by ligation

Fragmentation, adaptors and immobilization

Emulsion PCR and bead separation

Bead deposition

# SOLID

- sequencing length up to 75nt
- up to 4800M sequences - high coverage
- - 300 Gb per run
- sequence size limitation
- several sequencing rounds
- Every base is called twice

# SOLID

- sequencing length up to 75nt
- up to 4800M sequences - high coverage
- - 300 Gb per run
- sequence size limitation
- several sequencing rounds
- Every base is called twice

Ligation cycle   1    2    3    4    5    6    7 ... (n cycles)

# SOLID

- sequencing length up to 75nt
- up to 4800M sequences - high coverage
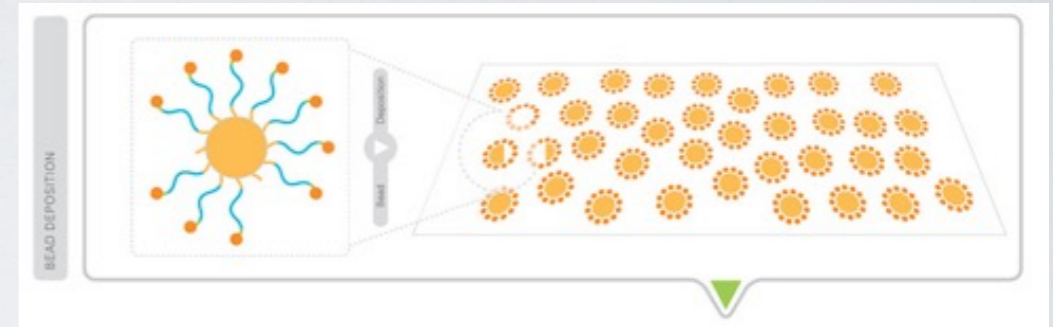- - 300 Gb per run
- sequence size limitation
- several sequencing rounds
- Every base is called twice



# SOLID

- [ ] Sequence fragmentation
- [ ] Adaptor ligation
- [ ] Sequence immobilization
- [ ] PCR amplification (emulsion or bridge PCR)
- [ ] Real-time sequencing (by synthesis or ligation)
- [ ] huge amount of short sequence reads
- [ ] high coverage
- [ ] difficulties with assembling

# SUMMARY

# PACIFIC BIOSCIENCES

- Polymerase immobilized in a nano well
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (75K reads)
- Sequencing by synthesis (fluorescence)
- Read length up to 10000nt average >1000
- Fast sample preparation and sequencing (8h)

# PACIFIC BIOSCIENCES

- Polymerase immobilized in a nano well
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (75K reads)
- Sequencing by synthesis (fluorescence)
- Read length up to 10000nt average >1000
- Fast sample preparation and sequencing (8h)



Zero-mode
waveguide

# PACIFIC BIOSCIENCES

☐ Polymerase immobilized in a nano well
☐ NO amplification (true single molecule sequencing)
☐ Sequencing on flow cell (75K reads)
☐ Sequencing by synthesis (fluorescence)
☐ Read length up to 10000nt average >1000
☐ Fast sample preparation and sequencing (8h)

Zero-mode waveguide

# PACIFIC BIOSCIENCES

# ION TORRENT

- Immobilized in a nano well (semiconductor)
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (1M reads)
- Sequencing by synthesis (H+ release)
- Read length up to 10000nt average >400
- Fast sample preparation and sequencing (8h)

# ION TORRENT

- Immobilized in a nano well (semiconductor)
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (1M reads)
- Sequencing by synthesis (H+ release)
- Read length up to 10000nt average >400
- Fast sample preparation and sequencing (8h)



Micro-machined wells

Ion-sensitive layer

Proprietary Ion sensor

# ION TORRENT

Immobilized in a nano well (semiconductor)
NO amplification (true single molecule sequencing)
Sequencing on flow cell (1M reads)
Sequencing by synthesis (H+ release)
Read length up to 10000nt average >400
Fast sample preparation and sequencing (8h)

Micro-machined wells

Two bases
are incorporated

Two hydrogen ions
are released

H+ H+

Ion-sensitive layer

Proprietary Ion sensor

T  G  A  C  TT
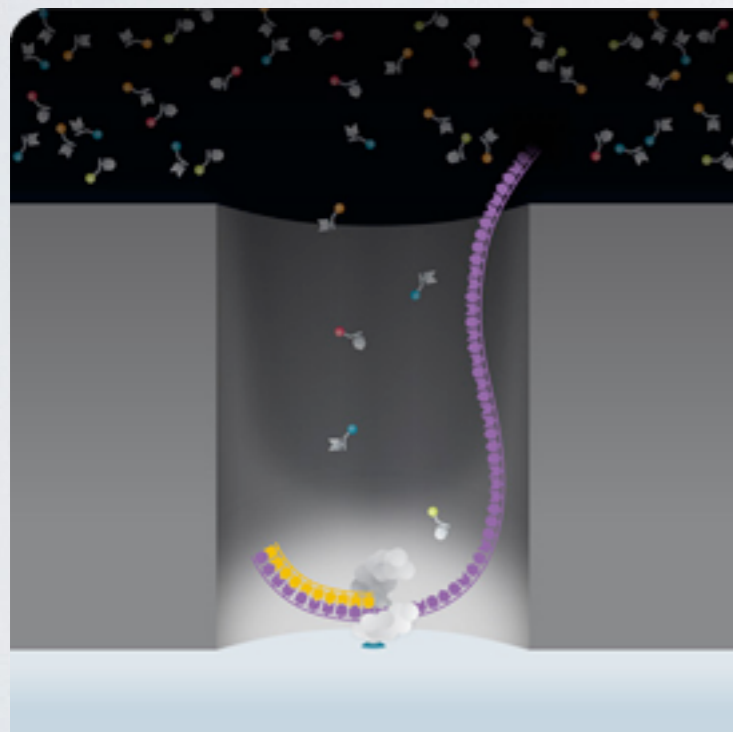
# ION TORRENT

- Immobilized in a nano well (semiconductor)
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (5M reads)
- Sequencing by synthesis (H+ release)

| Chip | Expected Sequencing Run Time | | | Expected Output | | |
|---|---|---|---|---|---|---|
| | 35 base reads | 100 base reads | 200 base reads | 35 base reads | 100 base reads | 200 base reads |
| Ion 314™ Chip | 0.5 hr | 1.5 hr | 2.4 hr | 3 Mb | 10 Mb | 20 Mb |
| Ion 316™ Chip | 0.7 hr | 1.7 hr | 3.1 hr | 30 Mb | 100 Mb | 200 Mb |
| Ion 318™ Chip | 0.9 hr | 2.4 hr | 4.5 hr | 300 Mb | 500 Mb | 1 GB |

(8h)

# ION TORRENT

# NANOPORE

- ☐ NO immobilization
- ☐ NO amplification (true single molecule sequencing)
- ☐ Sequencing through solid-state nanopore
- ☐ Sequencing by current disruption (8000 pores)
- ☐ Read length up to 100000nt
- ☐ in future 20 pores sequence human genome in 15min

# NANOPORE

- NO immobilization
- NO amplification (true single molecule sequencing)
- Sequencing through solid-state nanopore
- Sequencing by current disruption (8000 pores)
- Read length up to 100000nt
- in future 20 pores sequence human genome in 15min



# NANOPORE

- NO immobilization
- NO amplification (true single molecule sequencing)
- Sequencing through solid-state nanopore
- Sequencing by current disruption (8000 pores)
- Read length up to 100000nt
- in future 20 pores sequence human genome in 15min



# NANOPORE

- NO immobilization
- NO amplification (true single molecule sequencing)
- Sequencing through solid-state nanopore
- Sequencing by current disruption (8000 pores)
- Read length up to 100000nt
- in future 20 pores sequence human genome in 15min



current

NANOPORE

doi:10.1038/nature.2012.10051

- NO immobilization
- NO amplification (true single molecule sequencing)
- Sequencing through solid-state nanopore
- Sequencing by current disruption (8000 pores)
- Read length up to 100000nt
- in future 20 pores sequence human genome in 15min

