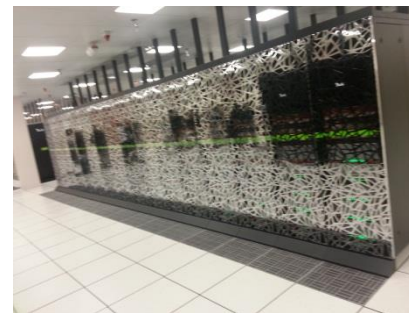


DE LA RECHERCHE À L'INDUSTRIE



[www.cea.fr](http://www.cea.fr)

## *RNA sequencing with the MinION at Genoscope*



**Jean-Marc Aury**



[jmaury@genoscope.cns.fr](mailto:jmaury@genoscope.cns.fr)



@J\_M\_Aury

December 13, 2017

RNA workshop, Genoscope

- Genoscope Overview
- MinION sequencing at Genoscope
- RNA-Seq using the Oxford Nanopore technology

<http://www.genoscope.cns.fr>

- French National Sequencing Center lead by Patrick Wincker, created in 1997 and part of the CEA since 2007
- Provide high-throughput sequencing data to the Academic community, and carry out in-house genomic projects
- Focus on biodiversity : *de novo* sequencing and metagenomic projects (TaraOceans)
- But.... it's not enough to just know one individual's DNA. A single reference genome is not compatible with resequencing approaches



*Triticum sp*  
(wheat)



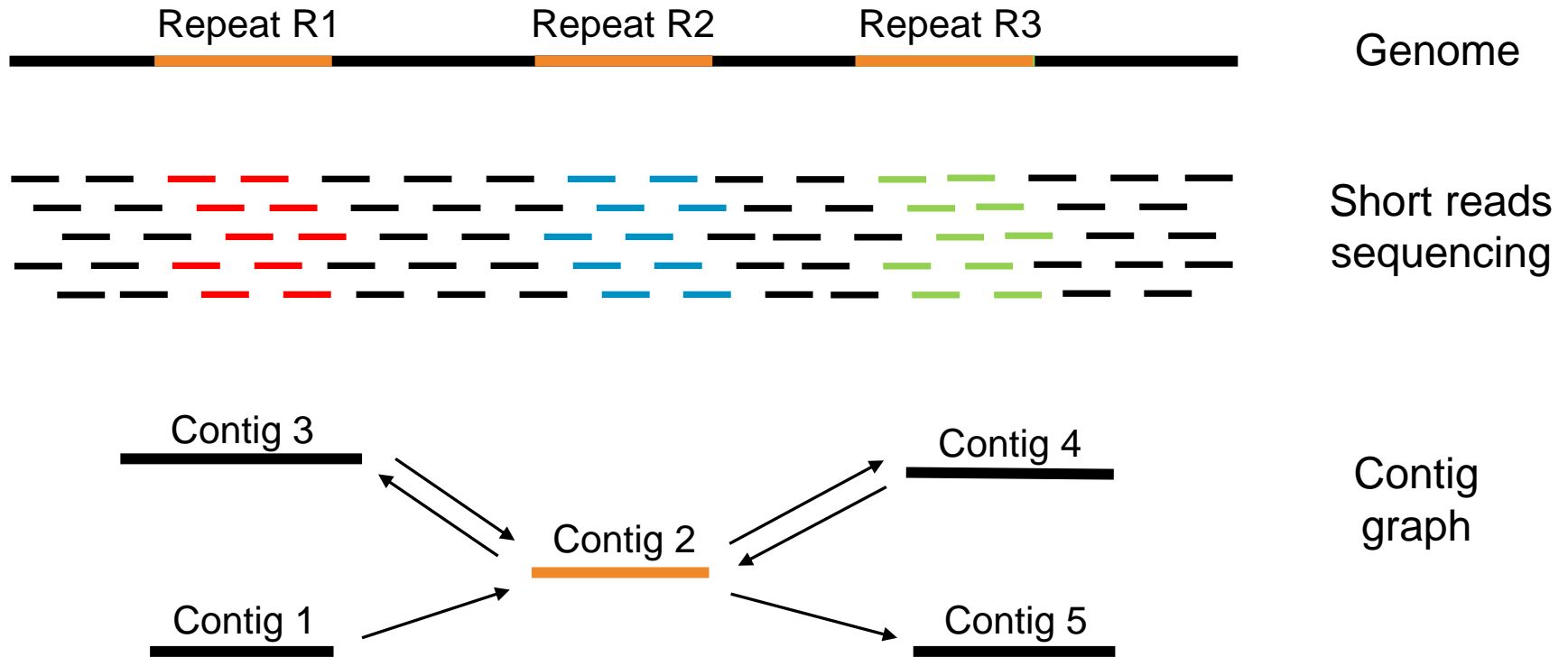
*Musa acuminata*  
(banana)



