

**Introduction of the webinar series**

The use and validation of High Throughput sequencing (HTS) tests for diagnostics of plant pests.

Webinar 1	What is High Throughput Sequencing (HTS)?	Friday 30 <sup>th</sup> April, 2 pm-3:30 pm
Webinar 2	How to prepare your laboratory to conduct HTS tests?	Monday 3 <sup>rd</sup> May, 2 pm-3 pm
Webinar 3	How to develop, validate and routinely use HTS for diagnostic purpose?	Tuesday 4 <sup>th</sup> May, 2 pm-3:30 pm
<b>Practical training activity</b>	How to apply the guidelines to your laboratory?	Wednesday 5 <sup>th</sup> May, 2 pm to 4:30 pm. Friday 7 <sup>th</sup> May, 2 pm to 4:30 pm.



## VALITEST Webinar series

What is High Throughput Sequencing (HTS)?

**Sebastien Massart**  
 Gembloux AgroBio Tech – Liège University  
 30/04/2021




# 1. Introduction : History

90'ies                          1.000  
2000'ies                          1.000.000  
2000'ies evolution  
2000'ies evolution HTS  
2020                              100.000.000.000.000

LIÈGE université Valitest Sébastien Massart – GxABT - 3

biochemistry  
Physiology

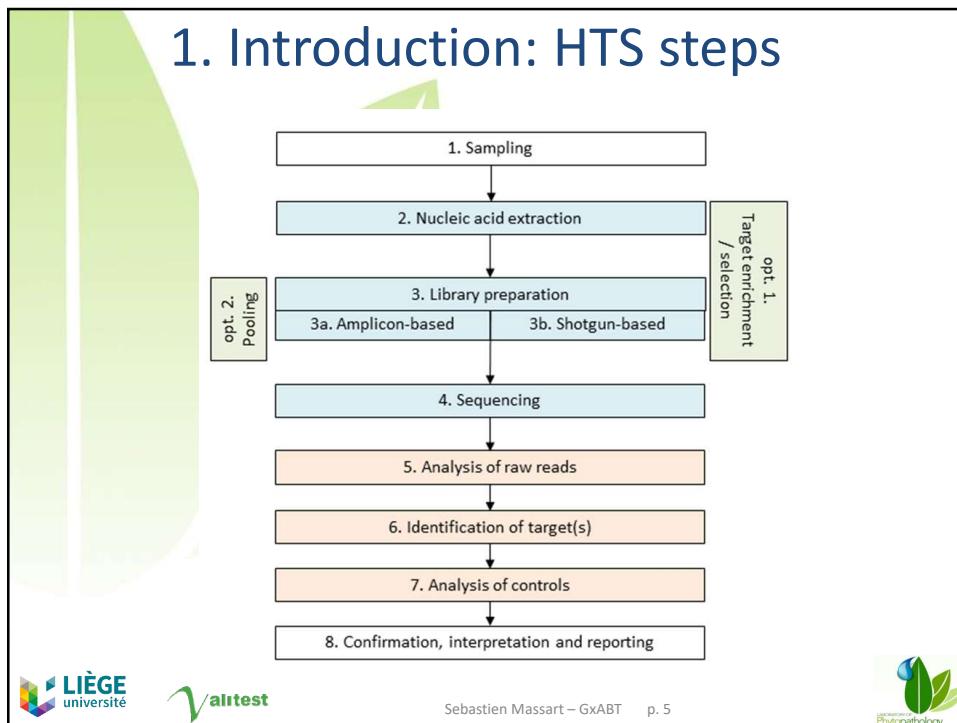
# 1. Introduction

www.wooclap.com/OCIPBU

QR code

LIÈGE université Valitest Sébastien Massart – GxABT - 4

biochemistry  
Physiology



## 2.1 Sampling

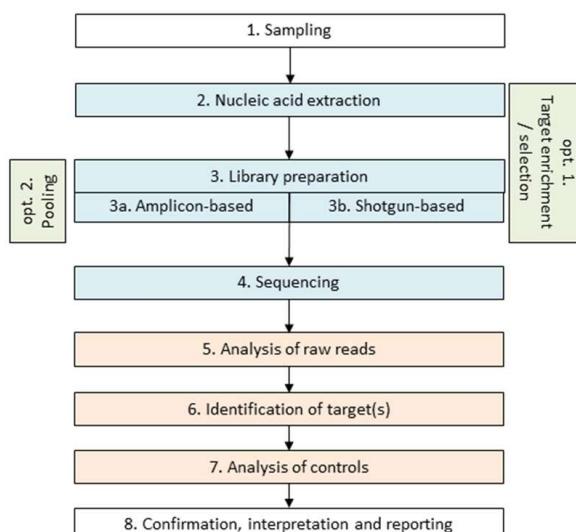
- Representative of sample status -> like any other test
- Matrix:
  1. Multiple organisms (plant, environmental samples, spore traps, insects with microbiota...)
  2. Isolated organisms (ufc)



Sebastien Massart – GxABT p. 7



## Steps of HTS technologies



Sebastien Massart – GxABT p. 8



## 2.2 Nucleic Acids extraction

- Huge diversity of protocols used
- DNA, RNA, smallRNA, dsRNA...
- Look into the M&m in the most recent publication you have read on THS (virus, fungi, bacteria...) – (5 minutes)
- What are the protocol used for N.A. extraction ?
- Are there specificities ?
- Send in the chat your observation



Sebastien Massart – GxABT - 9



## 2.2 Nucleic acids extraction

- ✓ Protocols used for NA extraction ?

[www.wooclap.com/OCIPBU](http://www.wooclap.com/OCIPBU)



Sebastien Massart – GxABT - 10



## 2.2 Nucleic Acids extraction

- Huge diversity of protocols used
- DNA, RNA, smallRNA, dsRNA...
- Not a best fit-for-all kit -> depending on your plant, tissue and experience
- Like (RT)-PCR but higher quality needed -> CONTROL !
- Quantity of RNA needed is going down  
5-10 µg -> 1 µg



Sébastien Massart – GxABT - 11

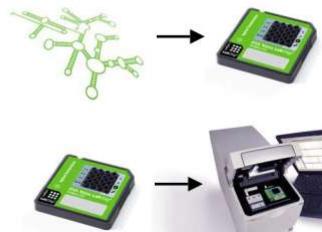


## 2.3 RNA quality control

### How the RNA Integrity Number (RIN) Works

**Step 1:**

Researchers deposit their total RNA sample into an RNA Nano LabChip.



Source: www.agilent.com

**Step 2:**

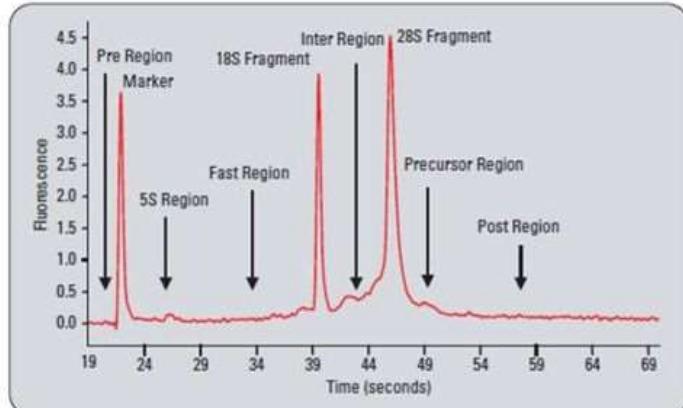
They insert the LabChip into the Agilent bioanalyzer and let the analysis run, generating a digital electropherogram.



Sébastien Massart – GxABT - 12



## 2.3 RNA quality control



**Figure 3**  
Electropherogram detailing the regions that are indicative of RNA quality.

Source: <http://www.openwetware.org/wiki/BioMicroCenter:RIN>



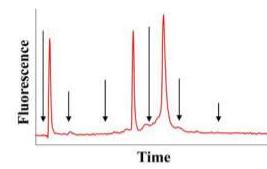
Sebastien Massart – GxABT - 13



## 2.3 RNA quality control

### Step 3:

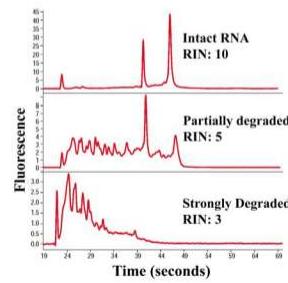
The new RIN algorithm then analyzes the entire electrophoretic trace of the RNA sample, including the presence or absence of degradation products, to determine sample integrity. Important elements of the electropherogram that are indicative of RNA quality are shown in the figure at right.



### Step 4:

The algorithm assigns a 1 to 10 RIN score, where level 10 RNA is completely intact. Because interpretation of the electropherogram is automatic and not subject to individual interpretation, universal and unbiased comparison of samples is enabled and repeatability of experiments is improved.

