



Anais POIRAUDEAU<sup>1\*</sup>, Maxime MANNO<sup>1\*</sup>, Céline VANDECASTEELE<sup>1</sup>, Alain ROULET<sup>1</sup>, Celine ROQUES<sup>1</sup>, Marie VIDAL<sup>1</sup>, Catherine ZANCHETTA<sup>1</sup>, Pauline HEUILLARD<sup>1</sup>, Fabrice ROUX<sup>2</sup>, Baptiste MAYJONADE<sup>2</sup>, Jérôme GOUZY<sup>2</sup>, Yann GUIGUEN<sup>3</sup>, Christophe KLOPP<sup>4</sup>, Pierre FRASSE<sup>5</sup>, Mohamed ZOUINE<sup>5</sup>, Cécile DONNADIEU<sup>1</sup>, Olivier BOUCHEZ<sup>1</sup>, Gérald SALIN<sup>1</sup> and Claire KUCHLY<sup>1</sup>

<sup>1</sup> INRA, US 1426, GeT-PlaGe, Genotoul, Castanet-Tolosan, France

<sup>2</sup> LIPM, UMR 441 INRA, 31326, Castanet-Tolosan, France

<sup>3</sup> INRA, UR 1037 LPGP Laboratoire de Physiologie et Génomique des Poissons, 35000, Rennes, France.

<sup>4</sup> MIAT, UR875 INRA, Plateforme bioinformatique, 31326, Castanet-Tolosan, France

<sup>5</sup> ENSAT, UMR 990 INRA/INPT-ENSAT, Laboratoire Génomique et Biotechnologie des Fruits, 31326, Castanet-Tolosan, France