*Quercus robur*  
(oak)



*Brassica napus*  
(seed rape)



=> Repetitive regions lead to fragmented assemblies and under-estimate repeat content

# Sequencing capacities



2  
2

Illumina HiSeq 2500  
Illumina HiSeq 4000



2

MiSeq



6  
1

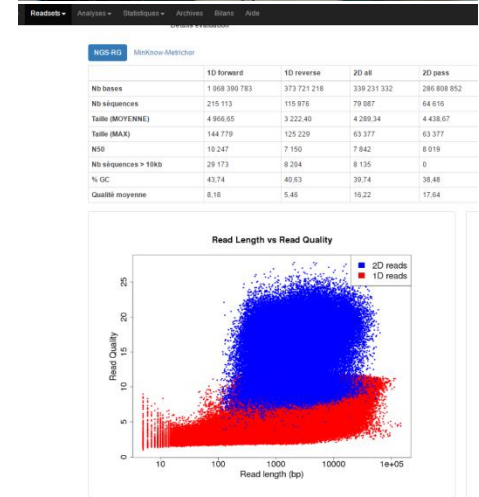
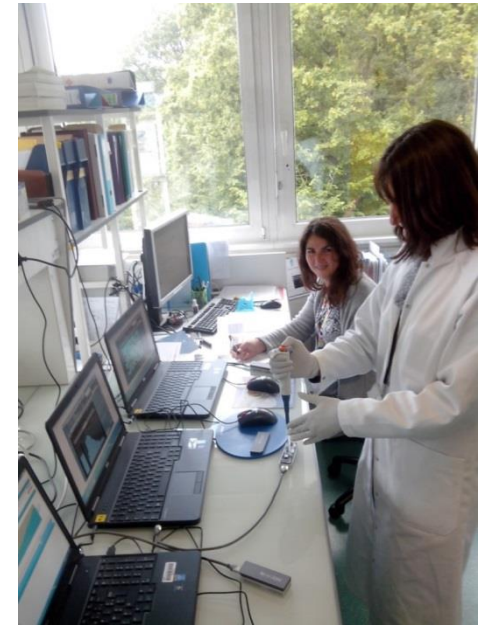
Oxford Nanopore Mk1  
PromethION