The RIN algorithm was developed using neural networks and adaptive learning in conjunction with a large database of eukaryote total RNA samples, which were obtained mainly from human, rat, and mouse tissues. The RIN score is largely independent of the amount of RNA used and the origin of the sample.



Source: [www.ag.unedu](http://www.ag.unedu)



Sebastien Massart – GxABT - 14



## 2.3 RNA quality control

What the RIN **can** do:

- ✓ Obtain a numerical assessment of the integrity of RNA.
- ✓ Directly compare RNA samples, e.g. before and after archival,
- ✓ compare integrity of same tissue across different labs.
- ✓ Ensure repeatability of experiments, e.g. if RIN shows a given value and is suitable for microarray experiments, then the RIN of the same value can always be used for similar experiments given that the same organism/tissue/extraction method is used.

What the RIN **cannot** do:

Tell a scientist ahead of time whether an experiment will work or not if no prior validation was done (e.g. RIN of 5 might not work for microarray experiments, but might work well for an appropriate RT-PCR experiment. Also, a RIN that might be good for a 3' amplification might not work for a 5' amplification).

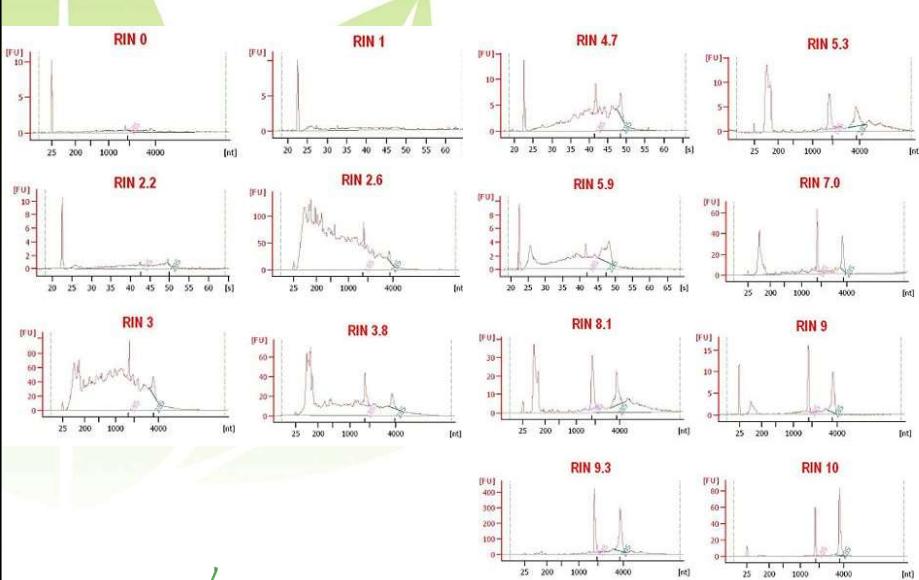
Minimal requested RIN = 7 or 8... complicated for plants !!!



Sebastien Massart – GxABT - 15



## 2.3 RNA quality control



Sebastien Massart – GxABT - 16



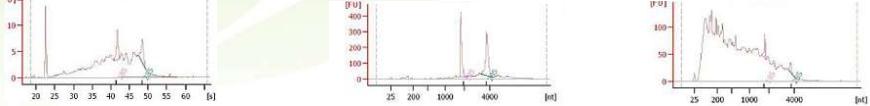
## 2.3 RNA quality control

✓ Evaluation of RIN:

[www.wooclap.com/OCIPBU](http://www.wooclap.com/OCIPBU)

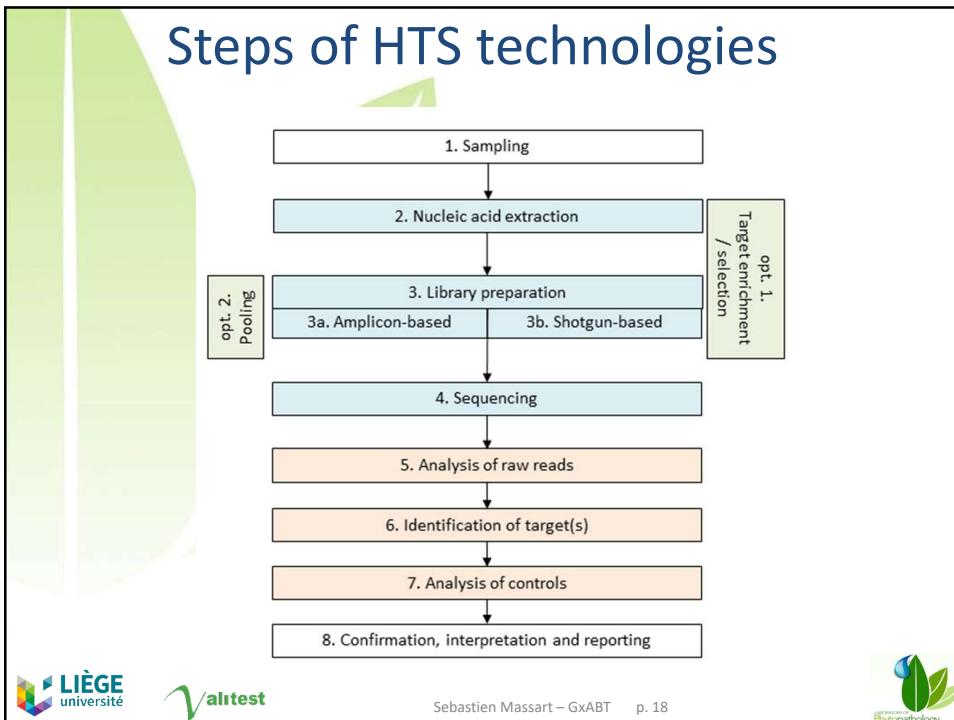


Fair                              Excellent                              Bad-very bad



Yet, viruses can still be identified from samples with low RIN quality

  Sébastien Massart – GxABT - 17 



## 2.4 Library preparation

- Format of nucleic acids **compatible** with the sequencing platform in sufficient quantity of the appropriate size



- Two groups of protocols:

1. AMPLICON SEQUENCING – PCR based

2. SHOTGUN SEQUENCING



Sebastien Massart – GxABT - 19



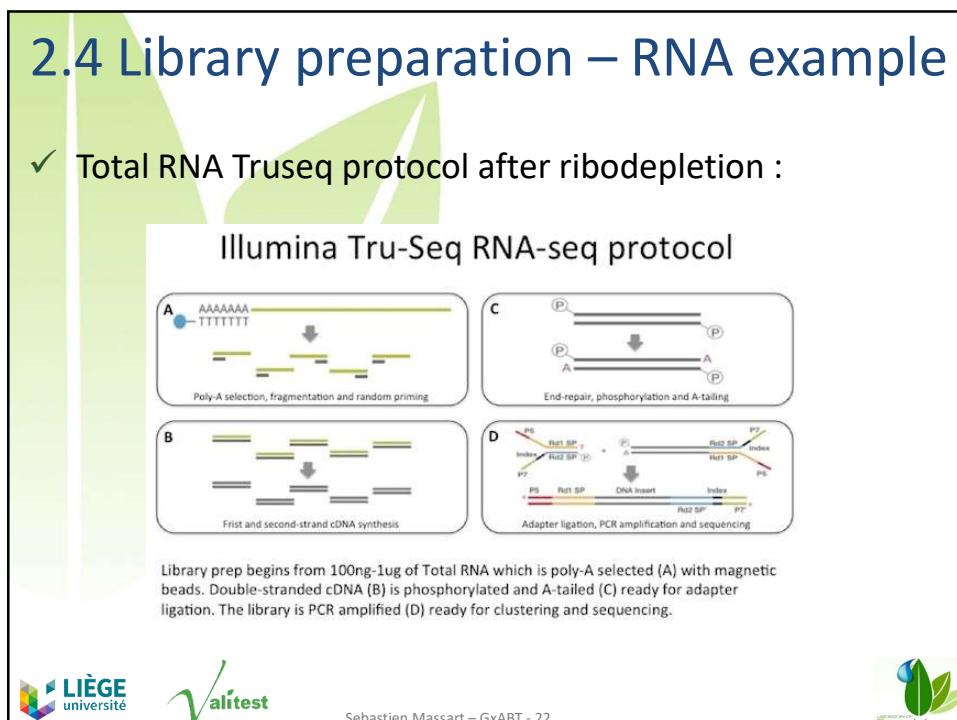
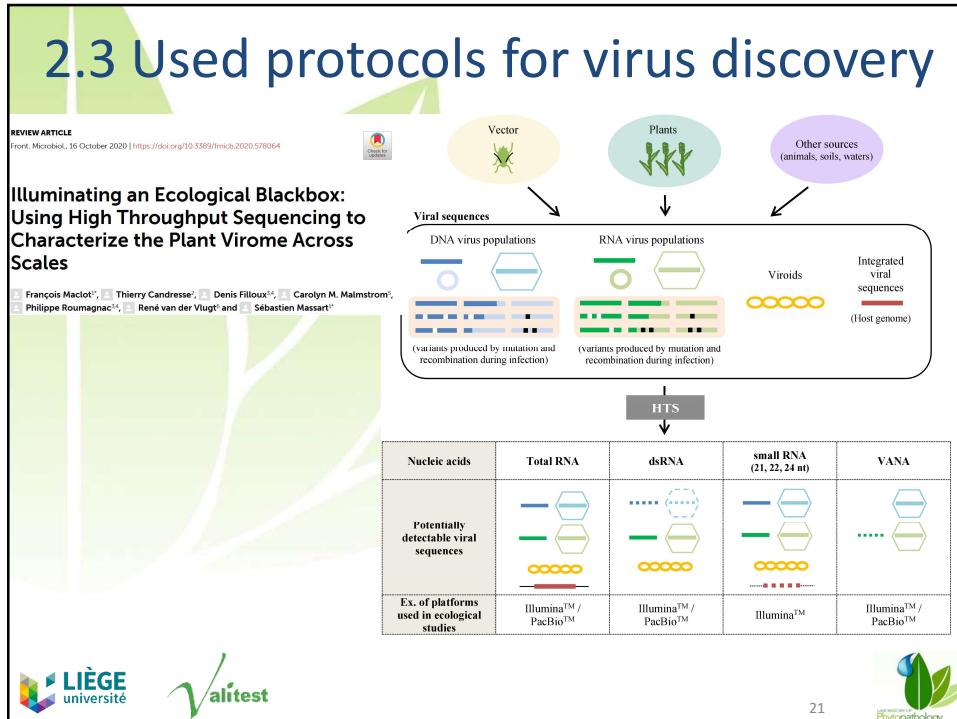
## 2.4 Library preparation - conclusion