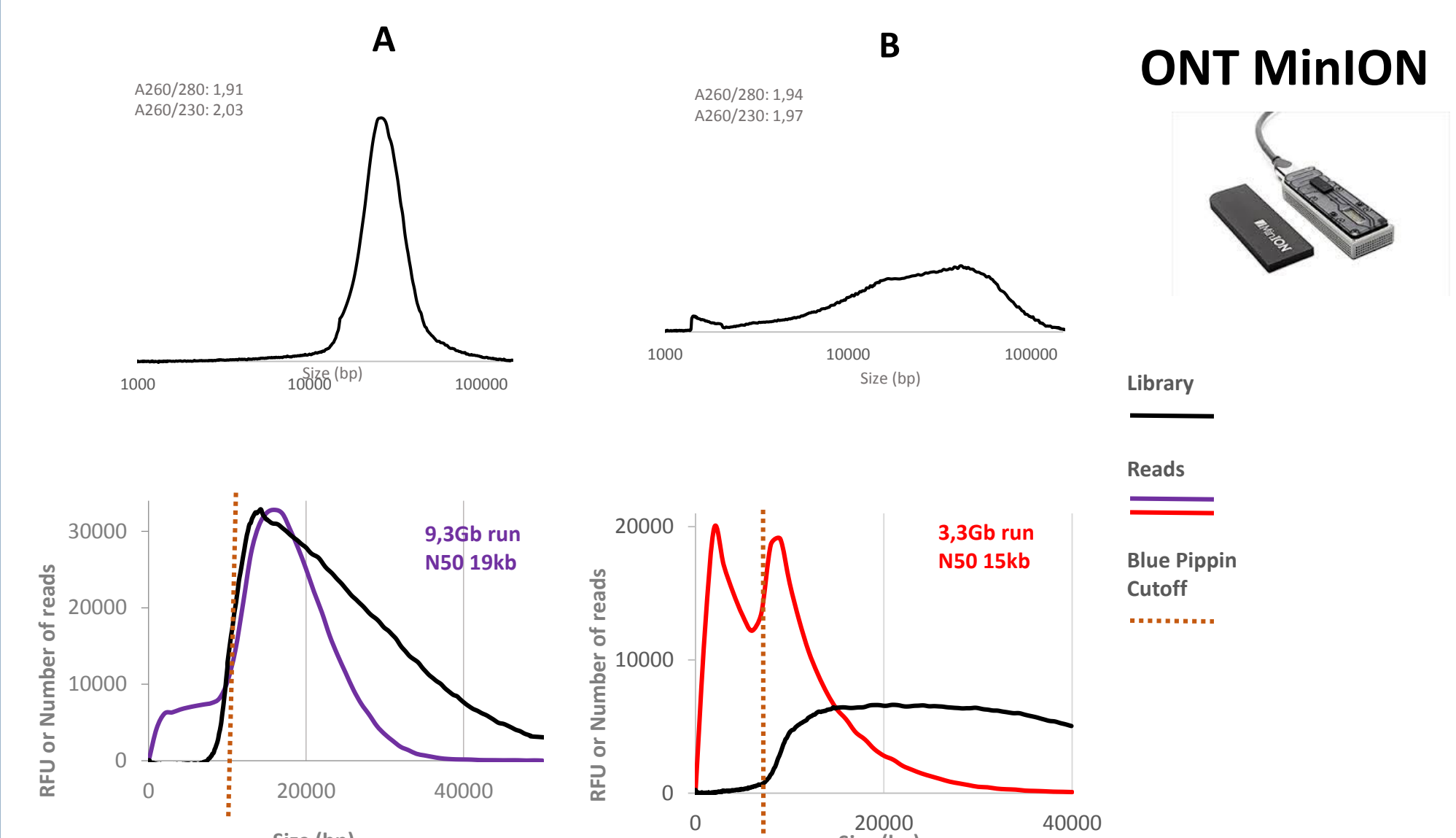
\* Equal contribution to this work

Thanks to its experience on short reads sequencing using the Illumina technologies, the GeT-PlaGe core facility began to evaluate and use long read technologies since the beginning of 2015: Pacific BioSciences RSII, Oxford Nanopore Technology MinION and 10XGenomics Chromium. Genomic issues such as complex genome assembly, structural variant discovery or phasing can be addressed by those long read technologies.

As DNA quality is the most important requirement to obtain an efficient sequencing, sample requirements for each technology, and quality controls performed on GeT-PlaGe will be detailed. For all the technology presented, DNA sample needs high quality and purity. For ONT MinION and 10XGenomics Chromium, the reads length have theoretically no limits compare to PacBio RSII (max around 50 kb) but the input DNA size is the key for all of them. The amount of DNA required for sequencing can be huge and challenging to obtain. The DNA quality is the cornerstone of a good bioinformatic analysis particularly for assemblies.

We are presenting current projects concerning *de novo* assembly results obtained using multiple Long Read technologies, for several genomes (bacteria, fungus, tomato and fish).