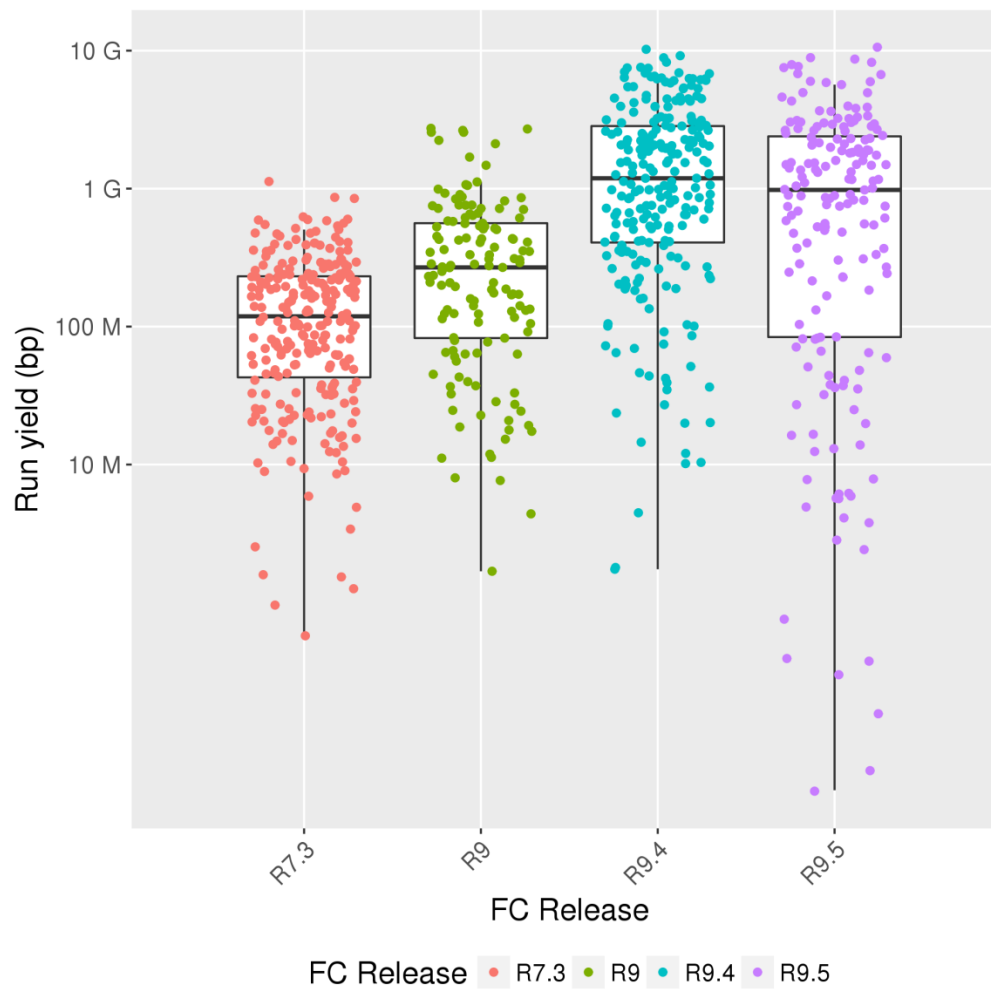
1

Irys System

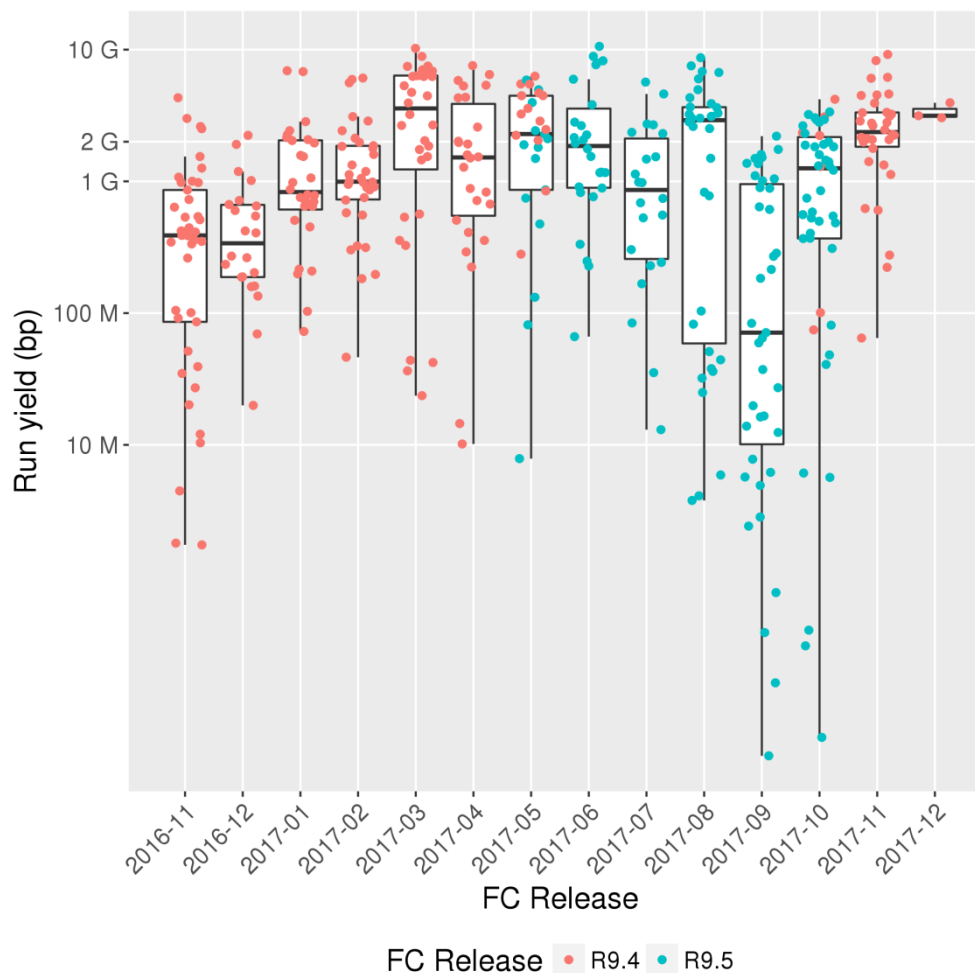
- 6 MinION devices
- >800 flowcells; >50 different organisms; ~700Gb of ONT reads ; DNA and RNA samples
- *de novo* assembly (22 yeast strains ~12Mb, 4 fungi genomes ~30Mb, several bacterial genomes, >10 plant genomes of 400-700Mb) and gene prediction
- Software development for the automation : management of the data flow, storing metrics in our LIMS
- Benchmark several DNA preparation protocols to obtain longer reads (size-selection using the blue pippin)