It is only about classical molecular biology: shearing, digesting, ligating, amplifying



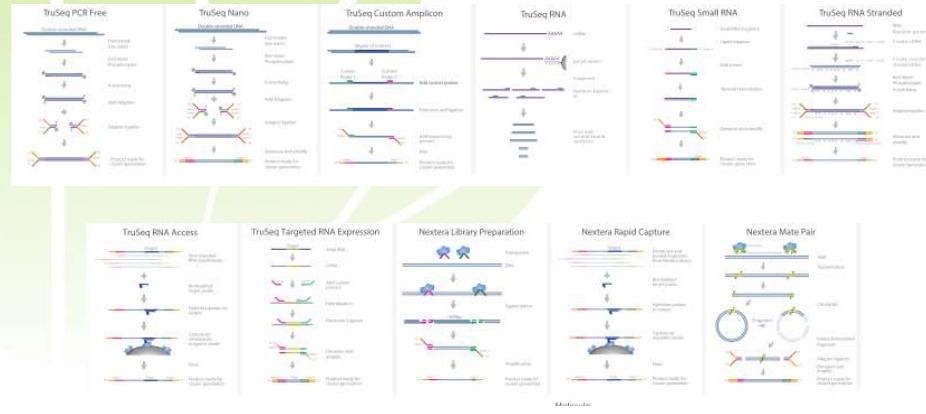
20





## 2.4 Library preparation – shotgun RNA

➤ For Illumina: 11 kits existing



Sebastien Massart – GxABT - 23



## 2.4 Library preparation – shotgun RNA

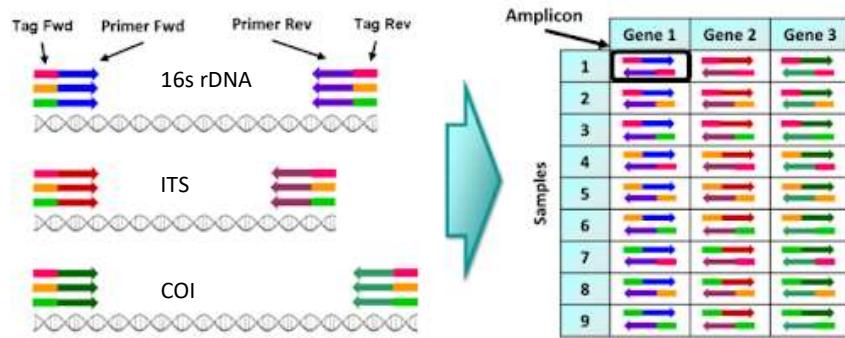
➤ For Illumina: >100 protocols



Sebastien Massart – GxABT - 24

## 2.4 Library prep. – amplicon example

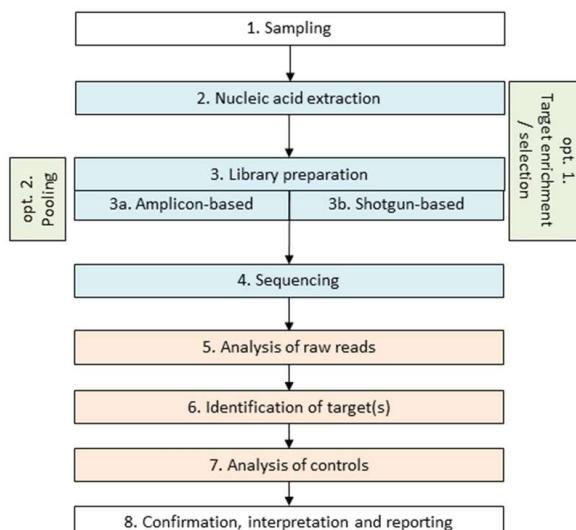
- ✓ Analysing bacterial, fungal and insect composition of a composite sample:



Sebastien Massart – GxABT - 25



## Steps of HTS technologies



Sebastien Massart – GxABT p. 26



## 2.5 Available sequencing technologies

- ✓ Illumina
- ✓ Ion Torrent
- ✓ Pacific biosciences
- ✓ Oxford Nanopore technologies



Sebastien Massart – GxABT - 27



## 2.5 Available sequencing technologies



Roche FLX, Helicos, Polonator, Complete Genomics, Solid,...



Sebastien Massart – GxABT - 28



## 2.5 Available sequencing technologies

- ✓ For all the sequencing technologies, specific protocols of library preparation have been developed
- ✓ Based on the same steps but some specificities make them incompatible



Sebastien Massart – GxABT - 29



## 2.5 Available sequencing technologies

- ✓ Four videos will be presented.
- ✓ Compare the technologies based on the following criteria
  - ✓ Detection method
  - ✓ Nucleotide used
  - ✓ Enzyme location
  - ✓ Amplification
  - ✓ Length of sequence
  - ✓ Fixation of DNA on support

[Illumina](#) - [Ion torrent](#) - [Pacific Bioscience](#) - [Oxford Nanopore](#)



Sebastien Massart – GxABT - 30



## 2.5 Available sequencing technologies

	Illumina	Ion torrent	PacBio	Minion
Detection method	Fluorescence	Electric signal	Fluorescence	Electric signal
Nucleotides	Mix	One by one	Mix	n.a.
Enzyme	Reagent	Reagent	Fixed	Fixed
Amplification	Bridge	Beads	No	No
Length	Max 300 nt	Max 300 nt	Up to 10 kb	Up to 10 kb
Fixation of DNA	Oligo on cell	Oligo on beads	Enzyme	Enzyme + nanopore



Sebastien Massart – GxABT - 31



## 2.5 Available sequencing technologies

Technology	Read length (bp)	Throughput	Reads	Runtime	Error profile	Instrument cost (US\$)	Cost per Gb (US\$, approx.)
<b>Sequencing by synthesis: CRT</b>							
Illumina MiSeq Mid output	150 (SE)	2.1–2.4 Gb	14–16 M	17 h	<1%, substitution	\$50	\$200–300
Illumina NextSeq 500/550 Mid output	75 (PE)	16–20 Gb	Up to 260 M (PE)	15 h	<1%, substitution	\$250	\$42
	150 (PE)	32–40 Gb		26 h			\$40
Illumina HiSeq X	150 (PE)	800–900 Gb per flow cell	2.6–3 B (PE)	<3 d	0.1%, substitution	\$1,000	\$7.0
<b>Sequencing by synthesis: SNA</b>							
Ion Proton	Up to 200 (SE)	Up to 10 Gb	60–80 M	2–4 h	1%, indel	\$224	\$80
Ion S5 540	200 (SE)	10–15 Gb	60–80 M	2.5 h	1%, indel	\$65	\$300
<b>Single-molecule real-time long reads</b>							
Pacific BioSciences RS II	~20 Kb	500 Mb–1 Gb	~55,000	4 h	13% single pass, ≤1% circular consensus read, indel	\$695	\$1,000
Pacific Biosciences Sequel	8–12 Kb	3.5–7 Gb	~350,000	0.5–6 h	?	\$350	?
Oxford Nanopore MK 1 MinION	Up to 200 Kb	Up to 1.5 Gb	>100,000	Up to 48 h	~12%, indel	\$1	\$750

