# NGS sequencing

Dr Francois Sabot & Christine Tranchant-Dubreuil

8th of October, 2018

IRD - UMR DIADE

# Introduction

From The economist



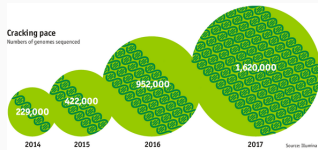From Bussiness Insider

3

- Genetic diversity

- Genetic diversity
- Gene discovery

- Genetic diversity

- Gene discovery

- Genomic structure

- Genetic diversity
- Gene discovery
- Genomic structure
- Contamination/pathogen detection

- Genetic diversity

- Gene discovery

- Genomic structure

- Contamination/pathogen detection

- Metagenomic

- Genetic diversity
- Gene discovery
- Genomic structure
- Contamination/pathogen detection
- Metagenomic
- Pangenomic

- Genetic diversity
- Gene discovery
- Genomic structure
- Contamination/pathogen detection
- Metagenomic
- Pangenomic
- And many other things...

# Methods

$2^{nd}$ Generation Sequencing

- DNA fragmentation (short)
- Matrix amplification
- Short reads
- Limited error rate
- High throughput

$2^{nd}$ Generation Sequencing

- DNA fragmentation (short)
- Matrix amplification
- Short reads
- Limited error rate
- High throughput

$3^{rd}$ Generation Sequencing

- DNA fragmentation (long)
- NO MATRIX AMPLIFICATION
- Long reads
- Important error rate
- Medium throughput
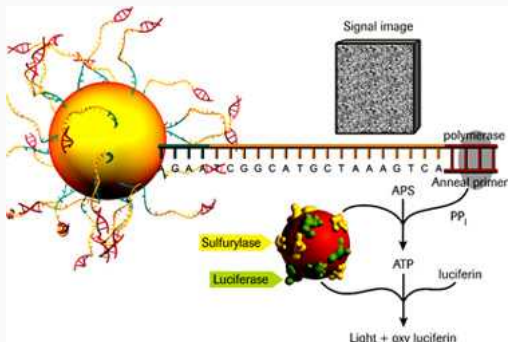
$2^{nd}$ Generation Sequencing

- DNA fragmentation (short)
- Matrix amplification
- Short reads
- Limited error rate
- High throughput

454
IonTorrent
Illumina

$3^{rd}$ Generation Sequencing

- DNA fragmentation (long)
- NO MATRIX AMPLIFICATION
- Long reads
- Important error rate
- Medium throughput

## Current Technologies

$2^{nd}$ Generation Sequencing

- DNA fragmentation (short)
- Matrix amplification
- Short reads
- Limited error rate
- High throughput

454
IonTorrent
Illumina

$3^{rd}$ Generation Sequencing

- DNA fragmentation (long)
- NO MATRIX AMPLIFICATION
- Long reads
- Important error rate
- Medium throughput

PacificBiosciences
Oxford Nanopore

**Advantages** : Length (400 - 750 bases)

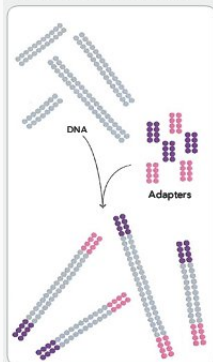**Limits** :

- Homopolymers
- Error rate (15%, non random)
- Output volume
- Price

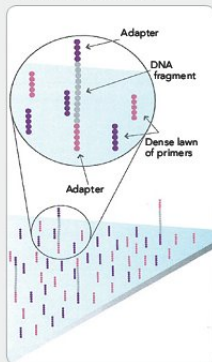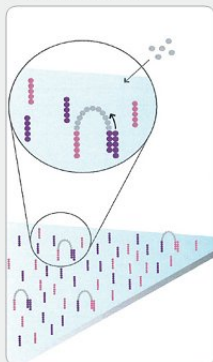Roche has stopped 454 dev, dying technology

**Advantages** :

- Output volume (200+ millions of 150b reads, HiSeq 2500)
- Accuracy (99.99 %)
- Run is cheap
- MySeq is cheap (around 60 000 USD per machine)

**Limits** : Size (150 + 150 in HiSeq4000 and X, but 400 for MySeq)

**Advantages** :

- Price (less than 200 USD per run)
- Lab sized machine (around 80 000 USD per machine)

**Limits** : Error rate (15%)

**Advantages** :

- Price: around 5,000 USD/run (machine at 125k USD)
- Illumina quality data

**Limits** : "Limited" length of long reads (10kb max)

## PacBio - Advantages and Limits

**Advantages** :

- Length (mean 10kb, more than 40kb regularly)
- Single strand direct sequencing, no amplification bias

**Limits** :

- Error Rate (15%, but can be corrected)
- Machine size and price (more than 900 000 USD)
- Run price (600 USD for 500 Mb)