- Yield improvement :  $\sim 100\text{Mb}$  to  $>1\text{Gb}$  but the throughput of R9.5 flowcells seems to be more erratic

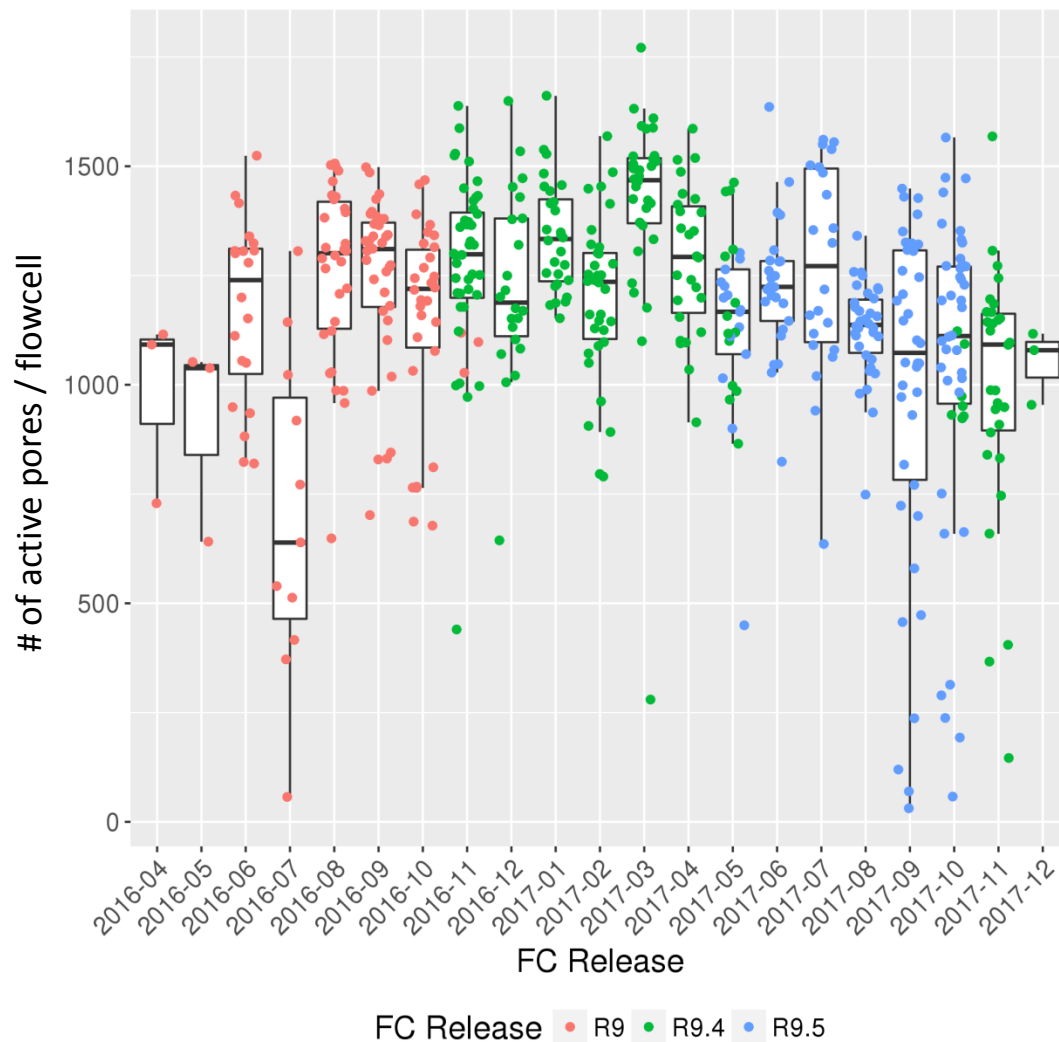


- The throughput dropped off in the last months, we now used R9.4 in production