Goodwin et al., 2016. Nature Review genetics

### Frequency of use in virus discovery and/or characterisation

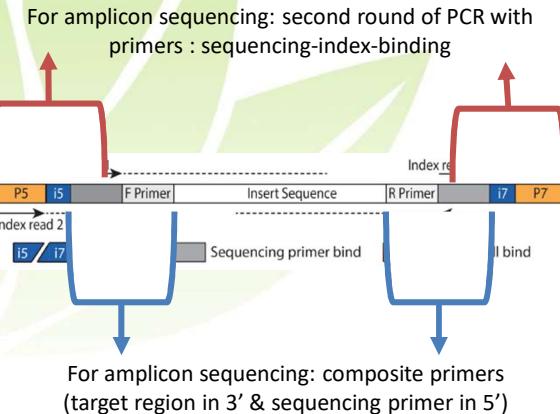


Sebastien Massart – GxABT - 32



## 2.6 Library preparation & sequencing

- ✓ Adapters at both ends for Illumina sequencing: amplicon

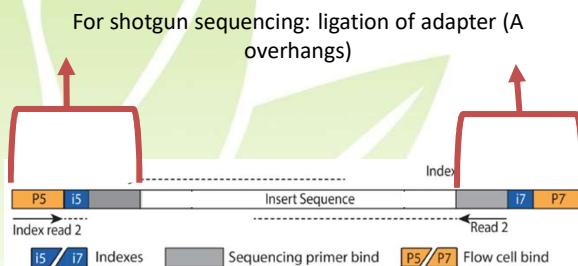


Sébastien Massart – GxABT p. 33



## 2.6 Library preparation & sequencing

- ✓ Adapters at both ends for Illumina sequencing: shotgun



Sébastien Massart – GxABT p. 34



# 3. HTS bioinformatics protocols



## 3.1 Introduction

- Balance between laboratory and bioinformatic



Sébastien Massart – GxABT - 36



### 3.1 Introduction

- Importance of bio-informatic in generation of bias:

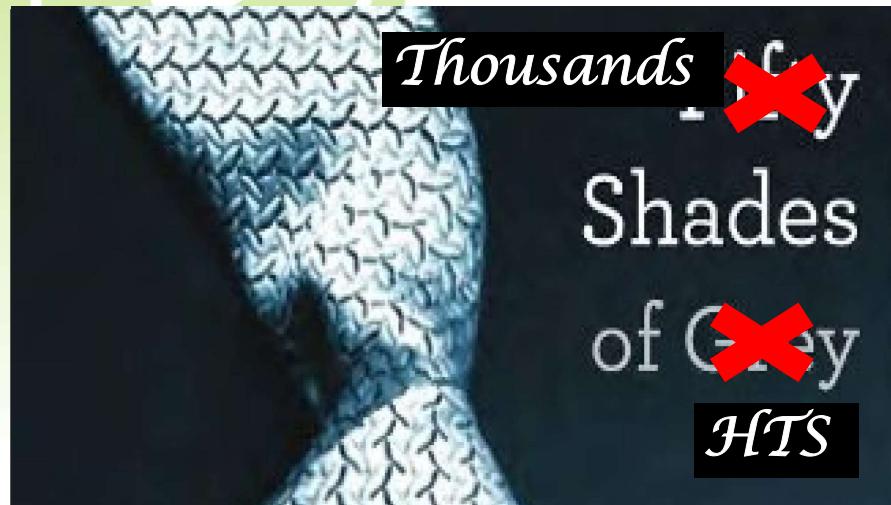


Sébastien Massart – GxABT - 37



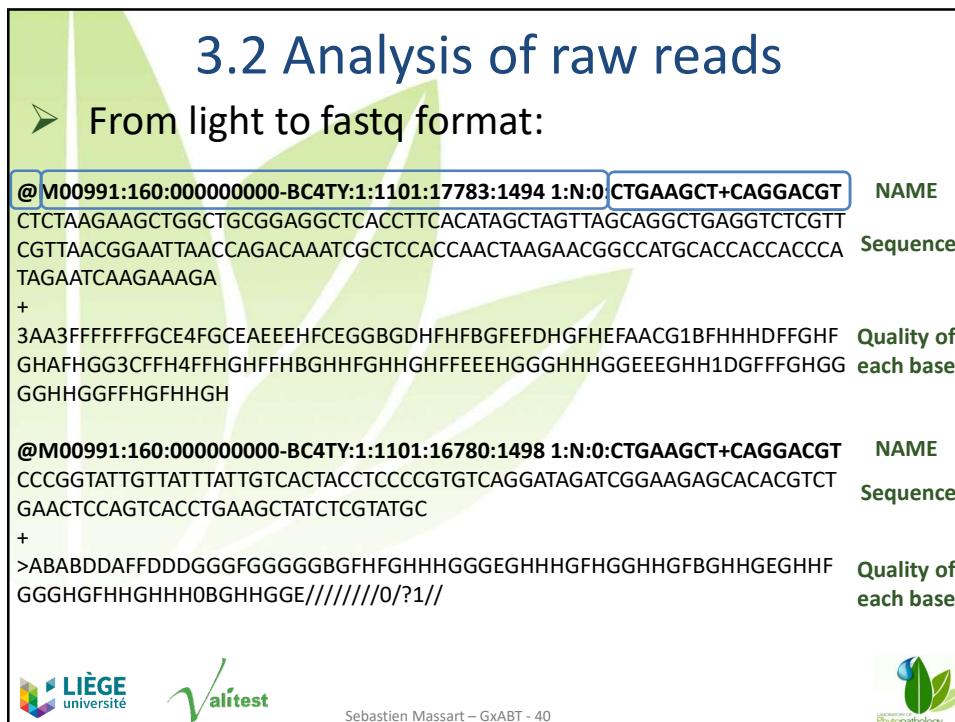
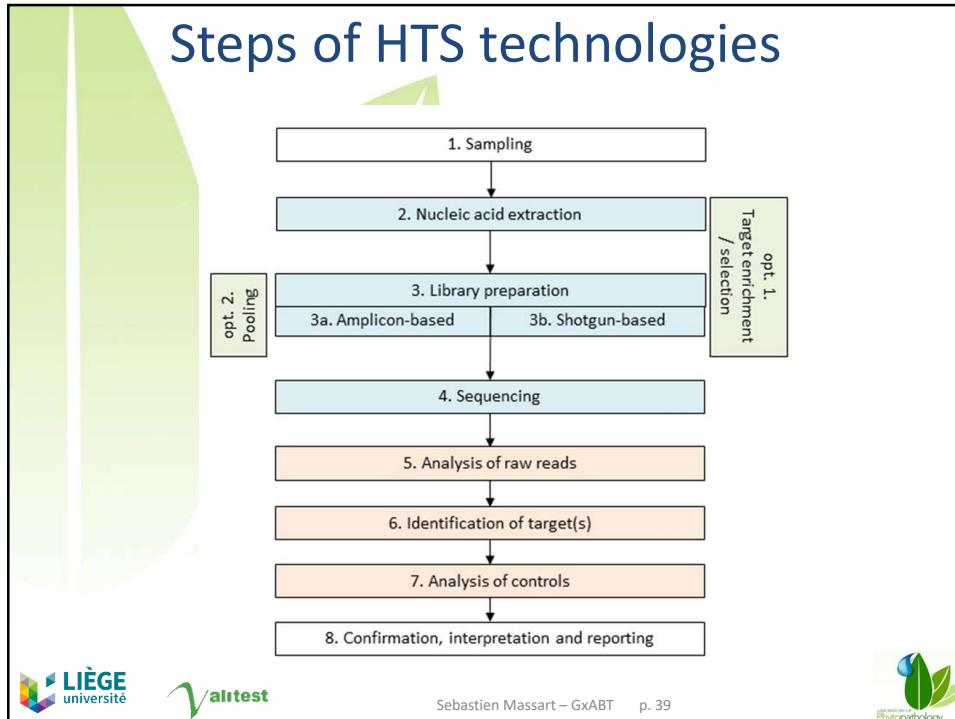
### 3.1 Introduction

- Parameter setting & result interpretation :



Sébastien Massart – GxABT - 38





## 6.2 Analysis of raw reads

- From light to fastq format:

```

MO0077:25:00000000-BI3C2:1:1101:12780:22431:N:0:C"GAAGCT*TATAAGCTT

      1       5       10      15
      |       |       |       |
Qualit: 32 32 32 32 35 35 35 36 32 32 34 35 37 34 38
      C  C  G  A  C  A  G  G  C  A  T  G  C  T  C

Qualit: 36 37 38 38 31 38 37 38 33 33 36 36 38 31 38
      A  C  A  C  T  C  G  A  A  C  C  C  T  T  C

Qualit: 37 38 39 35 31 33 37 38 35 37 38 37 16 14 37
      T  C  A  G  A  G  A  T  C  A  A  G  G  T

Qualit: 37 30 19 36 34 16 14 32 36 16 36 16 29 31 33
      C  G  G  T  C  G  G  C  G  G  T  G  C  A  C

Qualit: 16 36 35 38 14 36 36 37 39 19 39 38 38 31 14
      C  C  G  C  T  A  G  G  G  A  T  C  C  C  G

Qualit: 16 36 35 38 35 39 37 39 39 19 39 39 39 35 33
      C  C  A  A  T  C  A  G  C  T  T  C  C  T  T

```



Sebastien Massart - GxABT - 41



## 6.2 Analysis of raw reads

- Quality control of the sequences thanks to the quality score of each single base

-> elimination of reads  
-> elimination of nucleotides

- Importance ?

