



LEIDEN UNIVERSITY MEDICAL CENTER

ICT and research at the LUMC

Jeroen F. J. Laros

Leiden Genome Technology Center

Department of Human Genetics

Center for Human and Clinical Genetics



ICT and research

Division 5.

ICT and research

Division 5.

- Center for Human and Clinical Genetics.

ICT and research

Division 5.

- Center for Human and Clinical Genetics.
 - Human Genetics.

ICT and research

Division 5.

- Center for Human and Clinical Genetics.
 - Human Genetics.
 - Leiden Genome Technology Center.

ICT and research

Division 5.

- Center for Human and Clinical Genetics.
 - Human Genetics.
 - Leiden Genome Technology Center.

High throughput sequencing.

DNA sequencing



Figure 1: HiSeq 2000.

DNA sequencing



Figure 1: HiSeq 2000.



Figure 2: Flowcell.

DNA sequencing



Figure 1: HiSeq 2000.



Figure 2: Flowcell.

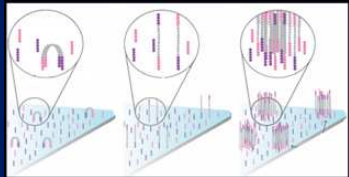


Figure 3: Amplification.

Raw data analysis

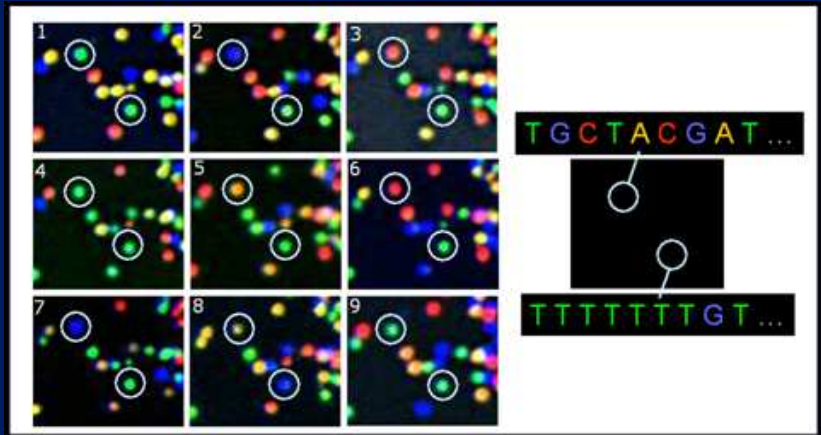


Figure 4: Base calling.

Result of base calling

```

1  @FCC0ADRACXX:2:1101:1684:1875#GCGGAACT/1
2  AGGGGATGCAACCTCGAGGAGGAAAGGAACGAAAGAGGAAGGGAG...
3  +
4  BS\ceeeegggggiihfhighihiiiiihhhihfhighiiii...
5  @FCC0ADRACXX:2:1101:1714:1885#GCGGAACT/1
6  ATTTTCTCAATGTCTACCACTGCGGGTAACACTTTGTGTTCCCA...
7  +
8  bbbbeeeegggggiihiiiiiiiiighiiiiiiiiiggfghiihi...

```

Listing 1: FastQ file.

These files are around 1T big.

Some figures

Runtime 2 weeks (5-10 full genomes).

Produces 6T of pictures per flowcell, 500G base calling.

Needs continuous connection to the storage, otherwise it stalls.

Some figures

Runtime 2 weeks (5-10 full genomes).

Produces 6T of pictures per flowcell, 500G base calling.

Needs continuous connection to the storage, otherwise it stalls.

Basecalling:

$$\begin{array}{r}
 24 \quad \text{time in hours for base calling per flowcell} \\
 12 \times \quad \text{number of cores} \\
 \hline
 2 \times \quad \text{number of flowcells} \\
 \hline
 576 \quad \text{CPU hours per flowcell}
 \end{array}$$

On a desktop machine this would take 48 days.

PacBio sequencing

Pacific Biosciences Single Molecule, Real-Time



Figure 5: PacBio RS.