# NANOPORE

☐ NO immobilization
☐ NO amplification (true single molecule sequencing)
☐ Sequencing through solid-state nanopore
☐ Sequencing by current disruption (8000 pores)
☐ Read length up to 100000nt
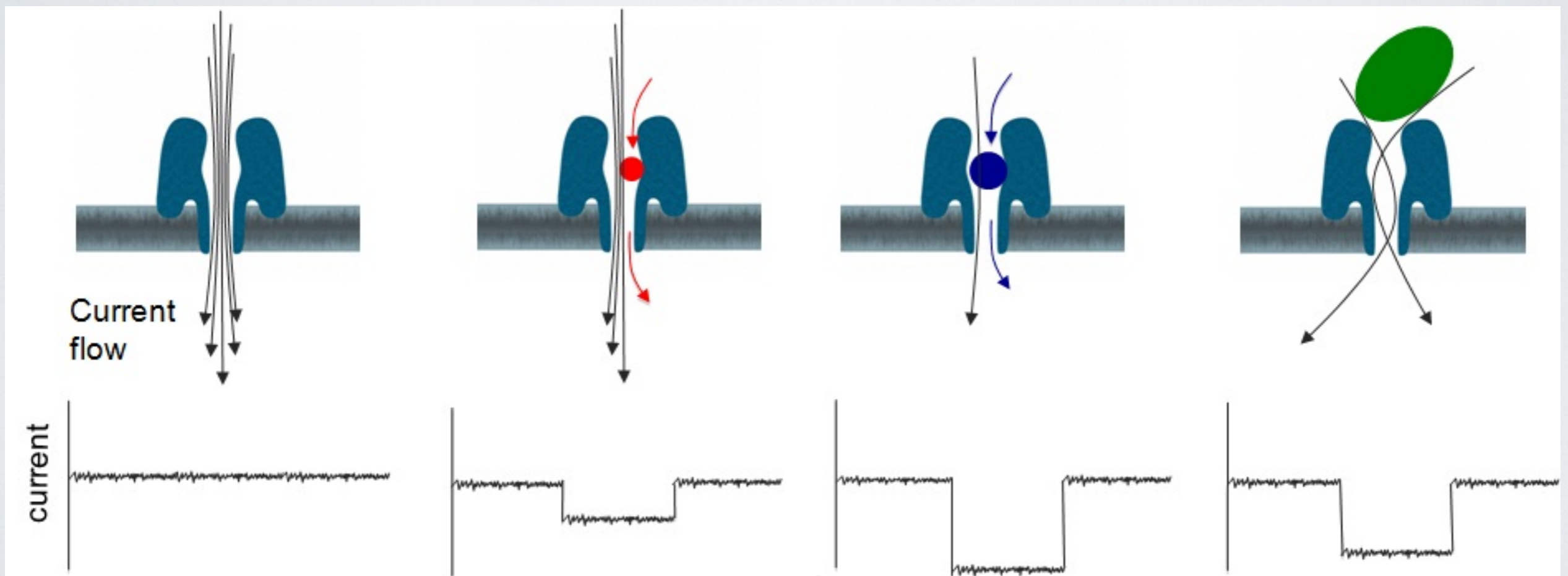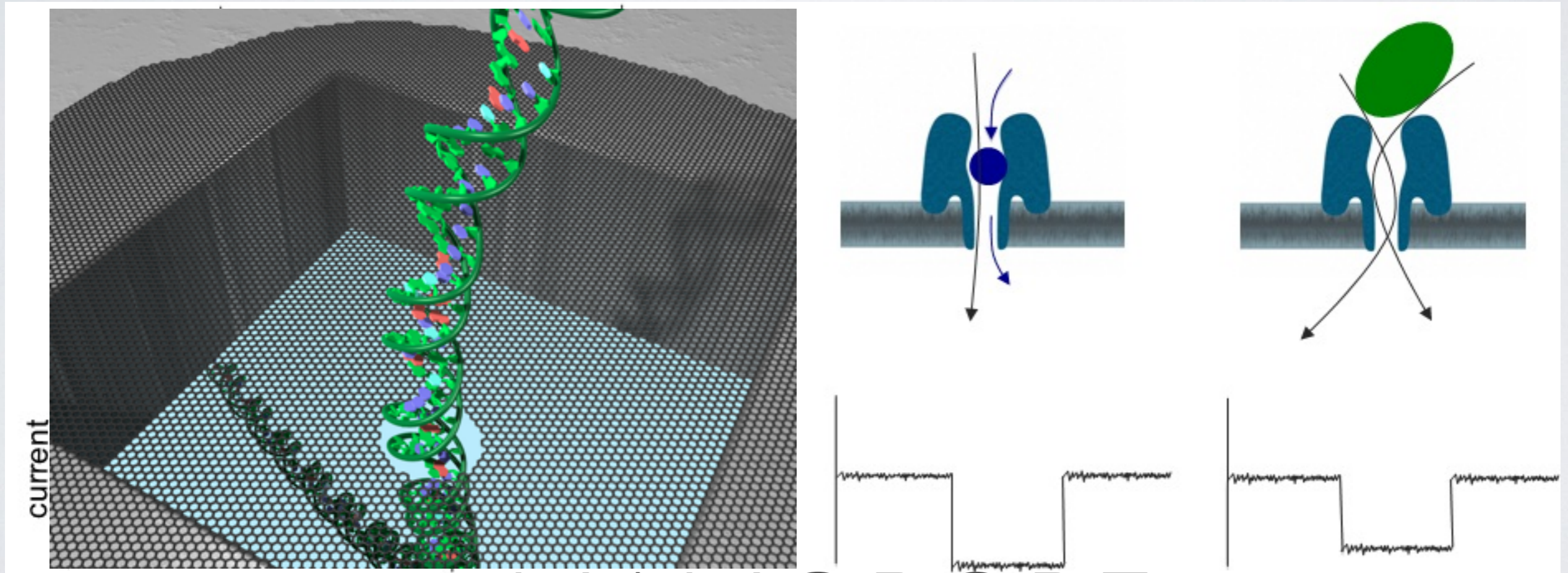☐ in future 20 pores sequence human genome in 15min
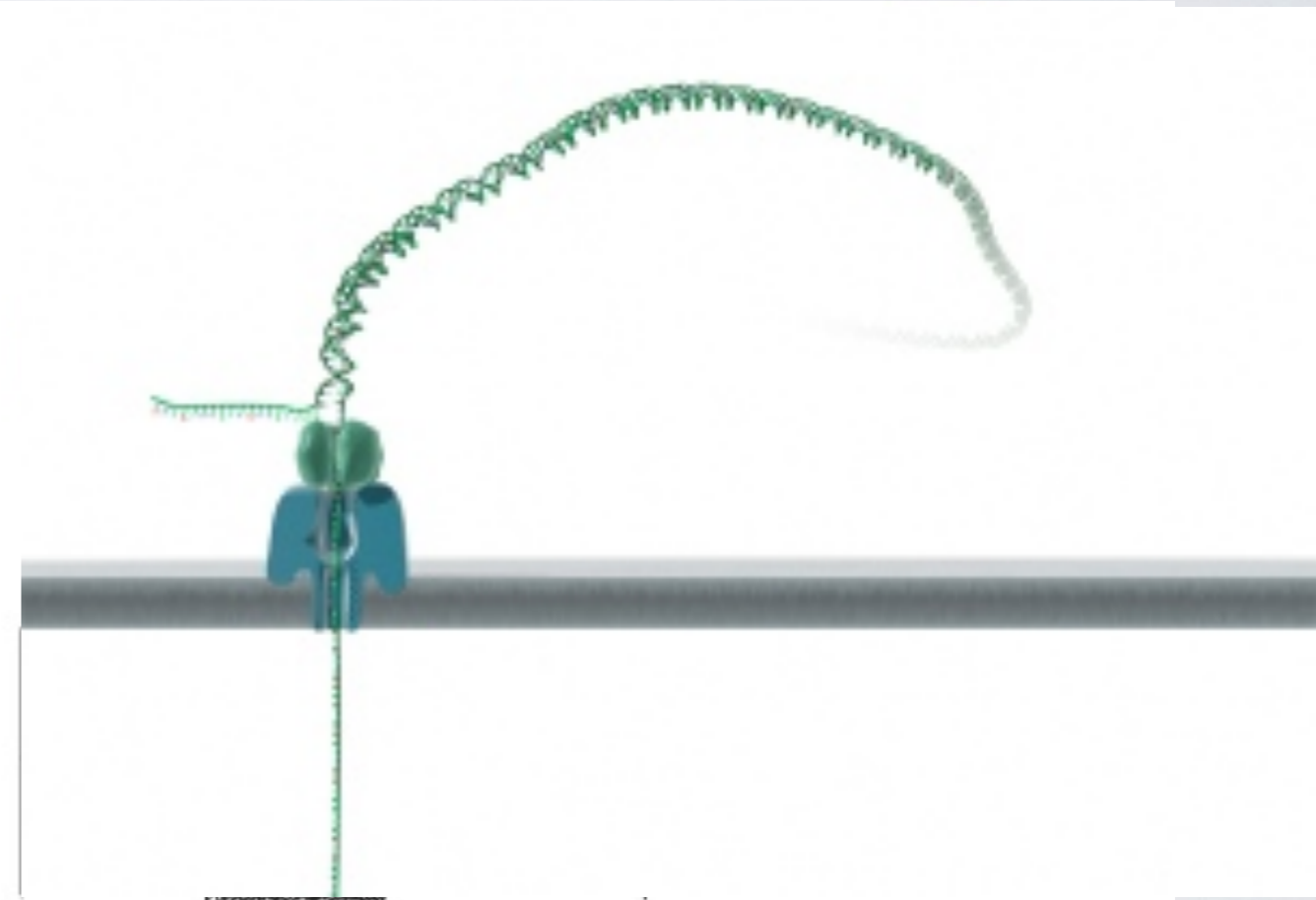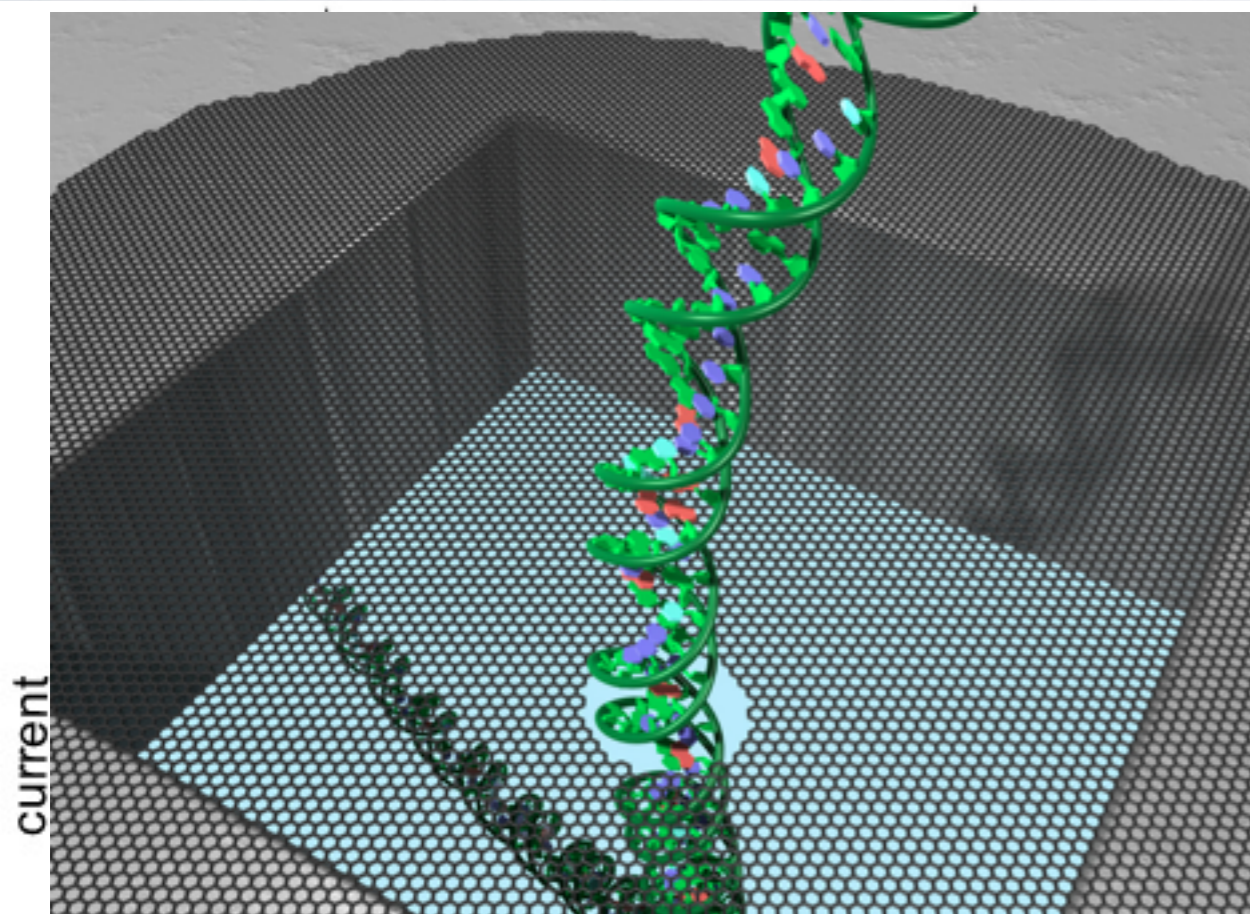