[www.wooclap.com/OCIPBU](http://www.wooclap.com/OCIPBU)

→ Stringency of quality control depends  
on the purpose of the HTS test

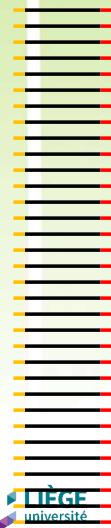


Sebastien Massart - GxABT - 42



## 6.2 Analysis of raw reads

1. Raw data  
with index



2. Raw data  
by sample

A

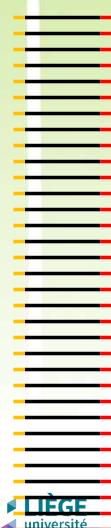
B

C

Sebastien Massart – GxABT - 43



1. Raw data  
with index



2. Raw data  
by sample

A

B

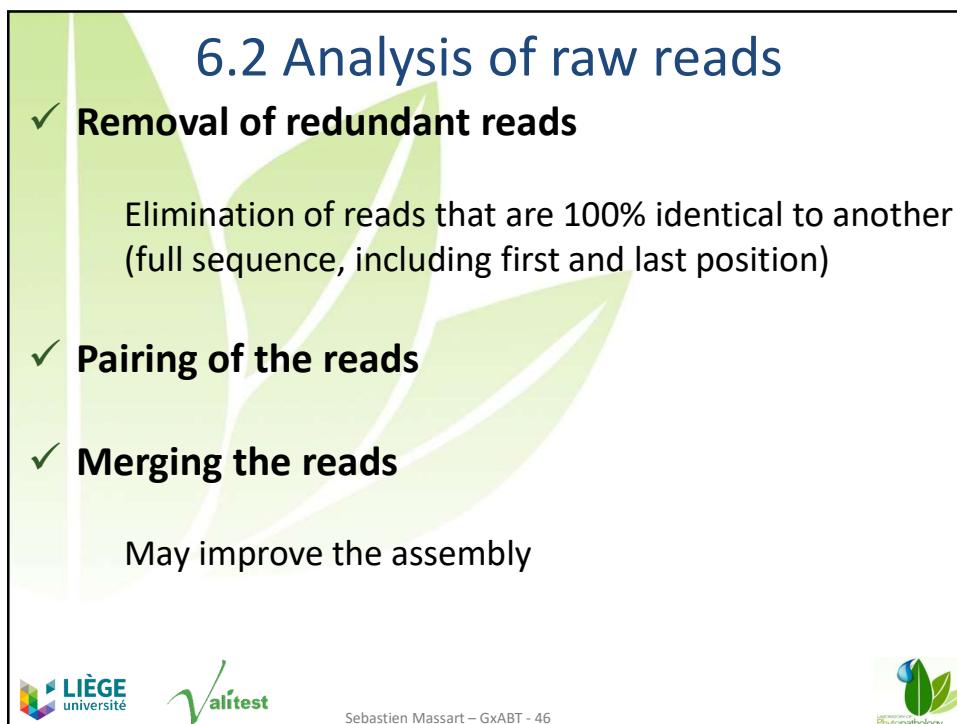
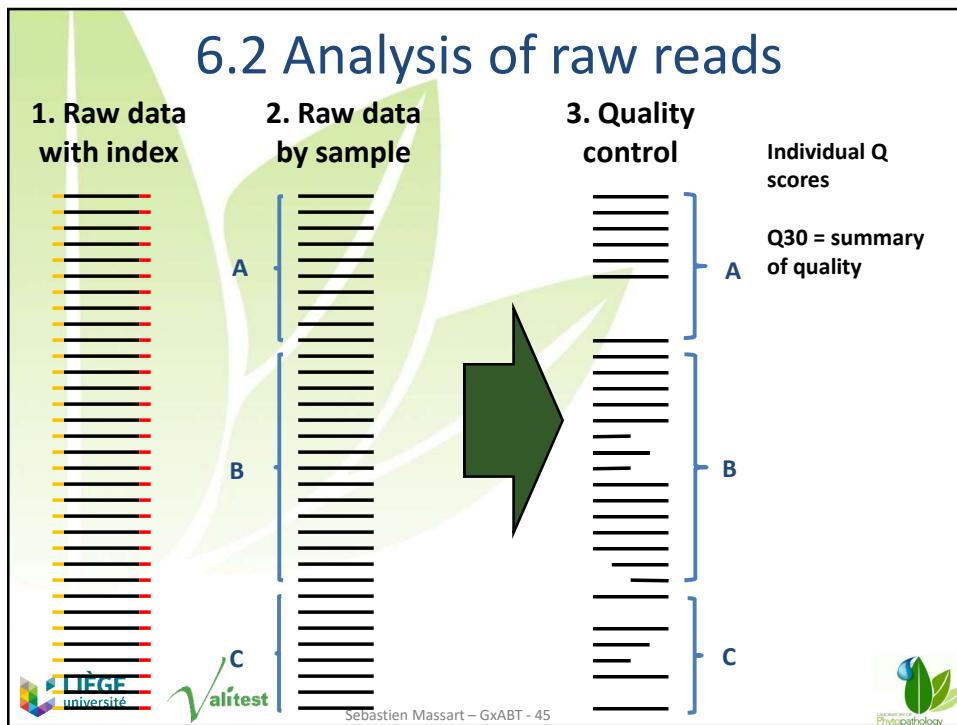
C

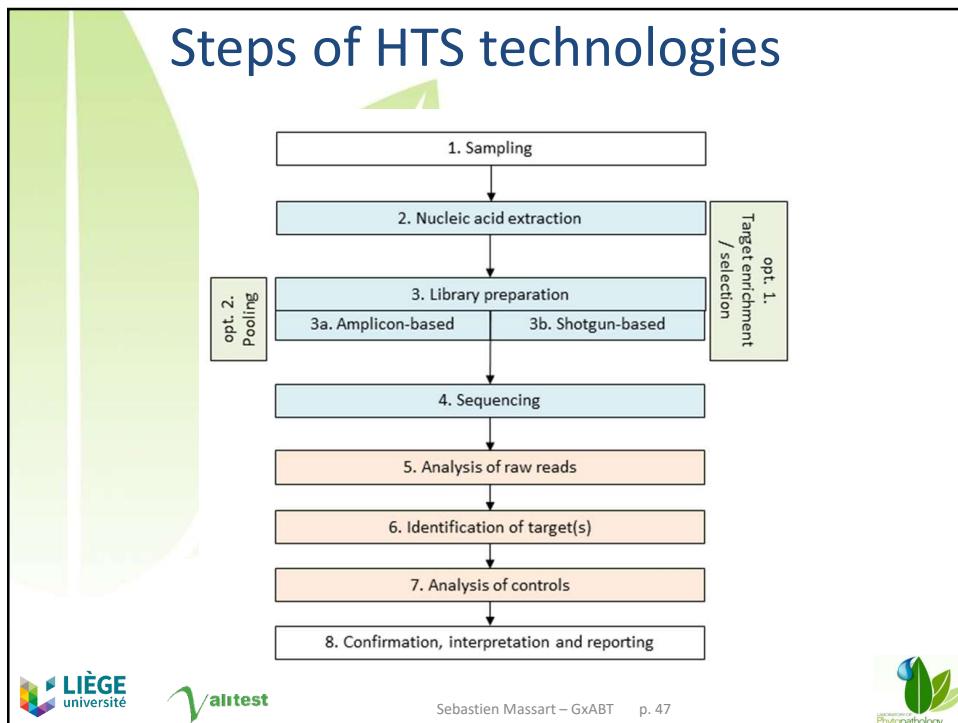
Sebastien Massart – GxABT - 44

@Sequence 1  
ATCCGACTAGATCAGATCAGAGAGGGAAA  
+  
40 40 40 40 30 20 10 10 10 40 38 10 20 24

@Sequence 10  
CGTTTAGAGAAACGACAGACTAGACAGAT  
+  
36 38 40 20 40 40 20 28 32 32 30 10 10 40







### 6.3 Identification of targets

From reads to contigs

✓ A contig (from contiguous) is a set of overlapping DNA segments that together represent a consensus region of DNA

✓ Two options:

- 1. Assembly vs. reference genome**

The diagram shows a reference genome represented by a long blue horizontal line. Below it, several short black horizontal lines represent raw reads. These reads are shown overlapping and being assembled into two larger black horizontal lines labeled 'CONTIG 1' and 'CONTIG 2'. Red 'X' marks at the ends of these contigs indicate where they do not align with the original reference genome, suggesting assembly errors or gaps.