Base calling

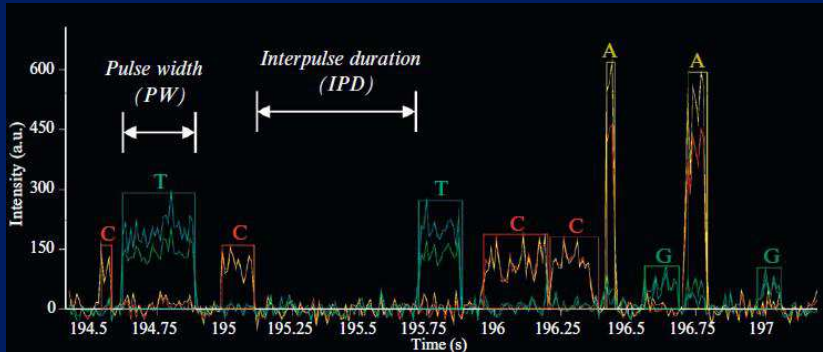


Figure 6: PacBio trace.

Some figures

Runtime 90 minutes.

Some figures

Runtime 90 minutes.

Per cell:

- 3G/sec.
 - 4 cameras.
 - 72 fps.
 - 10M pictures.
- 16T of raw data.
- 4G after base calling.

Some figures

Runtime 90 minutes.

Per cell:

- 3G/sec.
 - 4 cameras.
 - 72 fps.
 - 10M pictures.
- 16T of raw data.
- 4G after base calling.

There are 16 cells in this machine:

- 256T per run.
- 64G after base calling.

Alignment

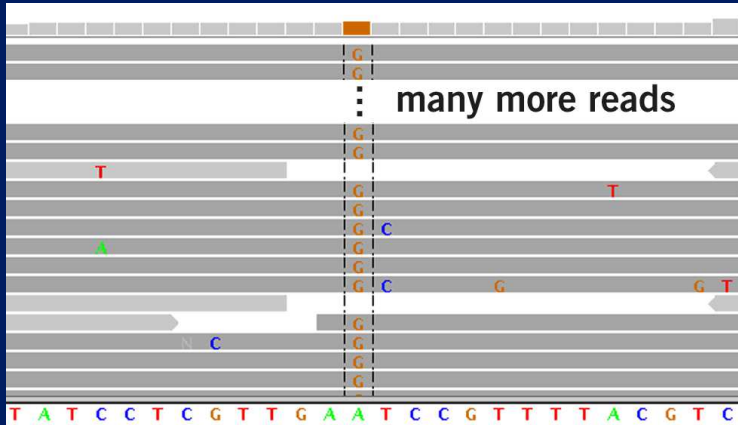


Figure 7: Result of an alignment.

Clusters

Figure 8: Dell M610 blade server

Local cluster

Some figures:

- 29 nodes.
- 368 cores.
- 94 users.

Local cluster

Some figures:

- 29 nodes.
- 368 cores.
- 94 users.

Funded by four departements:

- Molecular Epidemiology.
- Clinical Genetics.
- Human Genetics.
- Parasitology.

Storage

Funded by the same four departements.

Share	Size (TB)	Used
MolEpi	65	51
KG	27	13
HumGen	37	18
BMS	15	0
LGTC	105	89
SASC	5	0
GoNL	140	105
UCSC-bam	1	1
total	520	372

Table 1: Usage of the storage.

Resequencing

Exome:

- Only look at the genes.
- Will not detect everything.

Resequencing

Exome:

- Only look at the genes.
- Will not detect everything.

Full genome:

- Analyse everything.

Resequencing

Exome:

- Only look at the genes.
- Will not detect everything.

Full genome:

- Analyse everything.

type	desktop	cluster
exome	4 days	5 hours
genome	one year	3 days

Table 2: Gain of using a cluster.

Large projects

Genome of the Netherlands:

- 750 full genomes.

Large projects

Genome of the Netherlands:

- 750 full genomes.

Centre for Genome Diagnostics:

- 10 full genomes.

Large projects

Genome of the Netherlands:

- 750 full genomes.

Centre for Genome Diagnostics:

- 10 full genomes.

Geuvadis:

- QC for 667 samples done in two days.

These types of analysis would be impossible without the cluster.

Future

Sequencing in diagnostics (figures from Rotterdam):

- 100 exomes per week.
- $100 \times 2G$.

Future

Sequencing in diagnostics (figures from Rotterdam):

- 100 exomes per week.
- $100 \times 2\text{G}$.

We will switch to full genome sequencing soon:

- $\pm 40\text{G}$ per genome.
- 4T of data per week.

Future

Sequencing in diagnostics (figures from Rotterdam):

- 100 exomes per week.
- $100 \times 2G$.

We will switch to full genome sequencing soon:

- $\pm 40G$ per genome.
- 4T of data per week.

Needs to be more reliable.



Acknowledgements:

Michiel van Galen
Martijn Vermaat
Michel Villerius
Johan den Dunnen