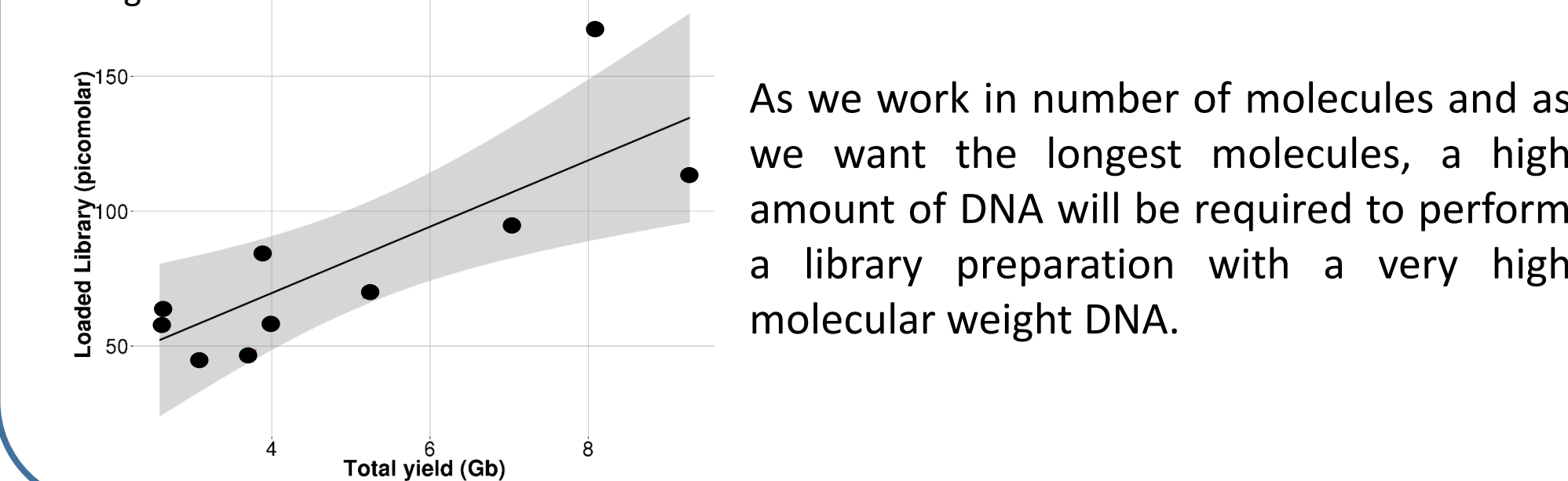
### Impact of the DNA quality on the long read sequencing efficiency



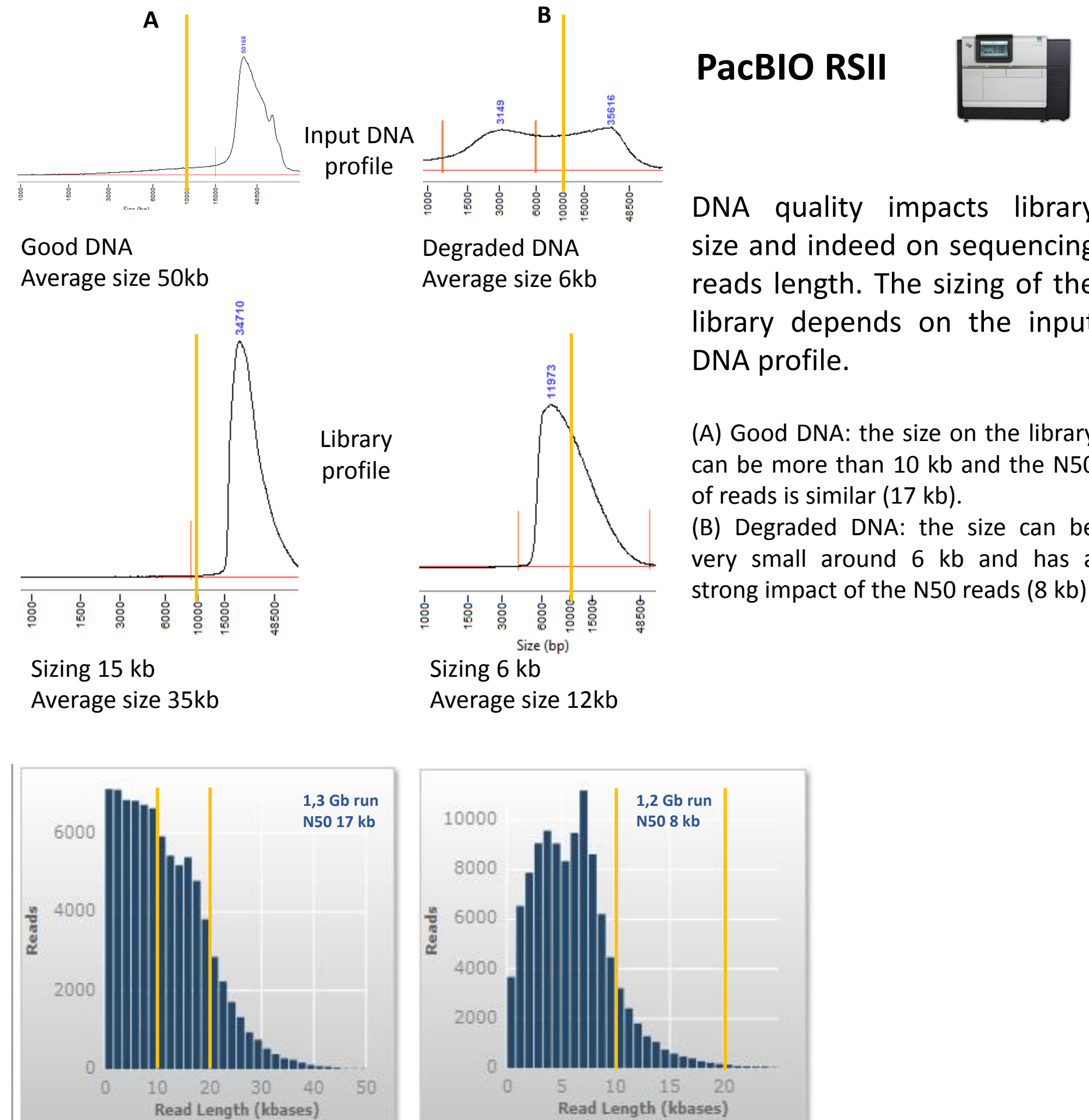
Impact of DNA quality (i.e. degradation) on sequencing efficiency