**Logos:**

- Université de Liège logo
- Valitest logo
- Sébastien Massart - GxABT - 48
- Pathogen Physiopathology logo

## 6.3 Identification of targets

From reads to contigs

- ✓ Two options:

- 1. Assembly vs. reference genome**
- 2. De novo assembly**

CONTIG 1                    CONTIG 2                    CONTIG 3

Sebastien Massart – GxABT - 49

## 6.3 Identification of targets

From reads to contigs

Liège université Valitest Université de Liège Sébastien Massart – GxABT - 50

## 6.3 Identification of targets

From reads to contigs

De novo assembly	Mapping on genome
+	+
Reflect real genomes in the sample	Quicker ( $n \leftrightarrow 1$ )
-	-
Computing power ( $n \leftrightarrow n$ )	Bias between strains/cultivars and reference sequences

There is no opposition between both methodologies



Sebastien Massart – GxABT - 51



## 6.3 Identification of targets

Grouping the reads together

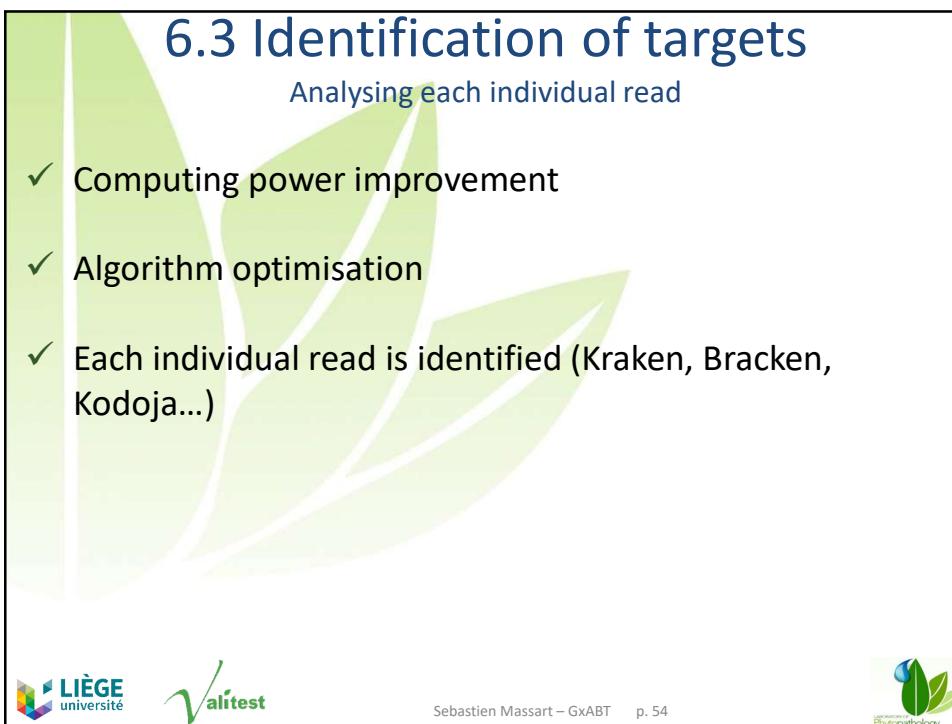
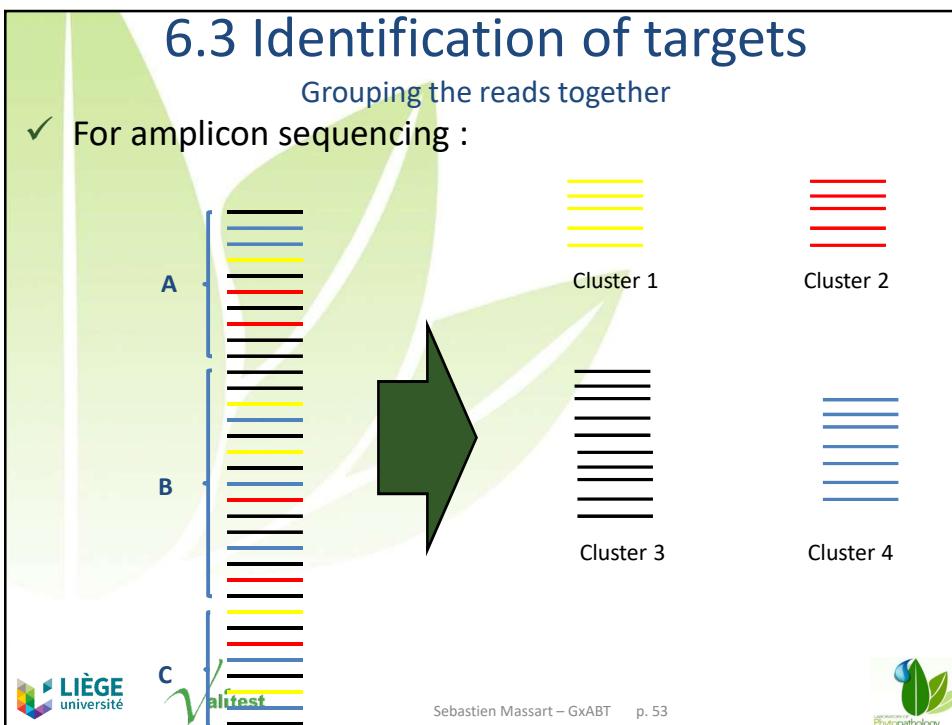
- ✓ For **amplicon sequencing**: grouping the sequences together based on their identity: example with the sequencing of 4 bacterial species in a mixed population:

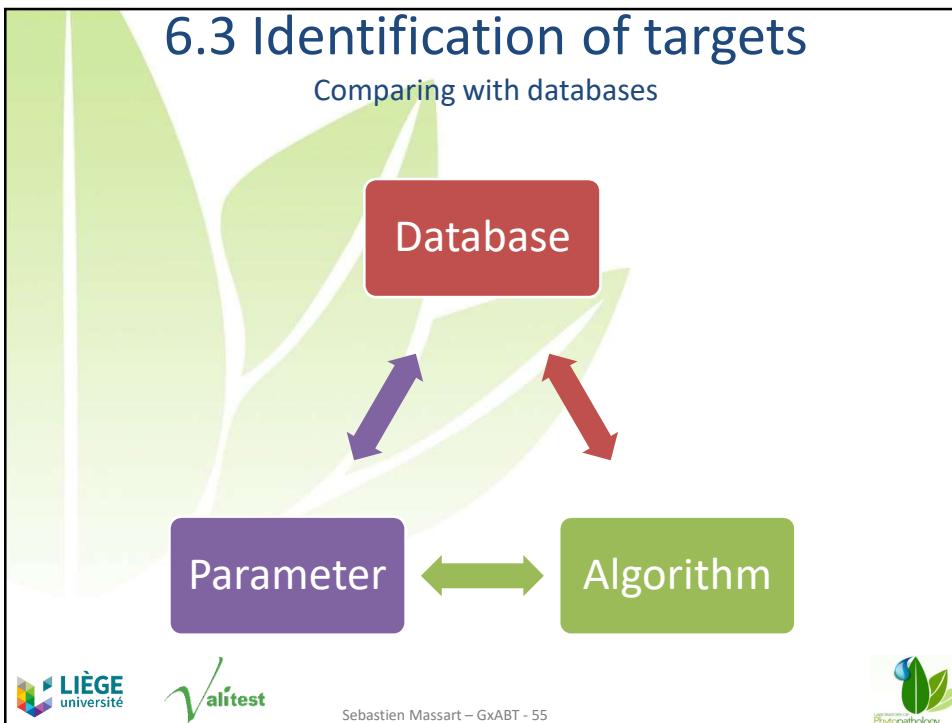
- Species 1     ——————
- Species 2     ——————
- Species 3     ——————
- Species 4     ——————



Sebastien Massart – GxABT    p. 52







## 6.3 Identification of targets

➤ How does it work ?      **Blast algorithm**

1. BLAST scans database for 'words' of a predetermined length (a 'hit') with some minimum threshold parameter, T.
2. BLAST then extends the hit until the score falls below the maximum score yet attained minus some value X.
3. From the 'hit', BLAST will extend the length of homology calculating score of homology for each additional nucleotide or aa and stop when the score is below a threshold
4. BLAST compute the probability that the homology is due to chance (depending on the lenght of the homology, the percentage of the homology and the size of the database)