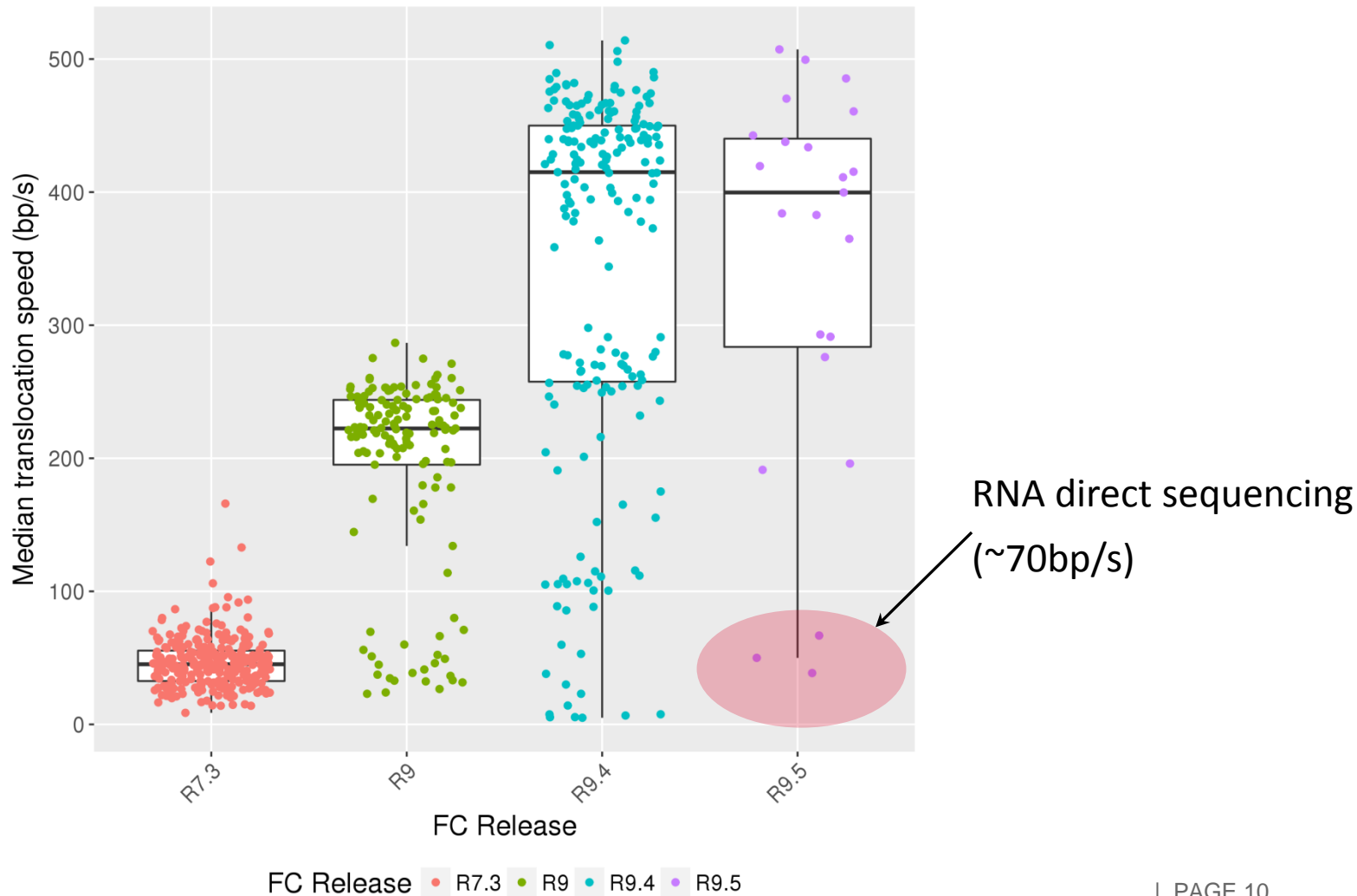




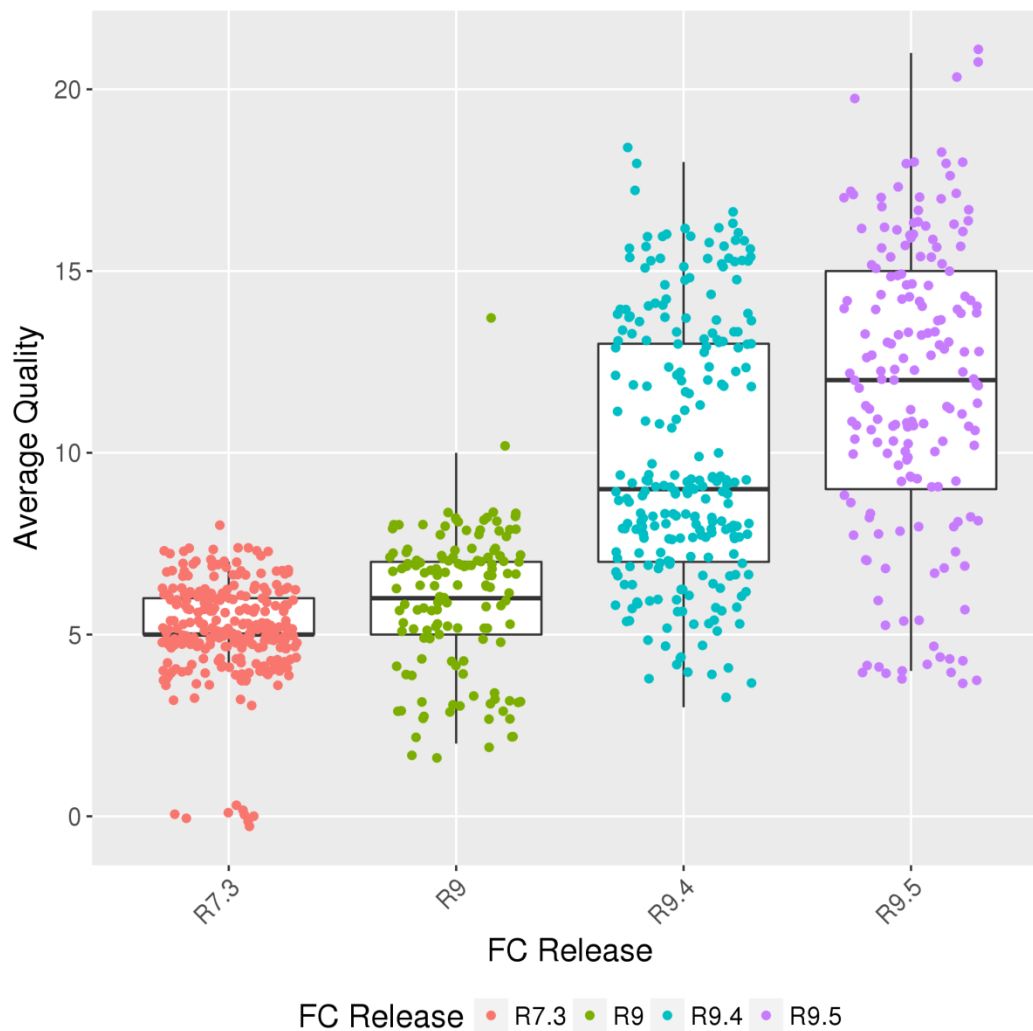
- The flowcell quality seems to be one of the issue