# NANOPORE

- NO immobilization
- NO amplification (true single molecule sequencing)
- Sequencing through solid-state nanopore
- Sequencing by current disruption (8000 pores)
- Read length up to 100000nt
- in future 20 pores sequence human genome in 15min



# NANOPORE

# DATA FORMAT

# DATA FORMAT

# ▢ Illumina

- SCARF (s_*_sequence.txt): Solexa Compact ASCII Read Format

```
HWI-EAS255_4_FC2010Y
1:43:110:790:TTAATCTACAGAATAGATAGCTAGCATATATTT:IIIIIIIIIIIIIIIAIIIIIIIII&;II&,I
HWI-EAS255_4_FC2010Y:
1:43:122:836:GATCGGAAGGCTCGTATGCCGTCTTCTTCTTTT:IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
```

- FASTQ (*.fastq)

```
@HWI-EAS255_4_FC2010Y_1_43_110_790
TTAATCTACAGAATAGATAGCTAGCATATATTT
+HWI-EAS255_4_FC2010Y_1_43_110_790
IIIIIIIIIIIIIIIIAIIIIIIIII&;II&,I
```

# DATA FORMAT

## ■ SOLiD

```
>1_51_64_F3
T1030103123033323320333000021122223
>1_51_127_F3
T2010323233203132310110100200310102
```

- QUAL (xxxx_.QV.qual):

```
>1_51_64_F3
12 7 21 16 6 2 25 5 25 26 6 7 2 8 5 2 3 2 6 21 5 2 3 9 4 2 2 2 17 6 2 2 2 5 3
>1_51_127_F3
3 18 15 4 11 2 6 4 4 6 2 7 2 9 4 3 2 6 18 2 2 4 3 2 2 2 2 2 2 4 2 3 4 4 2
```

# DATA FORMAT

# 454 Roche

- Roche 454 SFF Standard Flowgram Format (*.sff)

- FASTA (*.fna)

```
>E6PIHNP01B74B0
AACAGGATTAGCAGAGCGAGGTATGTAGGCGGTGCTACAGAGTTCTTGAAGTGGTGGCCTAACTACGGCTACACTAG
AAGAACAGTATTTGGTATCTGCGCTCTGCTGAAGCCAGTTACCTTCGGAAAAAGAGTTGGTAGCTCTTGATCCGGC
```

- QUAL (*.qual)

```
>E6PIHNP01B74B0
34 27 28 26 34 28 28 35 28 25 28 28 28 28 28 27 28 28 28 32 25 28 28 25 27 27
27 31 22 28 31 24 28 27 27 27 27 27 25 28 27 28 34 26 27 32 25 27 31 22 25 24
28 20 27 31 23 33 25 27 32 25 22 28 28 27 34 27 27 24 27 25 25 25 25 25 27 31
24 27 26 17 23 15 28 25 28 36 32 13 34 28 22 26 26 27 28 27 27 27 17 20 28 27
28 27 27 24 34 28 27 32 27 28 26 33 27 27 34 28 35 28 28 34 27 39 35 24 14 4
27 25 24 34 28 35 28 26 27 27 18 18 17 31 26 27 25 28 27 18 29 21 28
```

# DATA FORMAT

Quality scores are currently calculated to reliably call bases from a Sanger chromatogram; well-known as Phred scores.

They range from 0 to 93 (Illumina 0 - 40), even though rarely exceed 60; represented by ASCII code.

# QUALITY SCORE

Quality scores are currently calculated to reliably call bases from a Sanger chromatogram; well-known as Phred scores.

They range from 0 to 93 (Illumina 0 - 40), even though rarely exceed 60; represented by ASCII code.

```
@HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
TAAAGAAACTAAGAATAAGCAGATTATCTCGTAT
+HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
fffffdaadKccaccfffefdceffefefe`b`
```

# QUALITY SCORE
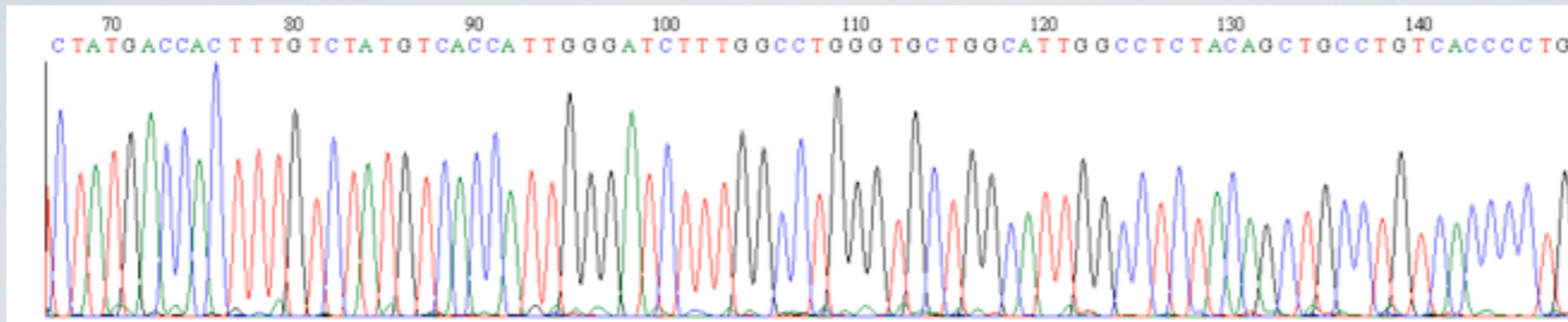
Quality scores are currently calculated to reliably call bases from a Sanger chromatogram; well-known as Phred scores.

They range from 0 to 93 (Illumina 0 - 40), even though rarely exceed 60; represented by ASCII code.

```
@HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
TAAAGAAACTAAGAATAAGCAGATTATCTCGTAT
+HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
fffffdaadKccaccfffefdceffefefe`b`
```

Sanger quality code (Phred): ASCII character code = phred quality value + 33

Illumina quality code: ASCII character code = phred quality value + 64
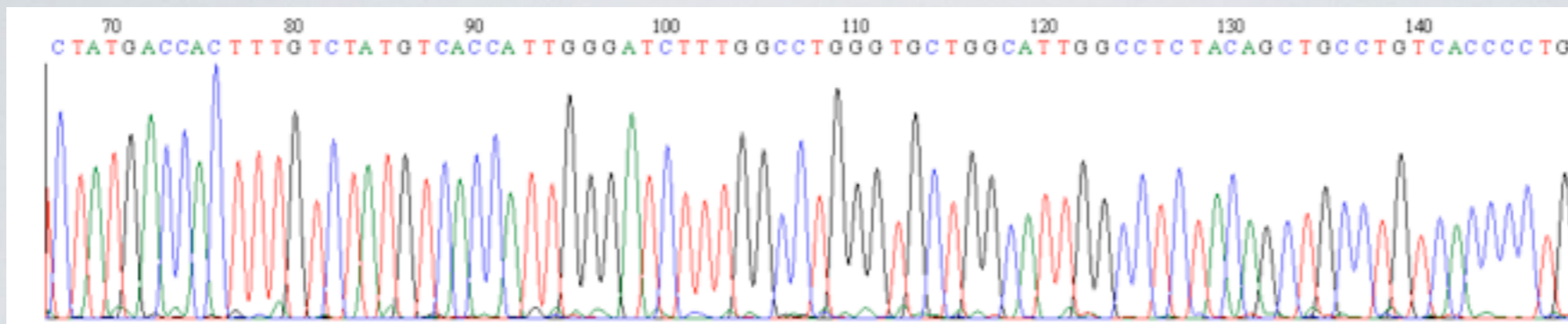
# QUALITY SCORE

Quality scores are currently calculated to reliably call bases from a Sanger chromatogram; well-known as Phred scores.
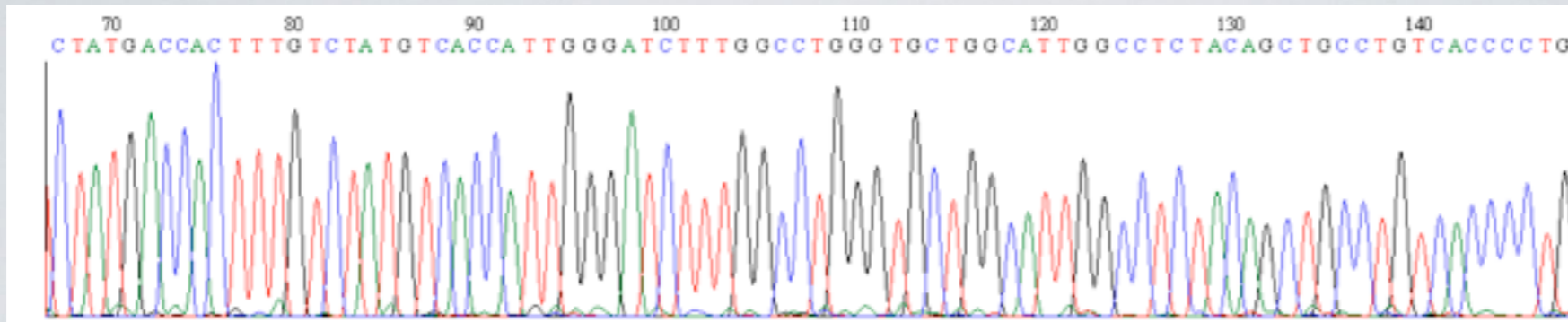
They range from 0 to 93 (Illumina 0 - 40), even though rarely exceed 60; represented by ASCII code.

```
@HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
TAAAGAAACTAAGAATAAGCAGATTATCTCGTAT
+HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
fffffdaadKccaccfffefdceffefefe`b`
```

Sanger quality code (Phred): ASCII character code = phred

Illumina quality code: ASCII character code = phred quality

| Quality Value | Error Probability | Probability Called Base is Correct |
|---|---|---|
| 10 | 0.1 | 0.9 |
| 20 | 0.01 | 0.99 |
| 30 | 0.001 | 0.999 |
| 40 | 0.0001 | 0.9999 |

$q = -10\log_{10}(p)$

# QUALITY SCORE

Quality scores are currently calculated to reliably call bases from a Sanger chromatogram; well-known as Phred scores.