## 6.3 Identification of targets

Blast algorithm

Query sequence:  
ACTGGCTGTCGTGTAAAGCACTAGATAGCTAGATCGATAG

Division in words of 5 letters

CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA

In total : 40 words of 5 letters (nt or aa) for the sequence

By default, word lenght is 28 letters (nucleotides)

Sebastien Massart – GxABT - 57

## 6.3 Identification of targets

Blast algorithm

- Millions sequences in database: same process -> cut sequences in words of same lenght
- Search for 100% homology

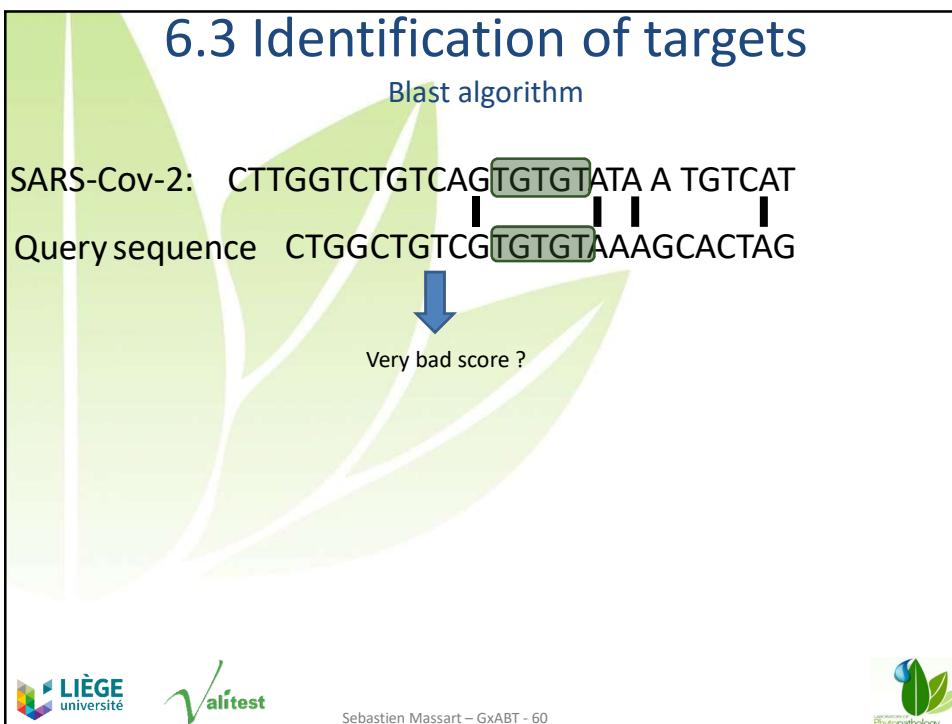
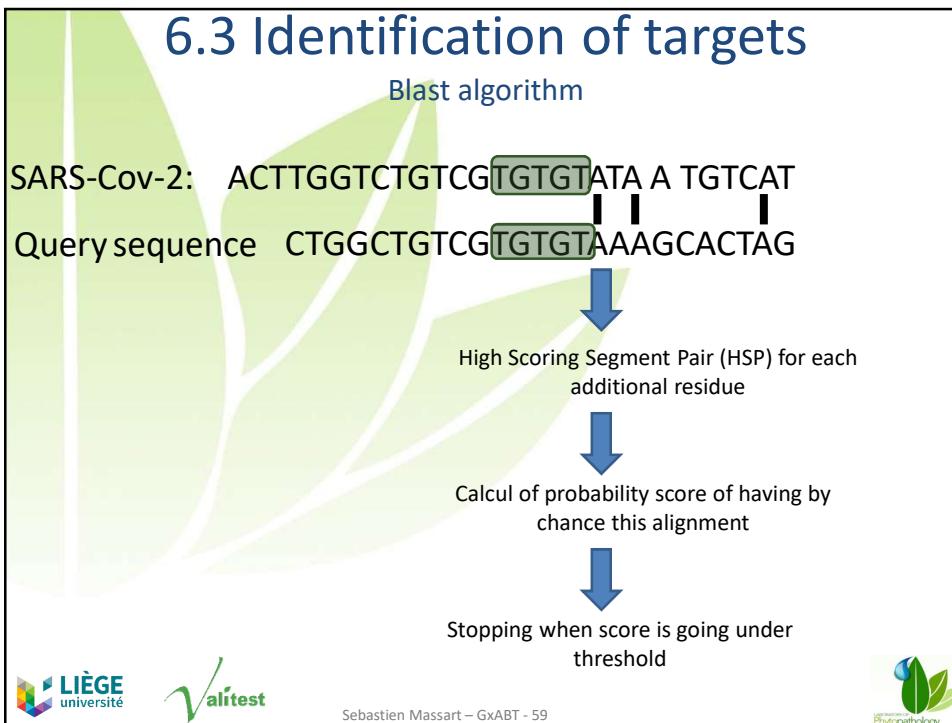
Database sequence:  
>SARS-CoV-2 genome sequence for E protein  
ACTGGCTGTCATGTCAG TGTGT ATAATGTCAT

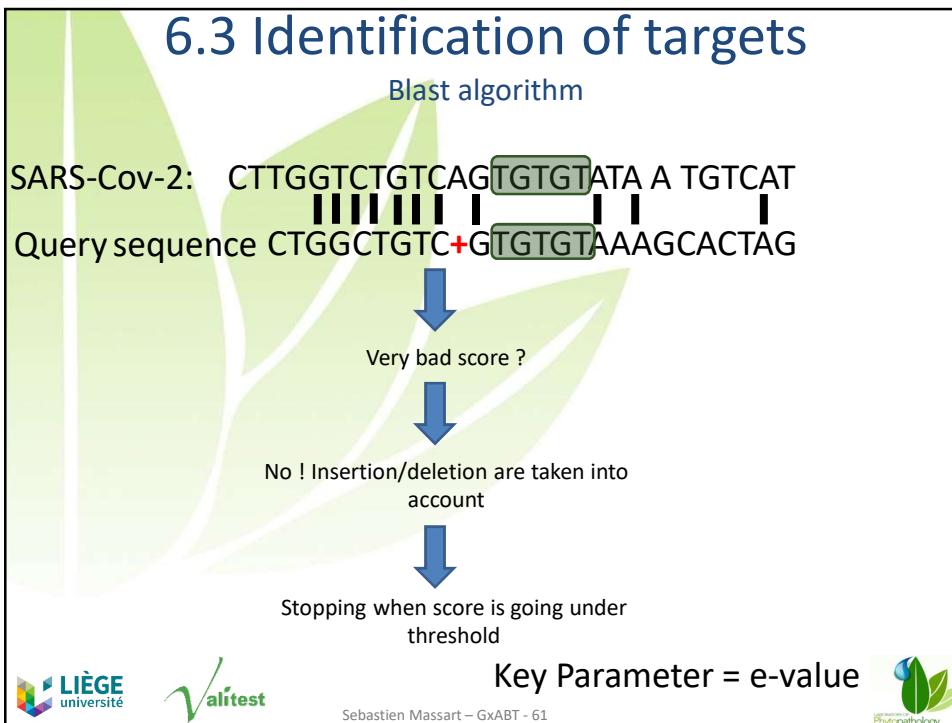
A hit is found !!!

CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA  
 CTGGC TGTCTGTGTAAAGC ACTAG ATAGCTAGATCGATA

It is a first anchoring for a further extended comparison (longer sequence)

Sebastien Massart – GxABT - 58





## 6.3 Identification of targets

Comparing with databases by blast

- Comparison of contigs (reads) vs. database

BLAST program	Query	Database
Nucleotide blast (blastn)	Nucleotide	Nucleotide
Protein blast (blastp)	Protein	Protein
Blastx	Translated nucleotide	Protein
tblastn	Protein	Translated nucleotide
tblastx	Translated nucleotide	Translated nucleotide

What are the pros and cons of translation of DNA and comparison with protein databases ?