(A) Non degraded DNA: when a suitable library preparation is performed, the read profile will be very similar to the library profile

(B) Degraded DNA: we will be able to perform a library preparation but a high amount of small reads will be generated



As we work in number of molecules and as we want the longest molecules, a high amount of DNA will be required to perform a library preparation with a very high molecular weight DNA.



#### PacBio RSII



DNA quality impacts library size and indeed on sequencing reads length. The sizing of the library depends on the input DNA profile.

(A) Good DNA: the size on the library can be more than 10 kb and the N50 of reads is similar (17 kb).

(B) Degraded DNA: the size can be very small around 6 kb and has a strong impact of the N50 reads (8 kb)

#### 10X Genomics Chromium

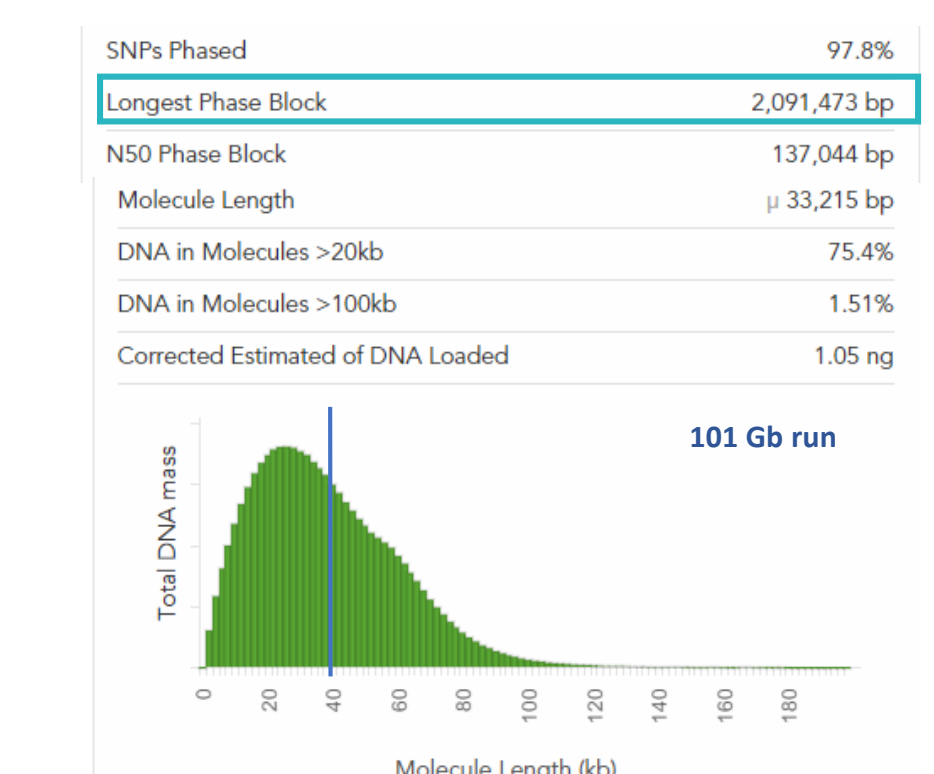
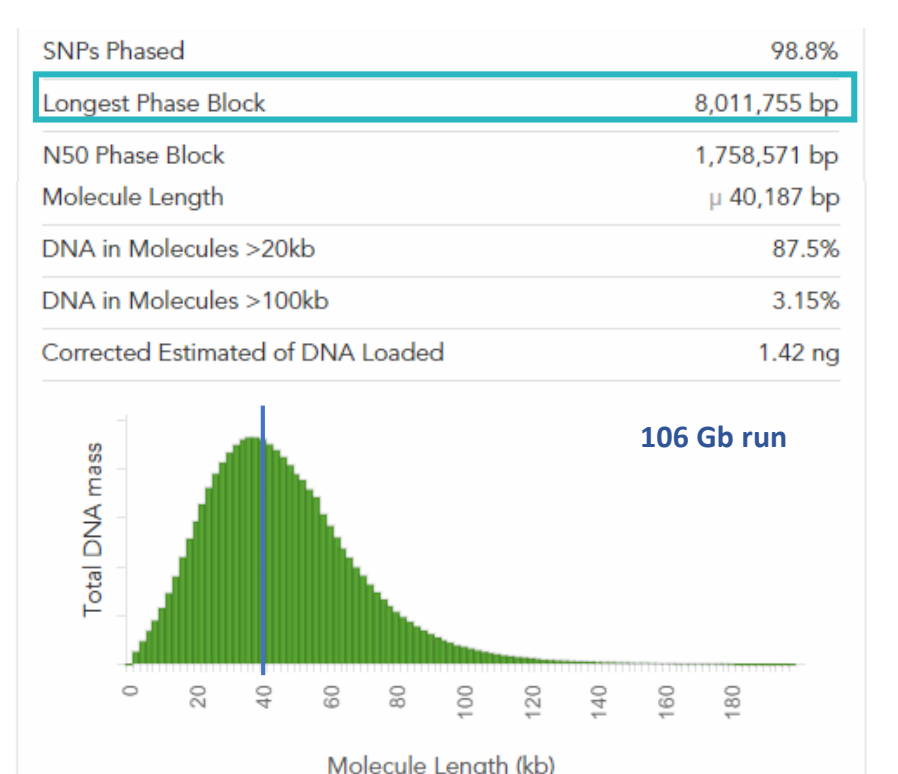
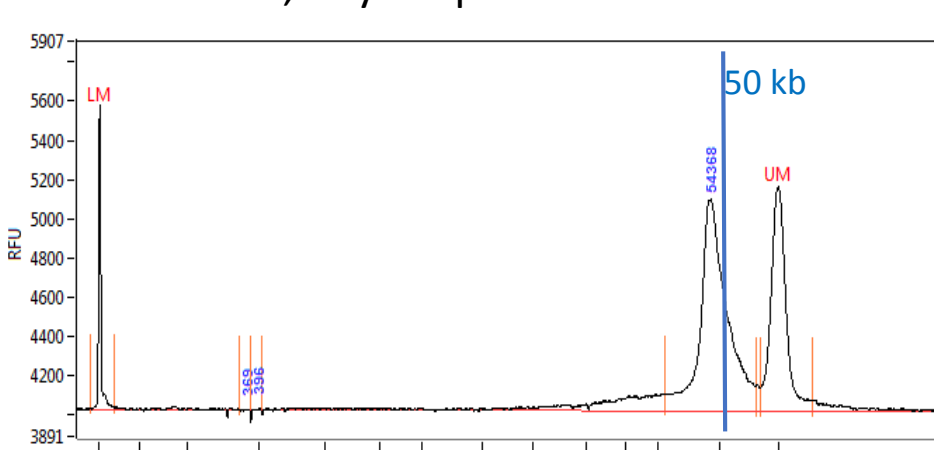
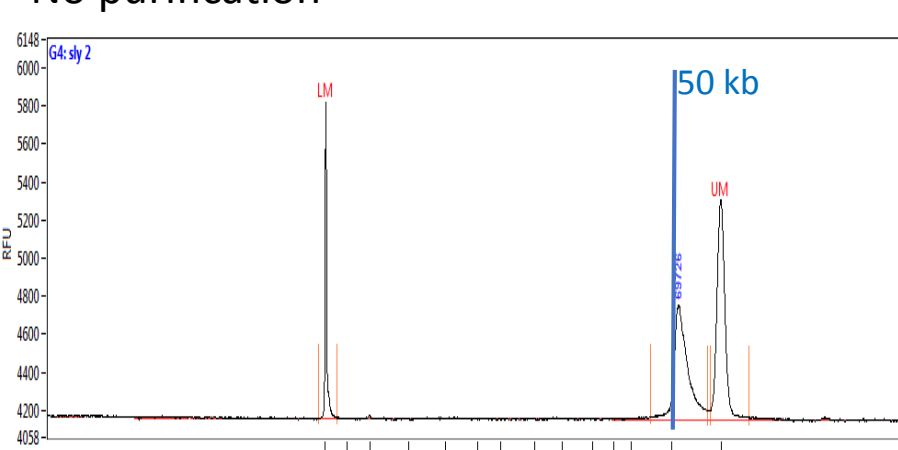