They range from 0 to 93 (Illumina 0 - 40), even though rarely exceed 60; represented by ASCII code.

```
@HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
TAAAGAAACTAAGAATAAGCAGATTATCTCGTAT
+HWIEAS210R_0008:6:1:1600:1545#NNCANC/1
fffffdaadKccaccfffefdceffefefe`b`
```
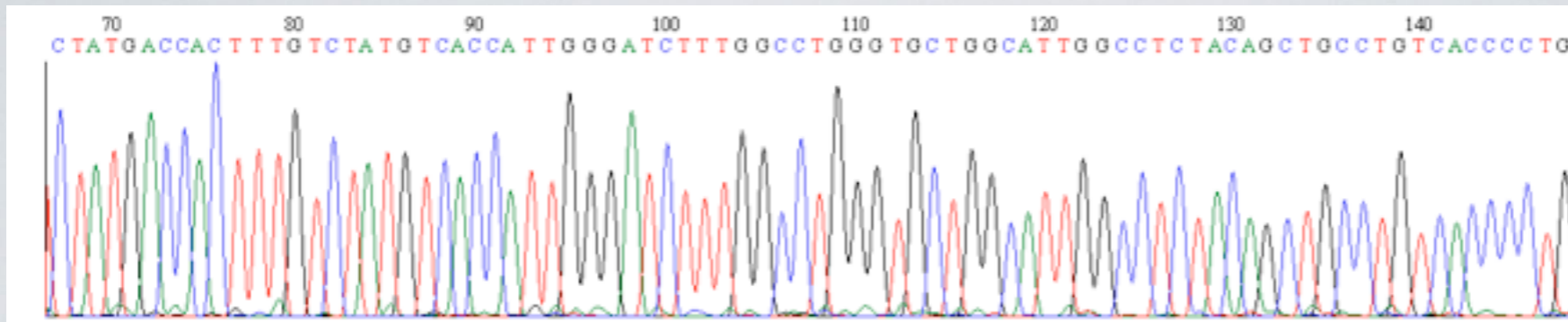
Sanger quality code (Phred): ASCII character code = phred

Illumina quality code: ASCII character code = phred quality

Illumina: f (ASCII 102) => 102 - 64 = 38
Phred:  f  (ASCII 102) => 102 - 33 = 69
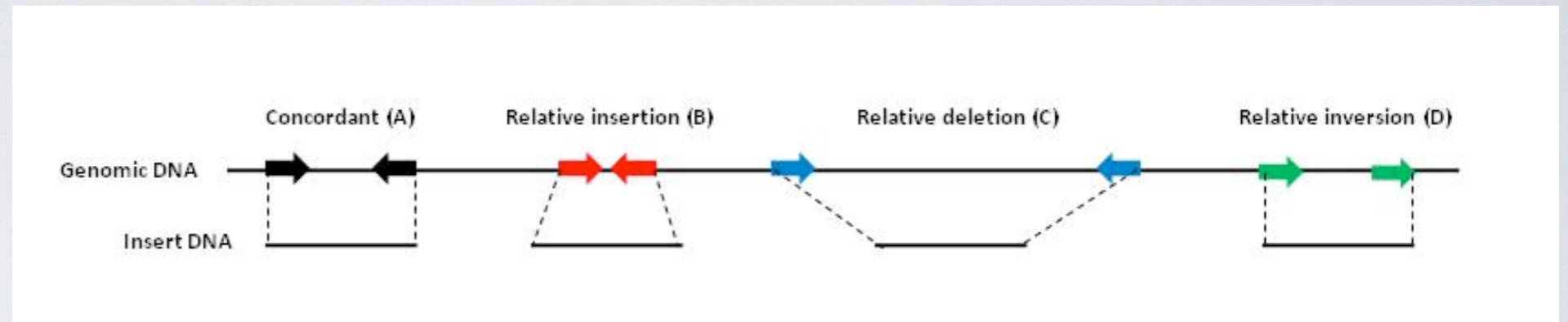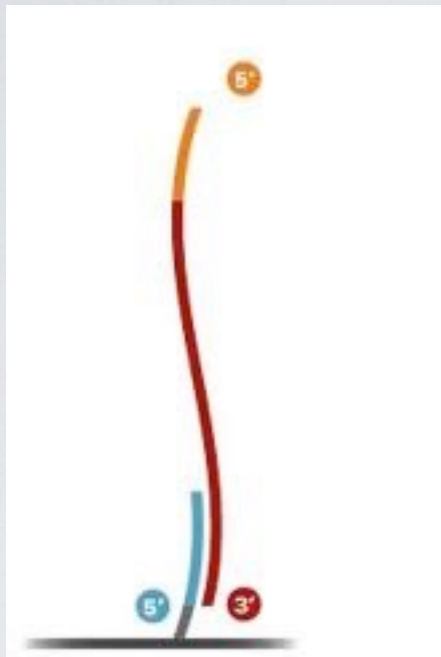Illumina: ` (ASCII 96) => 96 - 64 = 32

| Quality Value | Error Probability | Probability Called Base is Correct |
|---|---|---|
| 10 | 0.1 | 0.9 |
| 20 | 0.01 | 0.99 |
| 30 | 0.001 | 0.999 |
| 40 | 0.0001 | 0.9999 |

$q = -10\log_{10}(p)$

# QUALITY SCORE

# Sequencing of the two end of the same DNA fragment



# PAIRED-END

# Sequencing of the two end of the same DNA fragment



single reads

paird-end reads
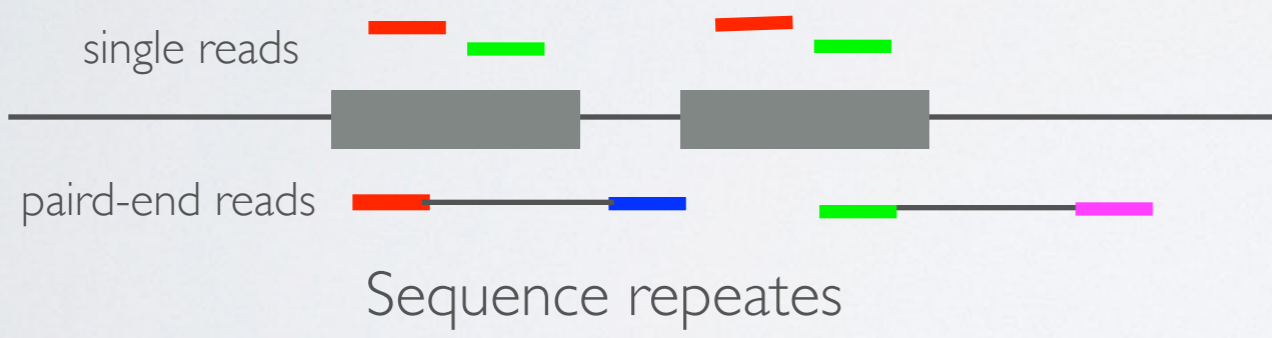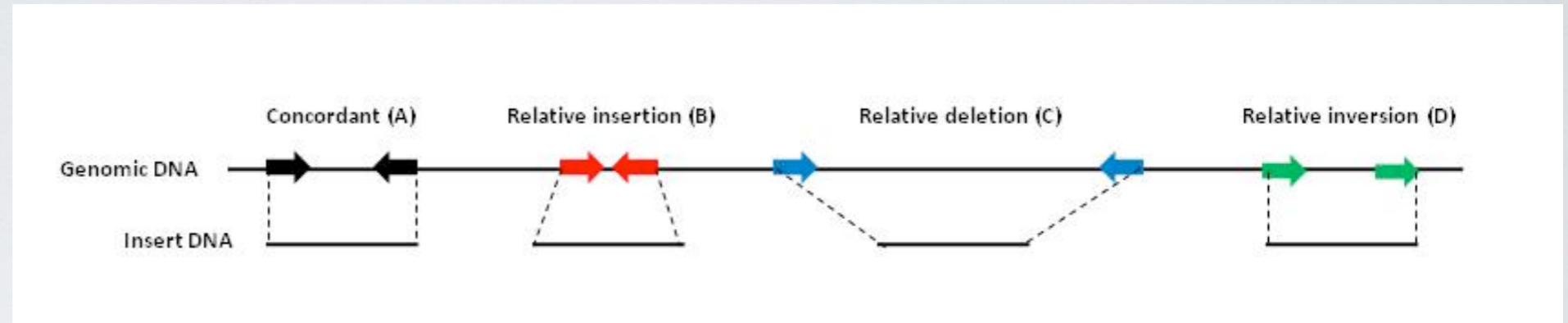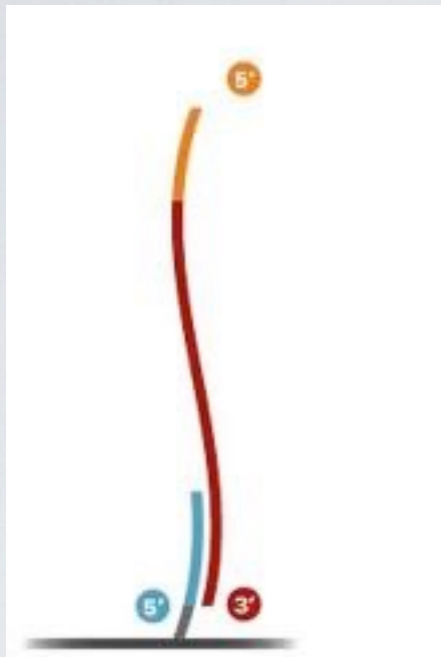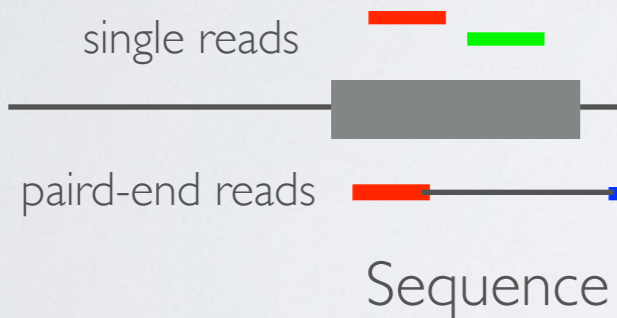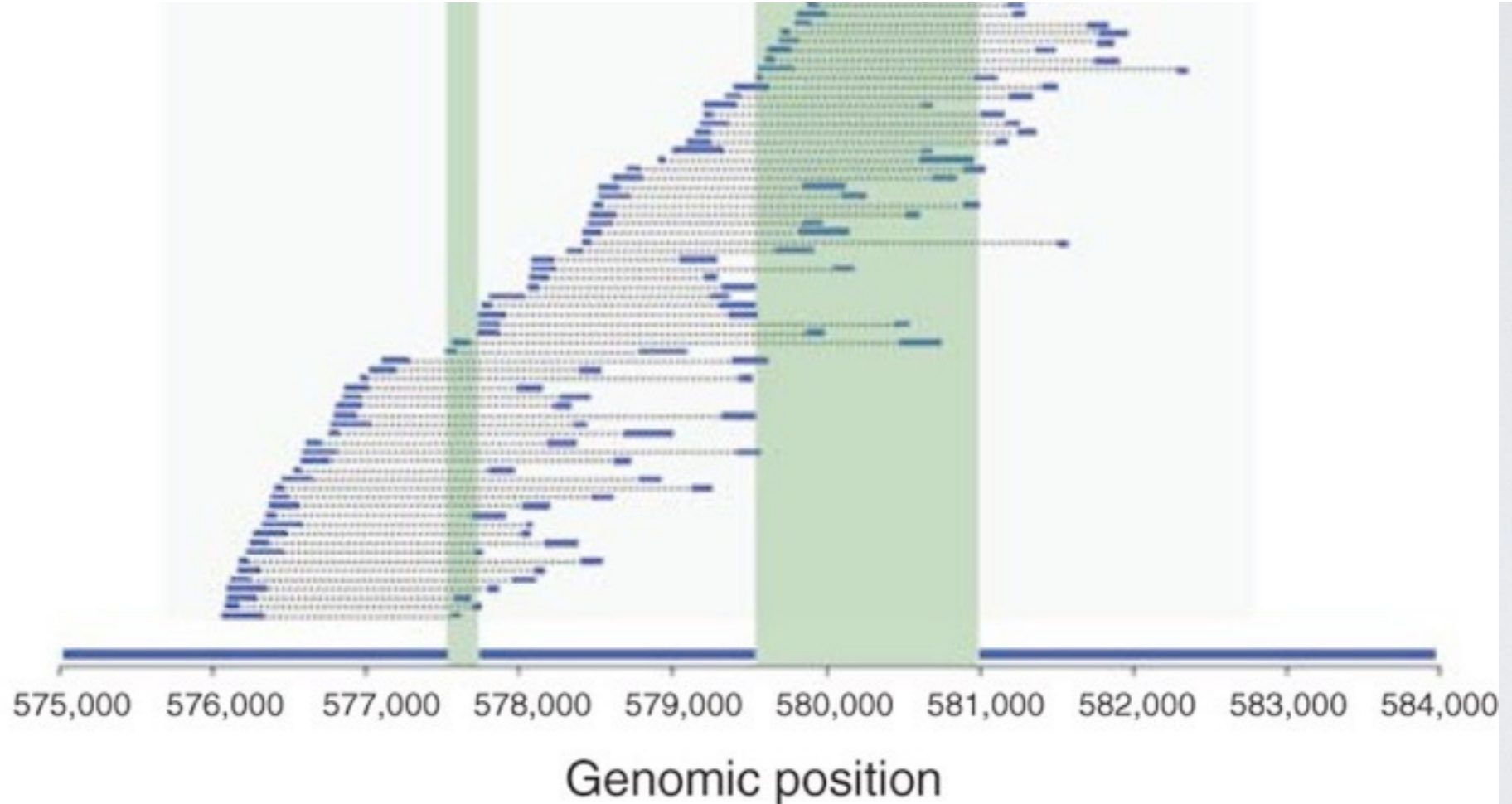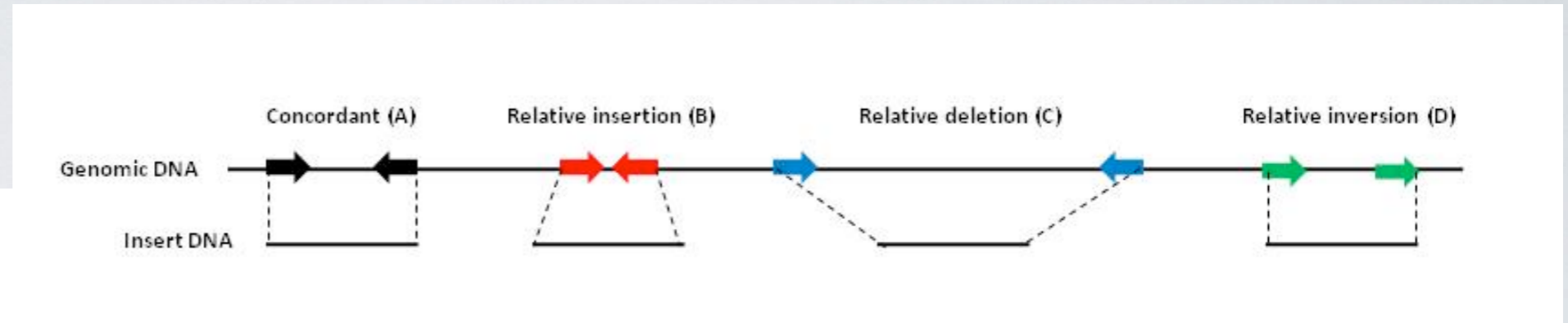
Sequence repeates