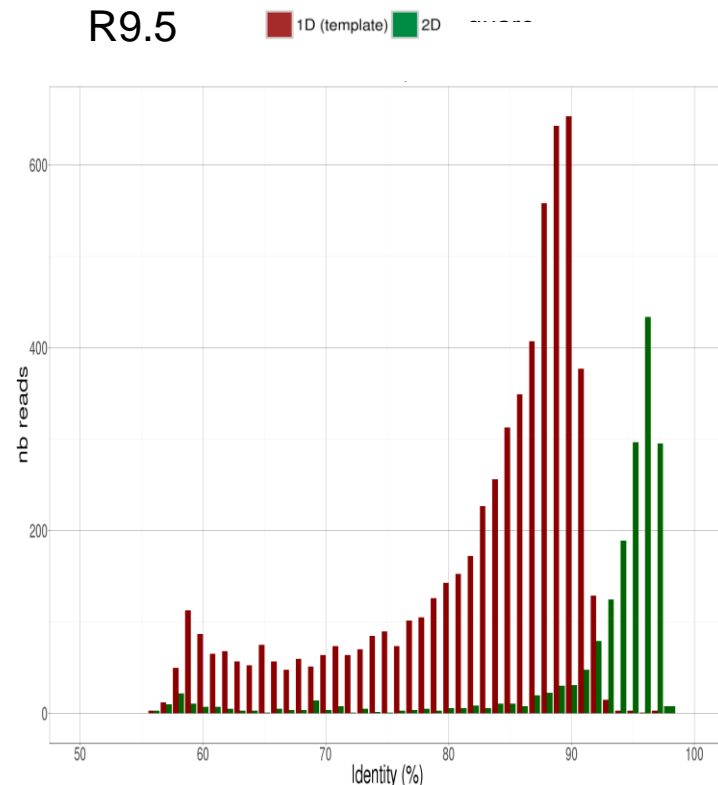
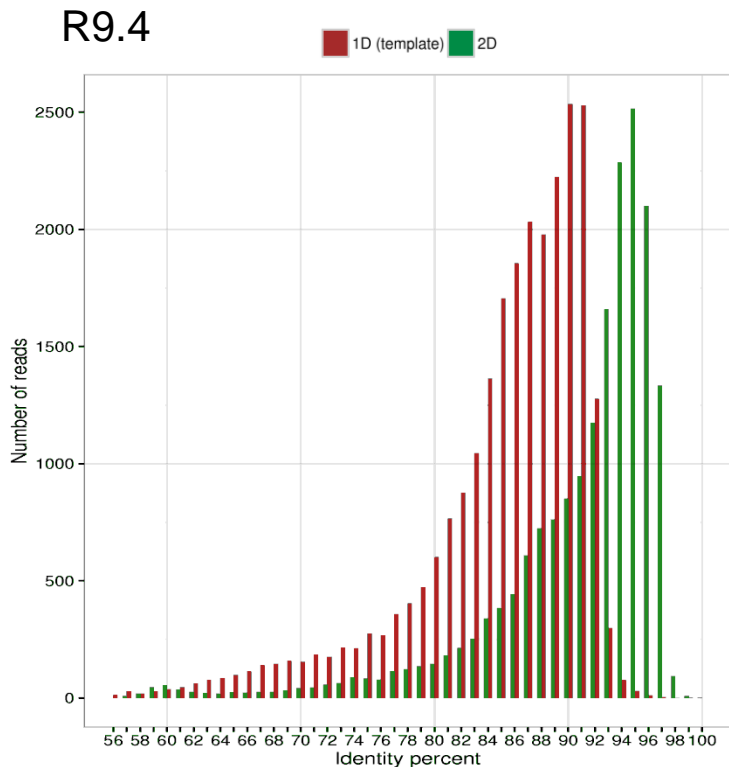
- Improvement of the DNA translocation speed through the pore