From Circulation Research

- No Amplification

- NO SYNTHESIS

- Very Long Length

From Circulation Research

- No Amplification
- NO SYNTHESIS
- Very Long Length

- Magnetic fields variation measure
- *Minion*: USB key - sized



From Nature Biotechnology

27

## ONT - Advantages and Limits

**Advantages** :

- Length (mean 10-50kb, more than 2Mb reported)
- Bases Modification detection in real-time
- Single strand direct sequencing
- Machine cheap (2,000 USD for Minion)
- Run cheap (1,000 USD for 30Gb by now minimum)
- Fast: 15mn library, 48-72h run

**Limits** :

- Error Rate (6-9%, but can be corrected)
- Quality of DNA limits the sequencing
- Heu...

- SOLiD †, because of too small sequence size and no new dev.
- Helicos †, because of too many errors and trouble in chemistry
- Polonator †, because of too small sequence size
- DNA nanoball sequencing (Complete Genomics©), nobody uses it but CG group.
- Single Sequence magnetic bead (ongoing development)
- Transmission electron microscopy DNA sequencing (ongoing development)

Comparison of high-throughput sequencing methods[63][64]

| Method | Read length | Accuracy (single read not consensus) | Reads per run | Time per run | Cost per 1 million bases (in US$) | Advantages | Disadvantages |
|--------|-------------|--------------------------------------|---------------|--------------|-----------------------------------|------------|---------------|
| Single-molecule real-time sequencing (Pacific Biosciences) | 30,000 bp (N50); maximum read length >100,000 bases[65][66][67] | 87% raw-read accuracy[68] (> 99.999% with CCS or consensus) | 500,000 per Sequel SMRT cell, 10–20 gigabases[65][69][70] | 30 minutes to 20 hours[65][71] | $0.05–$0.08 | Fast. Detects 4mC, 5mC, 6mA.[72] | Moderate throughput. Equipment can be very expensive. |
| Ion semiconductor (Ion Torrent sequencing) | up to 600 bp[73] | 99.6%[74] | up to 80 million | 2 hours | $1 | Less expensive equipment. Fast. | Homopolymer errors. |
| Pyrosequencing (454) | 700 bp | 99.9% | 1 million | 24 hours | $10 | Long read size. Fast. | Runs are expensive. Homopolymer errors. |
| Sequencing by synthesis (Illumina) | MiniSeq, NextSeq: 75-300 bp; MiSeq: 50-600 bp; HiSeq 2500: 50-500 bp; HiSeq 3/4000: 50-300 bp; HiSeq X: 300 bp | 99.9% (Phred30) | MiniSeq/MiSeq: 1-25 Million; NextSeq: 130-00 Million, HiSeq 2500: 300 million - 2 billion, HiSeq 3/4000 2.5 billion, HiSeq X: 3 billion | 1 to 11 days, depending upon sequencer and specified read length[75] | $0.05 to $0.15 | Potential for high sequence yield, depending upon sequencer model and desired application. | Equipment can be very expensive. Requires high concentrations of DNA. |
| Sequencing by ligation (SOLiD sequencing) | 50+35 or 50+50 bp | 99.9% | 1.2 to 1.4 billion | 1 to 2 weeks | $0.13 | Low cost per base. | Slower than other methods. Has issues sequencing palindromic sequences.[76] |
| Nanopore Sequencing | Dependent on library prep, not the device, so user chooses read length. (up to 500 kb reported) | ~92–97% single read (up to 99.96% consensus) | dependent on read length selected by user | data streamed in real time. Choose 1 min to 48 hrs | $500–999 per Flow Cell, base cost dependent on expt | Longest individual reads. Accessible user community. Portable (Palm sized). | Lower throughput than other machines, Single read accuracy in 90s. |
| Chain termination (Sanger sequencing) | 400 to 900 bp | 99.9% | N/A | 20 minutes to 3 hours | $2400 | Useful for many applications. | More expensive and impractical for larger sequencing projects. This method also requires the time consuming step of plasmid cloning or PCR. |

- DNA from plant, animal, microbial...

- DNA from plant, animal, microbial...
- RNA from various sources

- DNA from plant, animal, microbial...
- RNA from various sources
- smallRNA, idem

- DNA from plant, animal, microbial...

- RNA from various sources

- smallRNA, idem

- Environmental sample

- Organite DNA (mitochondria, chloroplast)

- Organite DNA (mitochondria, chloroplast)
- Subsample RNA (exon capture, 16S capture for Barcoding)

- Organite DNA (mitochondria, chloroplast)
- Subsample RNA (exon capture, 16S capture for Barcoding)
- Viral sample from infected tissue

## Alternative Samples

- Organite DNA (mitochondria, chloroplast)
- Subsample RNA (exon capture, 16S capture for Barcoding)
- Viral sample from infected tissue
- As many as you can extract...

# Thanks for your attention