The 10X Chromium technology requires fragments greater than 50 kb. For *de novo* assemblies, it recommends 100kb fragments.

DNA concentration: 3,70 ng/μL  
83,5% size > 50 kb  
No purification



DNA concentration: 11,6 ng/μL  
37,4% size > 50 kb  
Purification 1,5 by Ampure XP beads



With the 10X Chromium, DNA quality and large fragment size are essential for a good use of the technology and impact the size of the reconstructed fragments and the size of the phased blocks.

### ONT MinION vs PacBio RSII: let's fight !

Concerning the problematic of *de novo* assembly using long reads, PacBio technologies are challenged by ONT's. Can we compare them on the same input DNA (Bacteria and Fungus) ?

#### Bacterial raw data accuracy

ONT raw data accuracy seems to be better than PacBio's, but errors are not random (long homopolymer bias). It might be an obstacle to obtain a good assembly. Can we use the ONT technology alone or combined with Illumina data to get a genome with the equivalent quality than PacBio ?

#### Bacterial genome assembly (5 Mb)

Assembly sets	Input Cov <sup>1</sup>	nbrContigs	totNuc
MinION-CANU <sup>2</sup>	1156	1	5045252
MinION-CANU-PILON	71 <sup>3</sup>	1	5122723
RSII-HGAP <sup>3</sup>	81	1	5078836

<sup>1</sup> Raw data coverage (X)  
<sup>2</sup> ONT raw data (1D protocol), qscore > Q10 and size > 3 kb  
<sup>3</sup> Illumina data coverage  
<sup>4</sup> PacBio raw data, polymerase qscore > 0.80 and subread size > 3 kb

The ONT read length allow to get just one contig for bacteria like the other sets. ONT Assembly has been performed with all raw data but we didn't assess with less coverage. From basecalling to assembly, the PacBio RSII pipeline is more efficient than the ONT pipeline which is frequently updated.