# PAIRED-END

# Sequencing of the two end of the same DNA fragment



single reads

paird-end reads

Sequence

# PAIRED-END

# APPLICATIONS

**Sample Prep**

**Whole genome**
- Resequencing → sequence variations such as SNP, CNV, inserts, deletions, reversions
- De-novo → new genomes
- Targeted → sequence variations with higher coverage
- Metagenomics → environmental studies, community studies

**Transcriptome**
- RNA-Seq
- DGE
- Small RNA
- miRNA

**Regulation**
- Methylation
- ChIP-Seq

# APPLICATIONS

**Sample Prep**

**Whole genome**
- Resequencing  →  sequence variations such as SNP, CNV, inserts, deletions, reversions
- De-novo  →  new genomes
- Targeted  →  sequence variations with higher coverage
- Metagenomics  →  environmental studies, community studies

**Transcriptome**
- RNA-Seq  →  transcriptomics, splicing variants,
- DGE  →  digital gene expression
- Small RNA  →  non-coding RNA research
- miRNA  →  miRNA induced regulations

**Regulation**
- Methylation
- ChIP-Seq

# APPLICATIONS

**Sample Prep**

**Whole genome**
- Reasequencing → sequence variations such as SNP, CNV, inserts, deletions, reversions
- De-novo → new genomes
- Targeted → sequence variations with higher coverage
- Metagenomics → environmental studies, community studies

**Transcriptome**
- RNA-Seq → transcriptomics, splicing variants,
- DGE → digital gene expression
- Small RNA → non-coding RNA research
- miRNA → miRNA induced regulations

**Regulation**
- Methylation → epigenetics, methylation induced regulations
- ChIP-Seq → protein DNA interactions such as TFBS, histon, polymerase

# APPLICATIONS

☐ Quality control

☐ Mapping
☐ Assembly
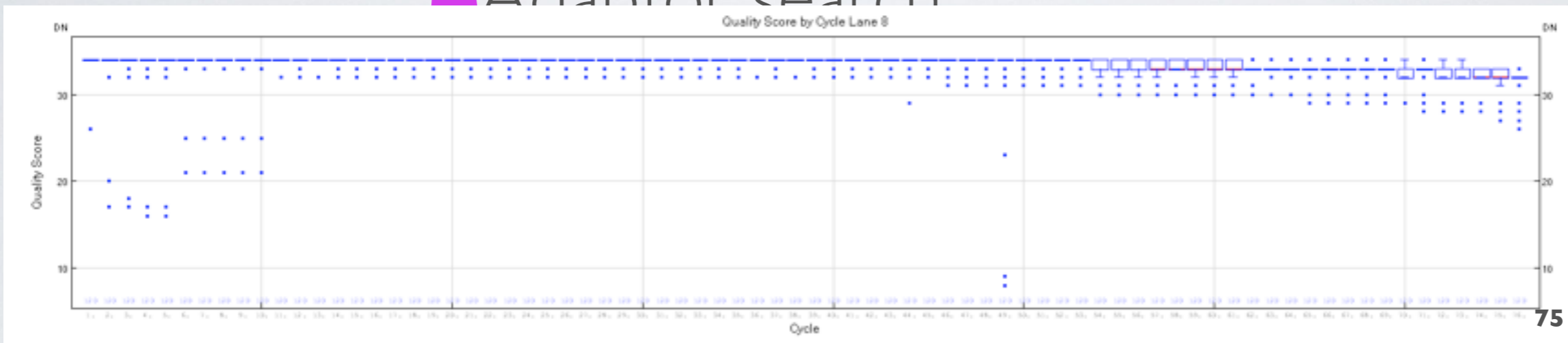☐ Digital gene expression

# DATA ANALYSIS

- Quality control
  - Quality score distribution
  - Sequence size distribution
  - Sequence coverage
  - Adaptor search

- Mapping
- Assembly
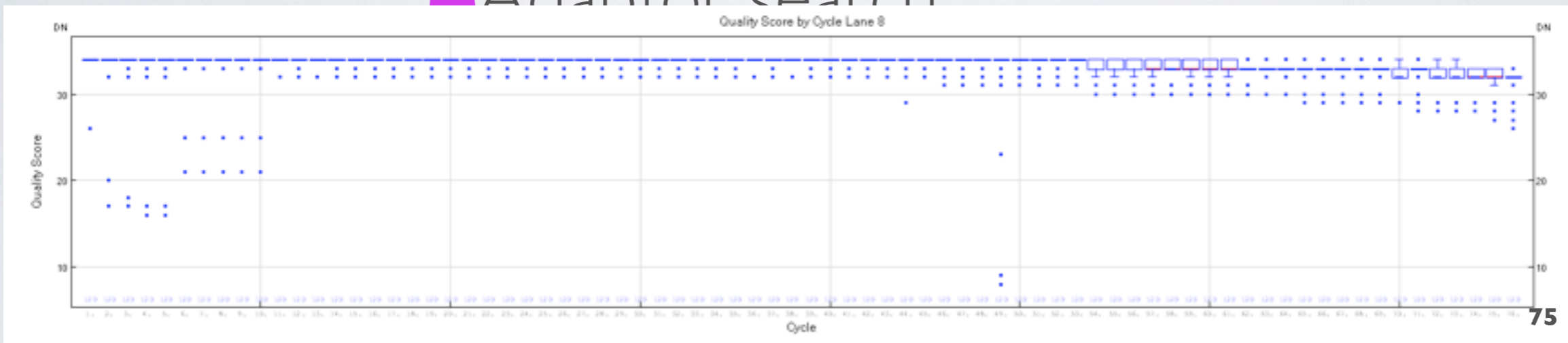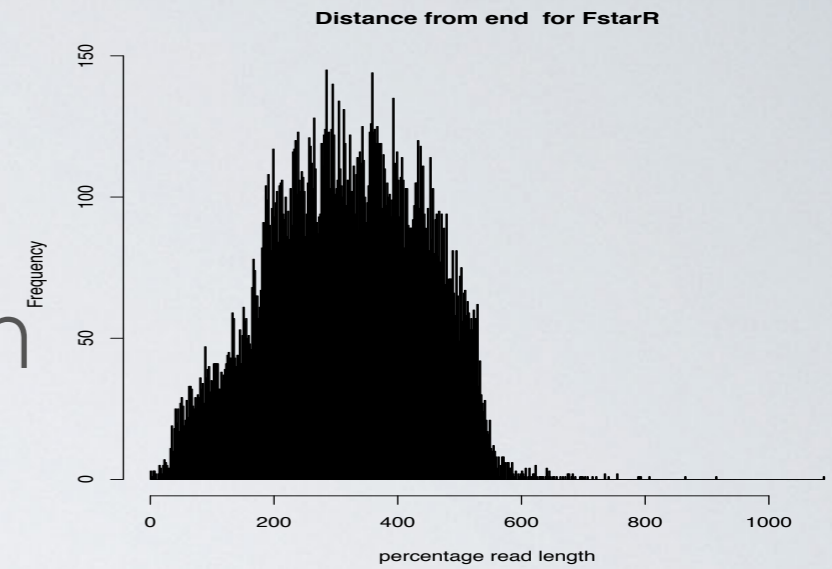- Digital gene expression