LIEGE université Valitest Institut de Pathologie

Sebastien Massart – GxABT - 62

## 6.3 Identification of targets

Comparing with databases by blast

For each contig: 110's of hits with databases

**Usual presentation of Blast results: best hit for each contig -> OK but.....**



**Look behind the tree hiding the forest !!!!**



Sebastien Massart – GxABT - 63



## 6.3 Identification of targets

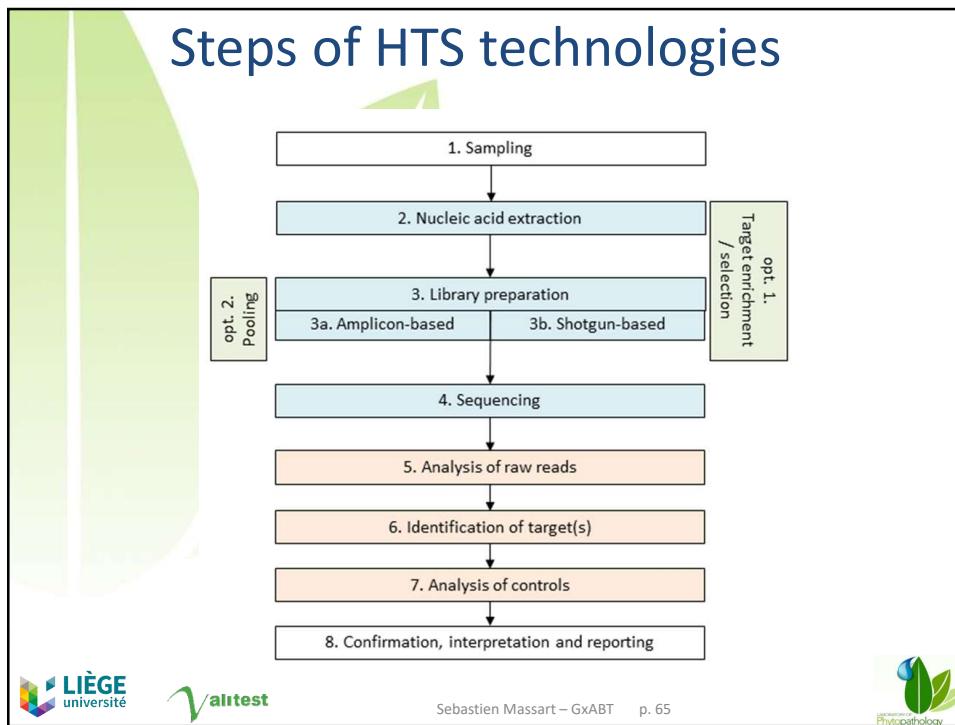
Comparing with databases by blast

➤ Some elements related to Blast :



Sebastien Massart – GxABT - 64





## 6.4 Biological interpretation

**frontiers**  
in Microbiology

PERSPECTIVE  
published: 24 January 2017  
doi: 10.3389/fmicb.2017.00045

**A Framework for the Evaluation of Biosecurity, Commercial, Regulatory, and Scientific Impacts of Plant Viruses and Viroids Identified by NGS Technologies**

**Sebastien Massart<sup>1\*</sup>, Thierry Candresse<sup>2</sup>, José Gil<sup>3</sup>, Christophe Lacomme<sup>4</sup>, Lukas Predajna<sup>5</sup>, Maja Ravnikar<sup>6</sup>, Jean-Sébastien Reynard<sup>7</sup>, Artemis Rumbou<sup>8</sup>, Pasquale Saldarelli<sup>9</sup>, Dijana Škorić<sup>10</sup>, Eeva J. Vainio<sup>11</sup>, Jari P. T. Valkonen<sup>12</sup>, Hervé Vandershuren<sup>13</sup>, Christina Varveri<sup>14</sup> and Thierry Wetzel<sup>15</sup>**

<sup>1</sup> Plant Pathology Laboratory, Gembloux Agro-Bio Tech, University of Liège, Gembloux, Belgium, <sup>2</sup> Institut National de la Recherche Agronomique (INRA), University of Bordeaux, CS20032 UMR 1332 BPF, Villeneuve d'Ornon, France, <sup>3</sup> Plant Biology, Linnaeus Centre for Plant Biology, Uppsala BioCentre, Swedish University of Agricultural Sciences, Uppsala, Sweden, <sup>4</sup> Virology and Zoology, Science and Advice for Scottish Agriculture, Edinburgh, UK, <sup>5</sup> Department of Plant Virology, Institute of Virology, Biomedical Research Center, Slovak Academy of Science (SAS), Bratislava, Slovakia, <sup>6</sup> Department of Biotechnology and Systems Biology, National Institute of Biology, Ljubljana, Slovenia, <sup>7</sup> Virology, Agroscope, Nyon, Switzerland, <sup>8</sup> Division Phytomedicine Lanzesella, Faculty of Life Sciences, Albrecht Daniel Thaer-Institute of Agricultural and Horticultural Sciences, Humboldt-University of Berlin, Berlin, Germany, <sup>9</sup> National Research Council Institute for Sustainable

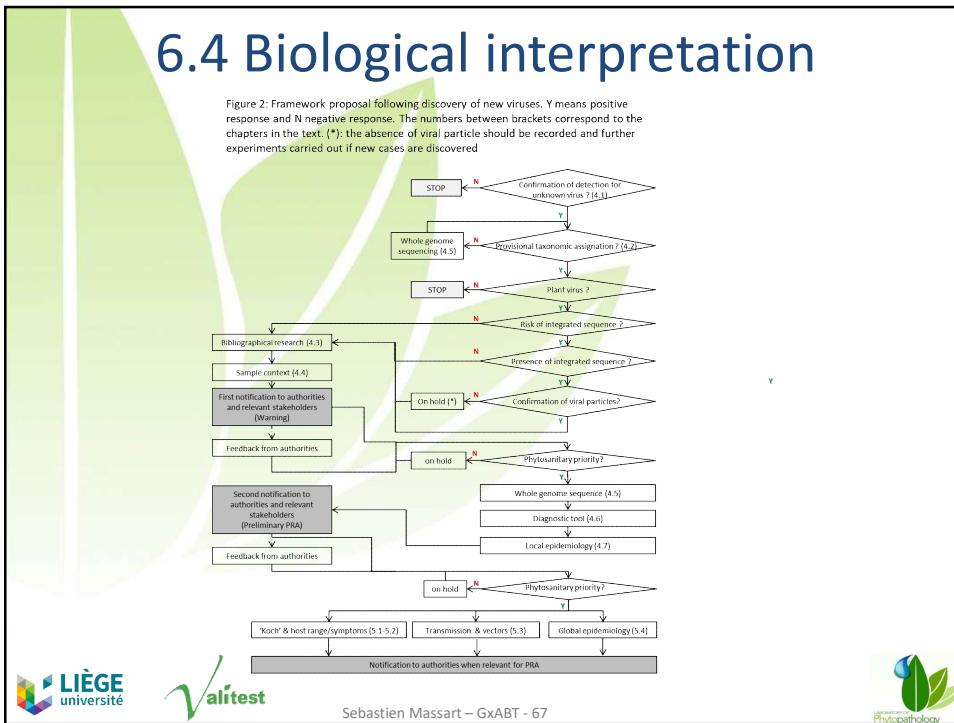
OPEN ACCESS

Edited by:  
David Gilmer,  
University of Strasbourg, France

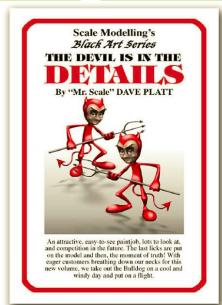
Liège université  
Valitest

Sébastien Massart – GxABT - 66

**Plant Pathology**  
Université de Liège



1. Reviewing all the steps at both lab and bioIT levels
2. Diversity of protocols causing all bias



**You can miss a pest in the sample because of library preparation method, sequencing depth, bioinformatic software and parameters !!**



Sebastien Massart – GxABT - 69



1. Reviewing all the steps at both lab and bioIT levels
2. Diversity of protocols causing all bias
3. Standardisation of protocols
  - Library preparation (Illumina kits) :



Sebastien Massart – GxABT - 70



1. Reviewing all the steps at both lab and bioIT levels
2. Diversity of protocols causing all bias
3. Standardisation of protocols
  - Library preparation (Illumina kits)
  - Bioinformatic (smallRNA):

METHODOLOGY ARTICLE | OPEN ACCESS

An internet-based bioinformatics toolkit for plant biosecurity diagnosis and surveillance of viruses and viroids

Roberto A. Barrero<sup>†</sup>, Kathryn R. Napier<sup>†</sup>, James Cunningham, Lia Liefting, Sandi Keenan, Rebekah A. Frampton, Tamas Szabo, Simon Bulman, Adam Hunter, Lisa Ward, Mark Whattam and Matthew I. Bellgard<sup>✉</sup>

<sup>†</sup>Contributed equally

BMC Bioinformatics BMC series – open, inclusive and trusted 2017 18:26 | <https://doi.org/10.1186/s12859-016-1428-4>  
© The Author(s). 2017

Received: 12 July 2016 | Accepted: 15 December 2016 | Published: 11 January 2017



Sebastien Massart – GxABT - 71



### 3. Standardisation of protocols:



Sebastien Massart – GxABT - 72



Thank you for your attention!

E-mail:

[sebastien.massart@uliege.be](mailto:sebastien.massart@uliege.be)

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N° 773139

The content of this presentation represents the views of the author only and is his/her sole responsibility; it cannot be considered to reflect the views of the European Commission and/or the Research Executive Agency or any other body of the European Union. The European Commission and the Agency do not accept any responsibility for use that may be made of the information it contains.