#### Bacterial assembly completeness<sup>1</sup> metrics

Assembly sets	Complete genes	Fragmented genes	Missing genes	Total genes groups searched
MinION-CANU	25 (16.9%)	53 (35.8%)	70 (47.3%)	148
MinION-CANU-PILON	141 (95.3%)	0	7 (4.7%)	148
RSII-HGAP3	141 (95.3%)	0	7 (4.7%)	148

<sup>1</sup> Completeness computed with BUSCO V2 on the bacteria\_odb9

The best BUSCO results are given by PacBio and ONT+Illumina data: 95.3% of complete genes. While ONT data allow us to generate only one contig, most of the genes are missing or fragmented. With these versions of MinKNOW and Albacore, ONT data have to be combined with Illumina's to get the same results as PacBio RSII.

#### Fungus genome assembly (50Mb)\*

Assembly sets <sup>1</sup>	N50 raw data	Quantity of raw data	Nbr contigs	N50 contigs	Total size of contigs	Complete genes
RSII	9.7Kb	2.9 Gb	915	166 kb	84 Mb	92%
MinION LR <sup>2</sup>	19Kb	2.3 Gb	408	312 kb	68 Mb	89%
MinION LR+8 Kb <sup>3</sup>	13.5Kb	3.8 Gb	525	160 kb	64 Mb	89%

<sup>1</sup> ONT and PacBio assembly with CANU 1.5 + PILON  
<sup>2</sup> ONT 1D Long Reads protocol, sizing from 15 Kb to 90 Kb.  
<sup>3</sup> ONT 1D Long Reads protocol, sizing from 15 Kb to 90 Kb + ONT 1D protocol, sharing at 8 Kb, size > 8 Kb  
<sup>4</sup> Completeness computed with BUSCOv2

With this more complex genome, it is difficult to conclude about the quality of the assembly. Regarding the PacBio data, we obtained a higher number of contigs, but also more complete genes. ONT data allow us to reduce the number of contigs despite of the integrity of genes. Moreover, most of the short ONT reads can be removed to improve the assembly metrics.

In case we have some Illumina data, ONT can be equivalent to PacBio concerning the *de novo* assembly. It strongly depends on the biological query and the genome complexity. Today, we have to combine ONT data with an other data set to exploit them, but that can change in a short term future.

\* Acknowledgments: Thanks to Frédéric Breton *et al* and Tan Joon Sheong *et al* to allow us to present their results

### PacBio RSII, 10X Chromium: a love story ?

These two approaches are long reads technologies but unlike PacBio RSII, 10X Genomics use synthetic reads. For assembly purposes, are PacBio RSII and 10X Chromium results similar or is there a benefit to use them together?

The analysis have been performed on fish and tomato genomes. For PacBio *de novo* assemblies, we assessed 5 long reads assembly softwares: Miniasm<sup>v0.2</sup>, Falcon<sup>v0.4.2</sup>, Canu<sup>v1.4</sup> and SmartDeNovo<sup>v1.0.0</sup>. We obtained the best results with Canu and SmartDeNovo, so only those will be presented here.

#### Genome assemblies metrics

Fish genome assembly	Cov (X)	#	L50 contig	N50 contig
10x Supernova <sup>1</sup>	104x	43476	3641	51157
PacBio SmartDeNovo	71x	701	55	4137750
PacBio Canu	71x	4062	126	1341974
10x + SmartDeNovo <sup>2</sup>		534	44	5470783

Tomato genome Assembly	Cov (X)	#	L50 contig	N50 contig
10x Supernova <sup>1</sup>	116x	24579	1995	104298
PacBio SmartDeNovo	81x	857	112	2062400
PacBio Canu	81x	508	47	4973822
10x + Canu <sup>2</sup>		284	19	13603907

<sup>1</sup> The assemblies are realized by Supernova (v.1.1.2) a 10xGenomics software.  
<sup>2</sup> To combine PacBio RS II and Chromium 10x assemblies, we used ARCS software.