# DATA ANALYSIS

Quality control
- Quality score distribution
- Sequence size distribution
- Sequence coverage
- Adaptor search



75

Digital gene expression

# DATA ANALYSIS

Quality control
Quality score distribution
Sequence size distribution
Sequence coverage
Adaptor search

Digital gene expression

# DATA ANALYSIS

Quality control
- Quality score distribution
- Sequence size distribution
- Sequence coverage
- Adaptor search

DATA ANALYSIS

☐ Quality control

☐ Mapping
▪ Aligning short sequences on a reference
☐ Assembly
☐ Digital gene expression

# DATA ANALYSIS

Quality control

Mapping
Aligning short sequenc
Assembly
Digital gene expression

# DATA ANALYSIS

☐ Quality control

☐ Mapping
☐ Assembly
☐ Digital gene expression
☐ Visualisation

# DATA ANALYSIS

☐ Quality control

☐ Mapping
☐ Assembly
  ◼ Assemble sequencing reads to recover the sequence in investigation
☐ Digital gene expression
☐ Visualisation

# DATA ANALYSIS

☐ Quality control

☐ Mapping
☐ Assembly
◼ Assemble sequencing reads to recover the
sequence in investigation
☐ Digital gene expression
☐ Visualisation

# DATA ANALYSIS

Quality control

Mapping

Assembly

Assemble sequencing reads to recover the sequence in investigation

Digital gene expression

Visualisation



# DATA ANALYSIS

☐ Quality control

☐ Mapping
☐ Alignment
☐ Digital gene expression
◼ Compare differentially expressed sequences using there frequency witin the data set
☐ Visualisation

# DATA ANALYSIS

Qua...

Map...
Alig...
Dig...
Co... using
there...
Visu...

# DATA ANALYSIS

Qua

Map
Alig
Dig
Co                                                                using
there
Visu

# DATA ANALYSIS

36

DATA ANALYSIS

Qua

Map
Alig
Dig
Co                                                                    using
there
Visu

# DATA ANALYSIS

36

Qua

Map
Alig
Dig
Co                                                      using
there
Visu

tree 1

tree 2

DATA ANALY

Log₂ rpm CTV-infected
Log₂ rpm mock-inoculated

36

☐ Quality control

☐ Mapping
☐ Assembly
☐ Digital gene expression
☐ Visualisation
◼ Specialized browsers to visualize the vast amount of mapped sequences

# DATA ANALYSIS

# DATA ANALYSIS

- Quality control
- Mapping
- Assembly
- Digital gene expression
- Visualisation
  - Specialized browsers to visualize the vast amount of mapped sequences

# TOOLS

☐ Quality control

☐ Mapping
☐ Assembly
☐ Digital gene expression
☐ Visualisation

# TOOLS

# TOOLS

☐ Quality control
◢ in most cases incorporated in sequencing platform software
◢ GALAXY


☐ Mapping
☐ Assembly
☐ Digital gene expression
☐ Visualisation

# TOOLS

☐ Quality control

☐ Mapping
  ◼ read indexing with hash table
  ◼ genome indexing with hash table
  ◼ genome indexing with suffix array
  ◼ SAM/BAM format
  ◼ http://lh3lh3.users.sourceforge.net/NGSalign.shtml
☐ Assembly
☐ Digital gene expression
☐ Visualisation

# TOOLS

# SAMTools

- SAM Tools provide various utilities for manipulating alignments in the SAM format, including sorting, merging, indexing and generating alignments in a per–position format.
- SAM (Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments. SAM aims to be a format that:

```
HWIEAS210R_0008:6:1:1118:15625#NNTANG/1_18 0    scaffold_2   19205786 255 18M *    0    0
AGACCGGTAGACTTGAAC    d\ddd^a``^G_\bT_dd    XA:i:0   MD:Z:18   NM:i:0
```

| Col | Field | Description |
|-----|-------|-------------|
| 1 | QNAME | Query (pair) NAME |
| 2 | FLAG | bitwise FLAG |
| 3 | RNAME | Reference sequence NAME |
| 4 | POS | 1-based leftmost POSition/coordinate of clipped sequence |
| 5 | MAPQ | MAPping Quality (Phred-scaled) |
| 6 | CIAGR | extended CIGAR string |
| 7 | MRNM | Mate Reference sequence NaMe ('=' if same as RNAME) |
| 8 | MPOS | 1-based Mate POSistion |
| 9 | ISIZE | Inferred insert SIZE |
| 10 | SEQ | query SEQuence on the same strand as the reference |
| 11 | QUAL | query QUALity (ASCII-33 gives the Phred base quality) |
| 12 | OPT | variable OPTional fields in the format TAG:VTYPE:VALUE |

samtools view -bt ref_list.txt -o aln.bam aln.sam.gz

samtools sort aln.bam aln.sorted

samtools index aln.sorted.bam

samtools idxstats aln.sorted.bam

samtools view aln.sorted.bam chr2:20,100,000-20,200,000

samtools merge out.bam in1.bam in2.bam in3.bam

samtools faidx ref.fasta

samtools pileup -vcf ref.fasta aln.sorted.bam

samtools mpileup -C50 -gf ref.fasta -r chr3:1,000-2,000 in1.bam in2.bam

samtools tview aln.sorted.bam ref.fasta

- BAM Binary version of SAM

# TOOLS

☐ Quality control

☐ Mapping
☐ Assembly
◼ Greedy
◼ Overlap Layout Consensus (OLC)
◼ de Bruijn graph based

☐ Digital gene expression
☐ Visualisation

# TOOLS

# Assembler

## Greedy

The greedy algorithms apply one basic operation: given any read or contig, add one more contig. The basic operation is repeated until no more operations are possible. Each operation uses the next highest-scoring overlap to make the next join.

## Overlap Layout Consensus (OLC)

step 1 overlap discovery
step 2 build and use the overlap graph
step 3 multiple sequence alignment

## de Bruijn graph bases

The de Bruijn graph approach circumvents the problems of overlap consensus assembly. Rather than using the reads 'as is' and trying to link them, the k-mers (all subsequences of length k within the reads) are computed and the reads are represented as a path through the k-mers. Such a paradigm handles redundancy better than the overlap consensus approach and makes the computation of paths more tractable.

# TOOLS

# Assembler

## Greedy

The greedy algorithms apply one basic operation: given any read or contig, add one more contig. The basic operation is repeated until no more operations are possible. Each operation uses the next highest-scoring overlap to make the next join.

## Overlap Layout C

step 1 overlap discovery
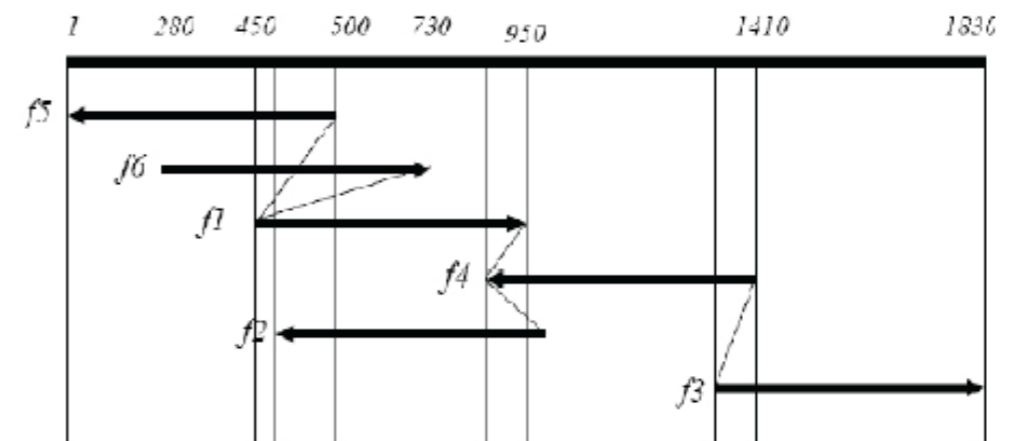step 2 build and use the over
step 3 multiple sequence alig

## de Bruijn graph b

The de Bruijn graph approach
Rather than using the reads 'a
length k within the reads) are
k-mers. Such a paradigm hand
and makes the computation o



Overlap:

$f_i$

$f_j$

| 1 | 280 | 450 | 500 | 730 | 950 | 1410 | 1836 |

f5

f6

f1

Layout:

f4

f2

f3

Consensus:

| R1 | ACGCTCCAACCGCTAATACG |
| R2 | ATCGCTAATCCACGCCCGCCCCGC |
| R2 | AAAC–CTCCAACCG |
| R3 | TGCGCGCCCGCCCCGAAACCGC |
| Consensus | AAAC–CTCCAACCGCTAATGCGCGCCCGCCCCGAAACCGC |