#### Genome assembly completeness<sup>1</sup> metrics

Fish Assembly <sup>2</sup>	Complete	Fragmented	Missing	Total groups searched
10x Supernova	2131 (82,4%)	328 (12,7%)	127 (4,9%)	2586
PacBio SmartDeNovo	2277 (88,1%)	224 (8,7%)	85 (3,2%)	2586
ARCS SmartDeNovo	2278 (88,1%)	223 (8,6%)	85 (3,2%)	2586

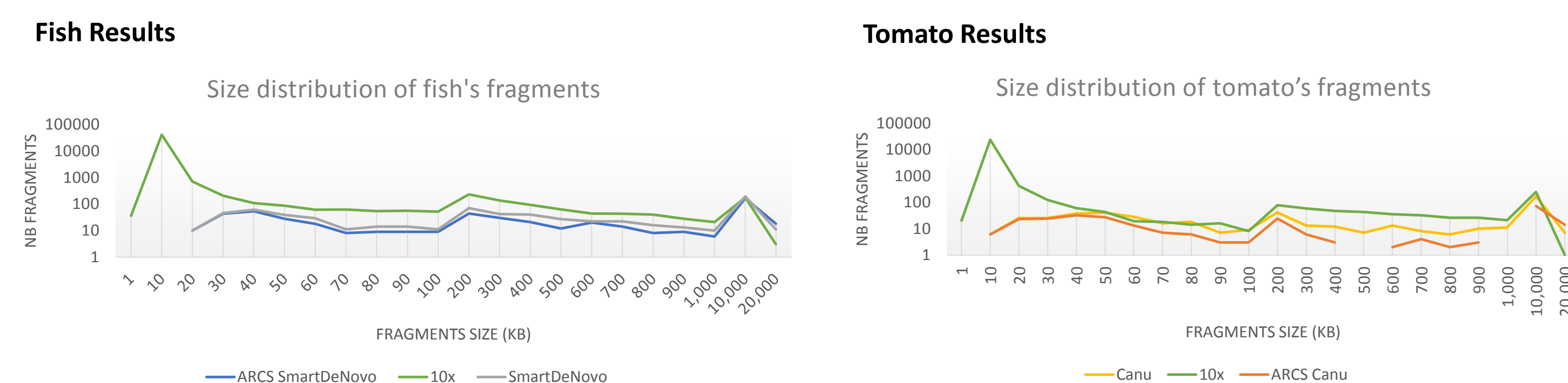
Tomato Assembly <sup>3</sup>	Complete	Fragmented	Missing	Total groups searched
10x Supernova	1303 (90,5%)	44 (3,1%)	93 (6,5%)	1440
PacBio Canu	1353 (94,0%)	25 (1,7%)	62 (4,3%)	1440
ARCS Canu	1357 (94,2%)	23 (1,6%)	60 (4,2%)	1440

<sup>1</sup> Completeness computed with BUSCO V2

<sup>2</sup> Lineage data: vertebrata\_odb9

<sup>3</sup> Lineage data: embryophyta\_odb9

#### Fragment distribution



In these two graphs, the fragment distribution of both samples are represented. The distributions follow the same trend except for 10X Chromium which are over represented at 10 kb. With ARCS software, we obtained less contigs but most of them are at 10.000 kb. Furthermore, there are more fragments at 20.000 kb in ARCS assemblies than the others. Indeed, the use of ARCS permits to obtain longer fragments than PacBio or 10X Chromium individual assemblies.

After the end of the analysis, we can not say there is a general « best assembler » for PacBio RSII because the result varies according to the samples. These differences can be explained by the genome complexity. The 10X Chromium Supernova results do not indicate better metrics than PacBio RSII. The combination of both improves significantly the PacBio quality assembly metrics.