# TOOLS

# Assembler

## Greedy

The greedy algorithms apply one basic operation: given any read or contig, add one more contig. The basic operation is repeated until no more operations are possible. Each operation uses the next highest-scoring overlap to make the next join.
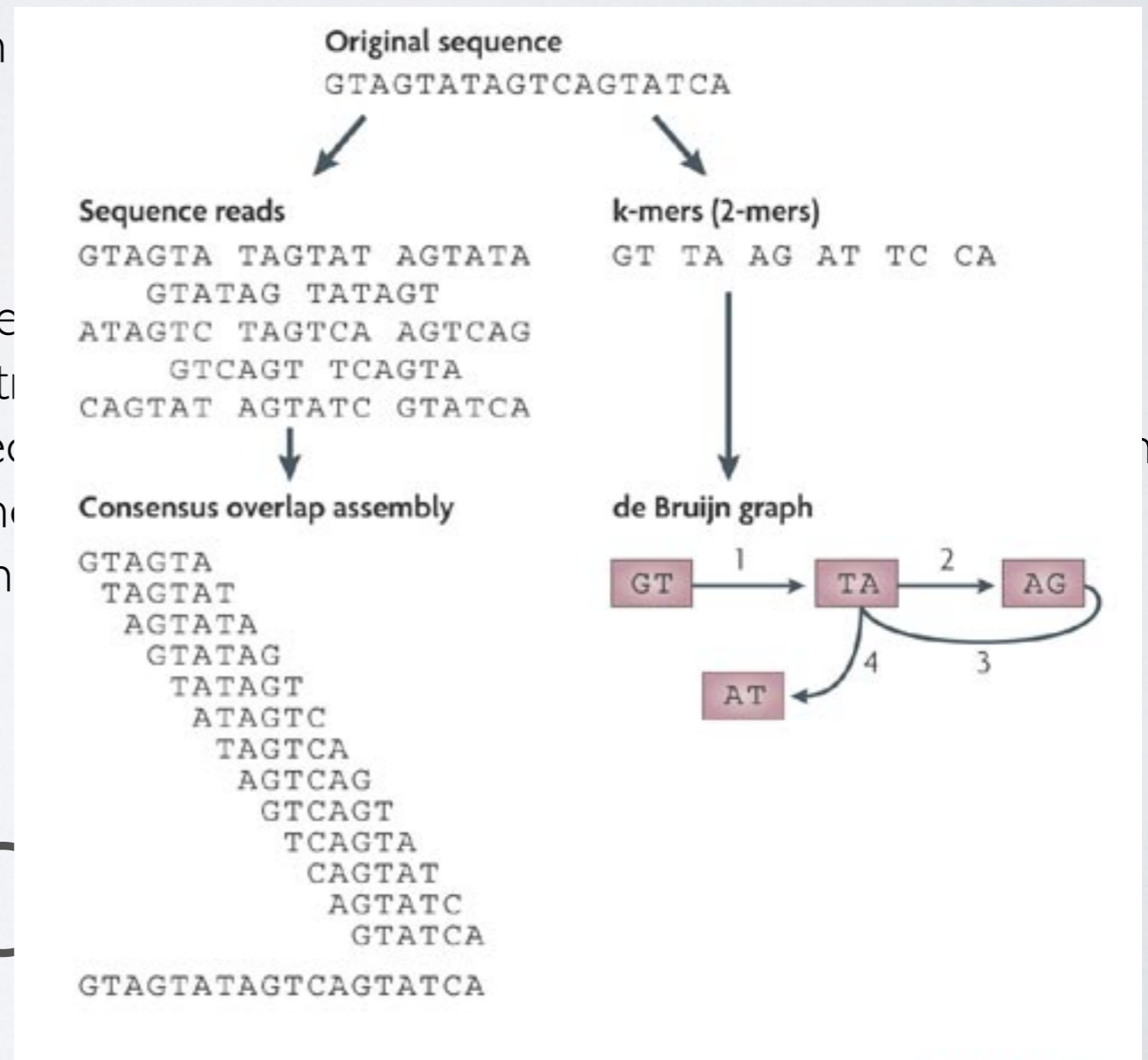
## Overlap Layout Consensus (OLC)

step 1 overlap discovery
step 2 build and use the overlap graph
step 3 multiple sequence alignment

## de Bruijn graph bases

The de Bruijn graph approach circumvents the problems of overlap consensus assembly. Rather than using the reads 'as is' and trying to link them, the k-mers (all subsequences of length k within the reads) are computed and the reads are represented as a path through the k-mers. Such a paradigm handles redundancy better than the overlap consensus approach and makes the computation of paths more tractable.

# TOOLS

# Assembler

## Greedy

The greedy algorithms apply one basic operation: given any read or contig, add one more contig. The basic operation is repeated until no more operations are possible. Each operation uses the next highest-scoring overlap to make the next join.

## Overlap Layout Consensus (OLC)

step 1 overlap discovery
step 2 build and use the overlap graph
step 3 multiple sequence alignment

## de Bruijn graph bases

The de Bruijn graph approach circumve
Rather than using the reads 'as is' and t
length k within the reads) are computed                                    he
k-mers. Such a paradigm handles redun
and makes the computation of paths m



Original sequence
GTAGTATAGTCAGTATCA

Sequence reads
GTAGTA  TAGTAT  AGTATA
    GTATAG  TATAGT
ATAGTC  TAGTCA  AGTCAG
    GTCAGT  TCAGTA
CAGTAT  AGTATC  GTATCA

k-mers (2-mers)
GT  TA  AG  AT  TC  CA

Consensus overlap assembly
GTAGTA
  TAGTAT
    AGTATA
      GTATAG
        TATAGT
          ATAGTC
            TAGTCA
              AGTCAG
                GTCAGT
                  TCAGTA
                    CAGTAT
                      AGTATC
                        GTATCA
GTAGTATAGTCAGTATCA

de Bruijn graph

# Assemblers

| Name | Algorithm | Author | Year |
|------|-----------|--------|------|
| Arachne WGA | OLC | Batzoglou, S. et al. | 2002 / 2003 |
| Celera WGA Assembler / CABOG | OLC | Myers, G. et al.; Miller G. et al. | 2004 / 2008 |
| Minimus (AMOS) | OLC | Sommer, D.D. et al. | 2007 |
| Newbler | OLC | 454/Roche | 2009 |
| Edena | OLC | Hernandez D., et al. | 2008 |
| SUTTA | B&B | NYU/Abraxis (unpublished) | 2009/2010 |
| TIGR | Greedy | TIGR | 1995 / 2003 |
| Phusion | Greedy | Mullikin JC, et.al. | 2003 |
| Phrap | Greedy | Green, P. | 2002 / 2003 / 2008 |
| CAP3, PCAP | Greedy | Huang, X. et al. | 1999 / 2005 |
| Euler | SBH | Pevzner, P. et al. | 2001 / 2006 |
| Euler-SR | SBH | Chaisson, MJ. et al. | 2008 |
| Velvet | SBH | Zerbino, D. et al. | 2007 / 2009 |
| ALLPATHS | SBH | Butler, J. et al. | 2008 |
| ABySS | SBH | Simpson, J. et al. | 2008 / 2009 |
| SOAPdenovo | SBH | Ruiqiang Li, et al. | 2009 |
| SHARCGS | Prefix-Tree | Dohm et al. | 2007 |
| SSAKE | Prefix-Tree | Warren, R. et al. | 2007 |
| VCAKE | Prefix-Tree | Jeck, W. et al. | 2007 |
| QSRA | Prefix-Tree | Douglas W. et al. | 2009 |
| Sequencher | - | Gene Codes Corporation | 2007 |
| SeqMan NGen | - | DNASTAR | 2008 |
| Staden gap4 package | - | Staden et al. | 1991 / 2008 |
| MIRA, miraEST | - | Chevreux, B. | 1998 / 2008 |
| NextGENe | - | Softgenetics | 2008 |
| CLC Genomics Workbench | - | CLC bio | 2008 / 2009 |
| CodonCode Aligner | - | CodonCode Corporation | 2003 / 2009 |

# TOOLS

# Assemblers

Grapevine clone: 6 lanes ($100bp$), insert size $200 \pm 50$
Coverage: $89 \times$

|  | AbySS | SOAPdenovo | CLC |
|---|---|---|---|
| # Scaf num | 289,854 (244k) | 127,648 (368k) | 151,288 (423k) |
| Tot Scaf. length (bp) | 562M (158M) | 257M (285M) | 339M (382M) |
| Max Scaf length (bp) | 89,700 (12k) | 59,054 (36k) | 69,474 (70k) |
| Mean Scaf lgth (bp) | 1942 (649) | 2014 (776) | 2241 (904) |
| N50 length | 2634 (872) | 3186 (2038) | 3328 (1823) |
| time | 18h 49m (12h) | 8h 57m (1d) | 6h 45m (7h) |
| RAM available (GB) | 130 (240) | 240 (120) | 120 (120) |
| RAM used (GB) | $\sim$ 90 (102) | 143 (70) | $\sim$ 80 (60) |
| CPUs | 80 (80) | 8 (8) | 8 (8) |

Grapevine genome size: 475Mb

Policriti et al per. com.

TOOLS

45

☐ Quality control

☐ Mapping
☐ Assembly
☐ Digital gene expression
◼ DESeq, BaySeq, edgeR are R package to analyse count data from high-throughput sequencing assays such as RNA-Seq and test for differential expression.
☐ Visualization

# TOOLS

☐ Quality control

☐ Mapping
☐ Assembly
☐ Digital gene expr...
■ DESeq, BaySeq, e... t
data from high-thro...
RNA-Seq and test
☐ Visualization



TOOLS

Quality control

Mapping
Assembly
Digital gene expression
Visualization
http://lh3lh3.users.sourceforge.net/NGSalnview.shtml

# TOOLS

# Tablet (http://bioinf.scri.ac.uk/tablet/)



# TOOLS

# GBrowse (http://gmod.org/wiki/GBrowse/)

Qual

Mapp

Asse

Digit

Visua

http: alnview.shtml

# Artemis (http://www.sanger.ac.uk/resources/software/artemis/)



# TOOLS

# http://www.ebi.ac.uk/ena/



# DATA REPOSITORY

# http://www.ebi.ac.uk/ena/



# DATA REPOSITORY

# THANKS!!!

# THANKS!!!

# HELICOS

☐ DNA template immobilized to a flow cell
☐ NO amplification (true single molecule sequencing)
☐ Sequencing on flow cell (1000M reads)
☐ Sequencing by synthesis (fluorescence)
☐ Read length up to 50nt average 32
☐ High error rate

# HELICOS

☐ DNA template immobilized to a flow cell
☐ NO amplification (true single molecule sequencing)
☐ Sequencing on flow cell (1000M reads)
☐ Sequencing by synthesis (fluorescence)
☐ Read length up to 50nt average 32
☐ High error rate

# HELICOS

- DNA template immobilized to a flow cell
- NO amplification (true single molecule sequencing)
- Sequencing on flow cell (1000M reads)
- Sequencing by synthesis (fluorescence)
- Read length up to 50nt average 32
- High error rate



**Position**

| | 1 | 2 | 3 |
|---|---|---|---|
| Cycle 1 | · | · | G |
| Cycle 2 | C | C | · |
| Cycle 3 | A | A | A |
| Cycle 4 | · | · | T |
| Cycle 5 | C | · | · |
| · | | | |
| Cycle X | G | G | · |

An image taken by the HeliScope Single Molecule Sequencer. Inset shows a close-up view of individual single molecules.



# HELICOS

☐ Quality control

☐ Mapping
☐ Alignment
  ◼ BLAST - Basic Local Alignment Search Tool
  ◼ GAST - Global Alignment for Sequence Taxonomy

☐ Assembly
☐ Digital gene expression
☐ Visualisation

# TOOLS

Quality control

Mapping
Alignment
■ BLAST - Basic L...
■ GAST - Global A...                                                          ...my

Assembly
Digital gene expr...
Visualisation

**Global Alignment for
Sequence Taxonomy (GAST)**

**Creating RefSSU and RefVx**

1. Full-length SSU rRNA reference sequences are downloaded from SILVA and low-quality sequences removed.

2. Taxonomy is assigned with RDP, additional taxonomy sources (e.g. reference genomes) are added.

3. The V3 and V6 regions are excised from the ARB alignment using primer locations.

4. Additional filters remove low-quality reference tags to create RefV3 or RefV6.
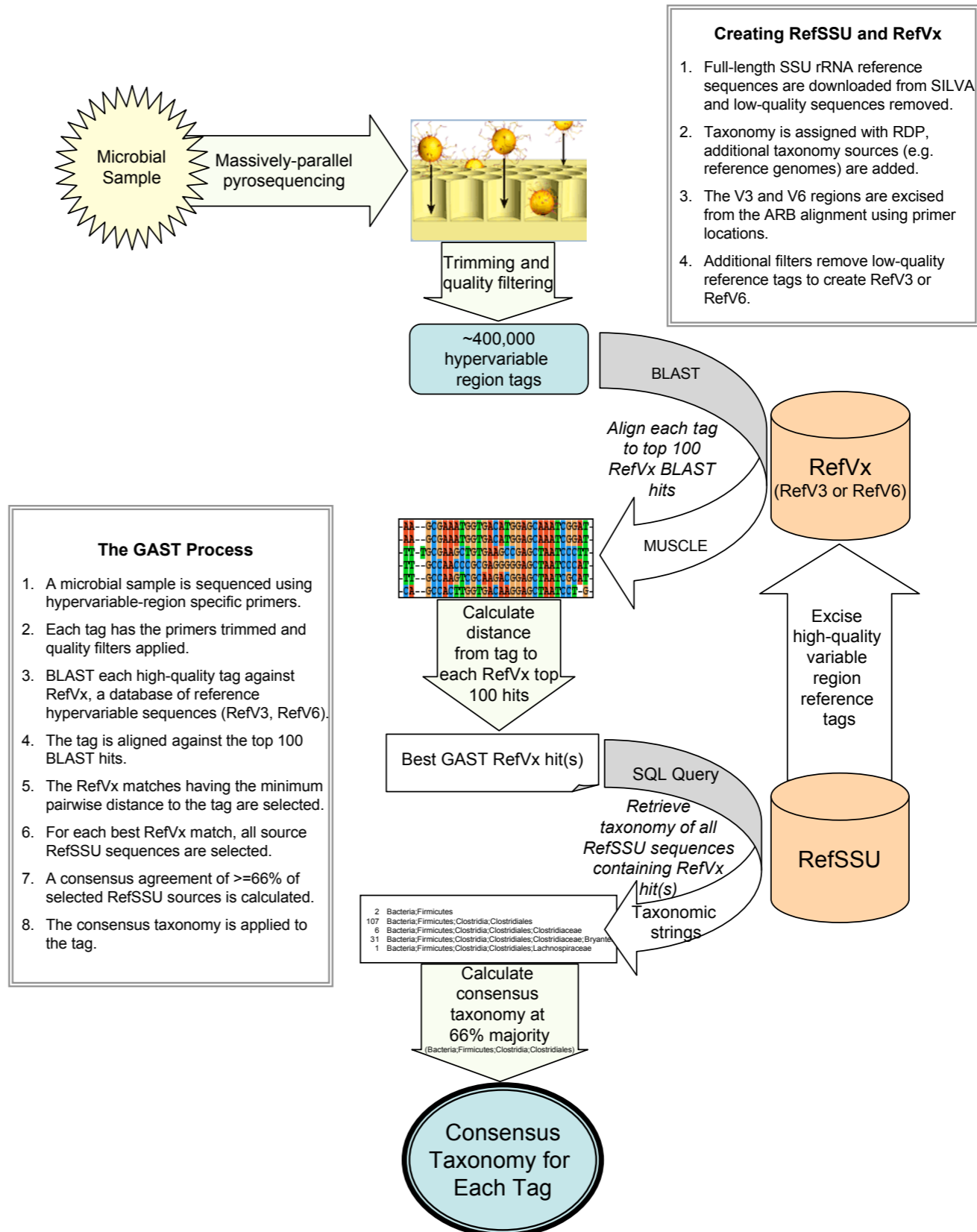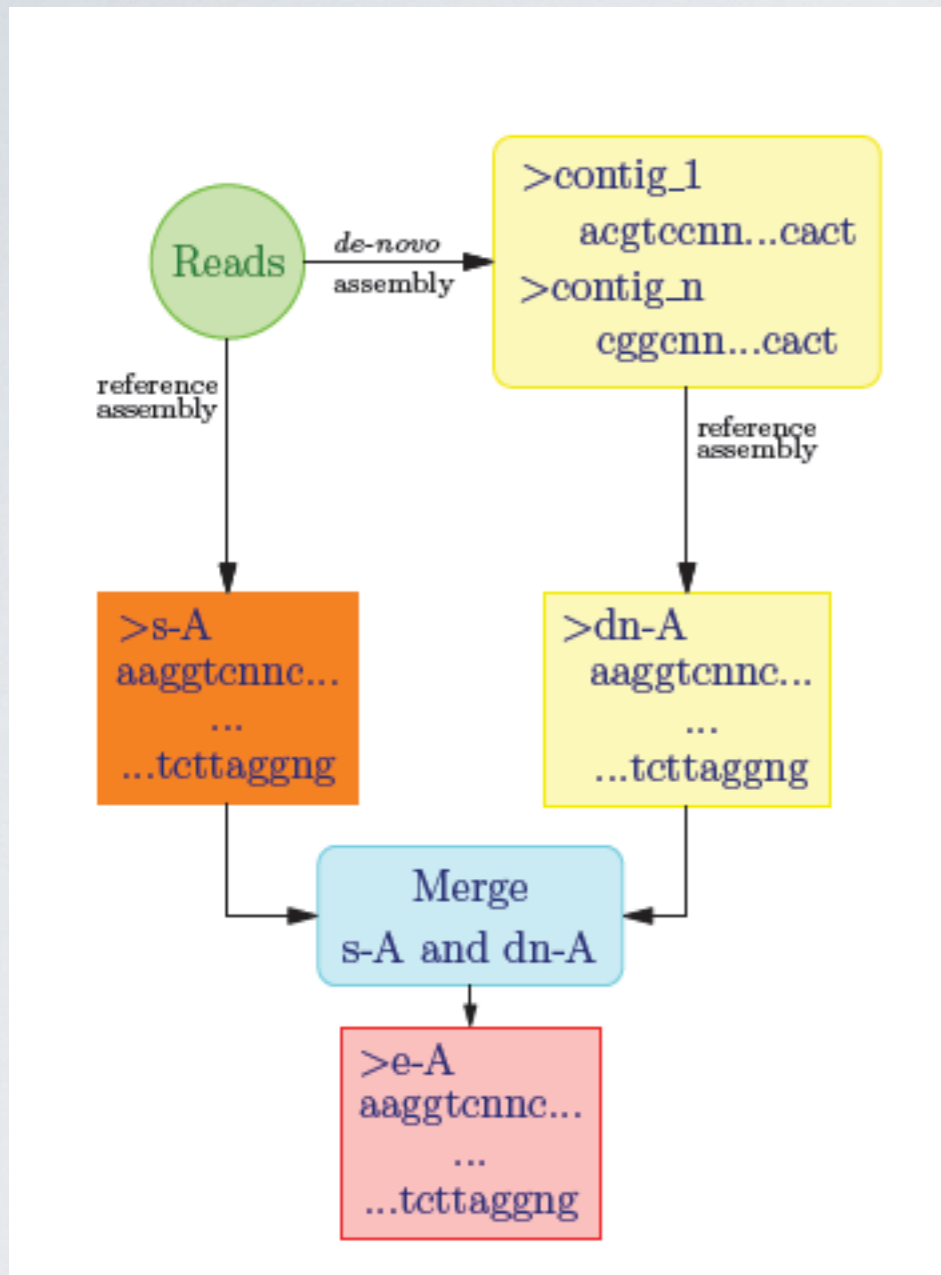
Microbial Sample

Massively-parallel pyrosequencing

Trimming and quality filtering

~400,000 hypervariable region tags

BLAST

*Align each tag to top 100 RefVx BLAST hits*

RefVx
(RefV3 or RefV6)

MUSCLE

**The GAST Process**

1. A microbial sample is sequenced using hypervariable-region specific primers.

2. Each tag has the primers trimmed and quality filters applied.

3. BLAST each high-quality tag against RefVx, a database of reference hypervariable sequences (RefV3, RefV6).

4. The tag is aligned against the top 100 BLAST hits.

5. The RefVx matches having the minimum pairwise distance to the tag are selected.

6. For each best RefVx match, all source RefSSU sequences are selected.

7. A consensus agreement of >=66% of selected RefSSU sources is calculated.

8. The consensus taxonomy is applied to the tag.

Calculate distance from tag to each RefVx top 100 hits

Excise high-quality variable region reference tags

Best GAST RefVx hit(s)

SQL Query

*Retrieve taxonomy of all RefSSU sequences containing RefVx hit(s)*

RefSSU

Taxonomic strings

```
  2  Bacteria;Firmicutes
107  Bacteria;Firmicutes;Clostridia;Clostridiales
  6  Bacteria;Firmicutes;Clostridia;Clostridiales;Clostridiaceae
 31  Bacteria;Firmicutes;Clostridia;Clostridiales;Clostridiaceae;Bryant...
  1  Bacteria;Firmicutes;Clostridia;Clostridiales;Lachnospiraceae
```

Calculate consensus taxonomy at 66% majority
(Bacteria;Firmicutes;Clostridia;Clostridiales)

Consensus Taxonomy for Each Tag

TOOLS

TOOLS

# TOOLS