- Average quality and error rate improvement



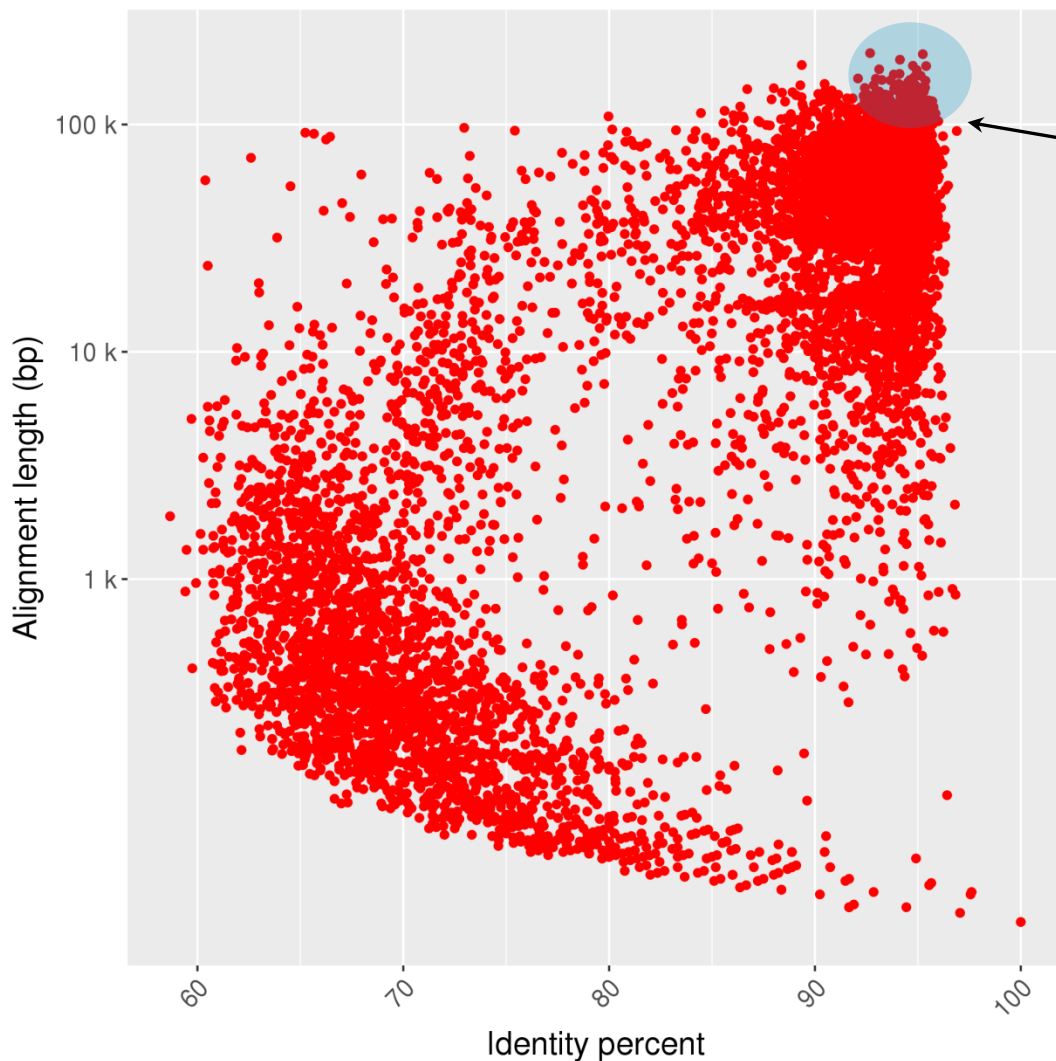
- Today error rate is even lower (in average 14% for 1D reads and <7% for 1D<sup>2</sup> reads),  
=> basecaller is a key component in the error rate drop off



Distribution of identity percent based on yeast reference (S288C). Alignments were performed using bwa-mem

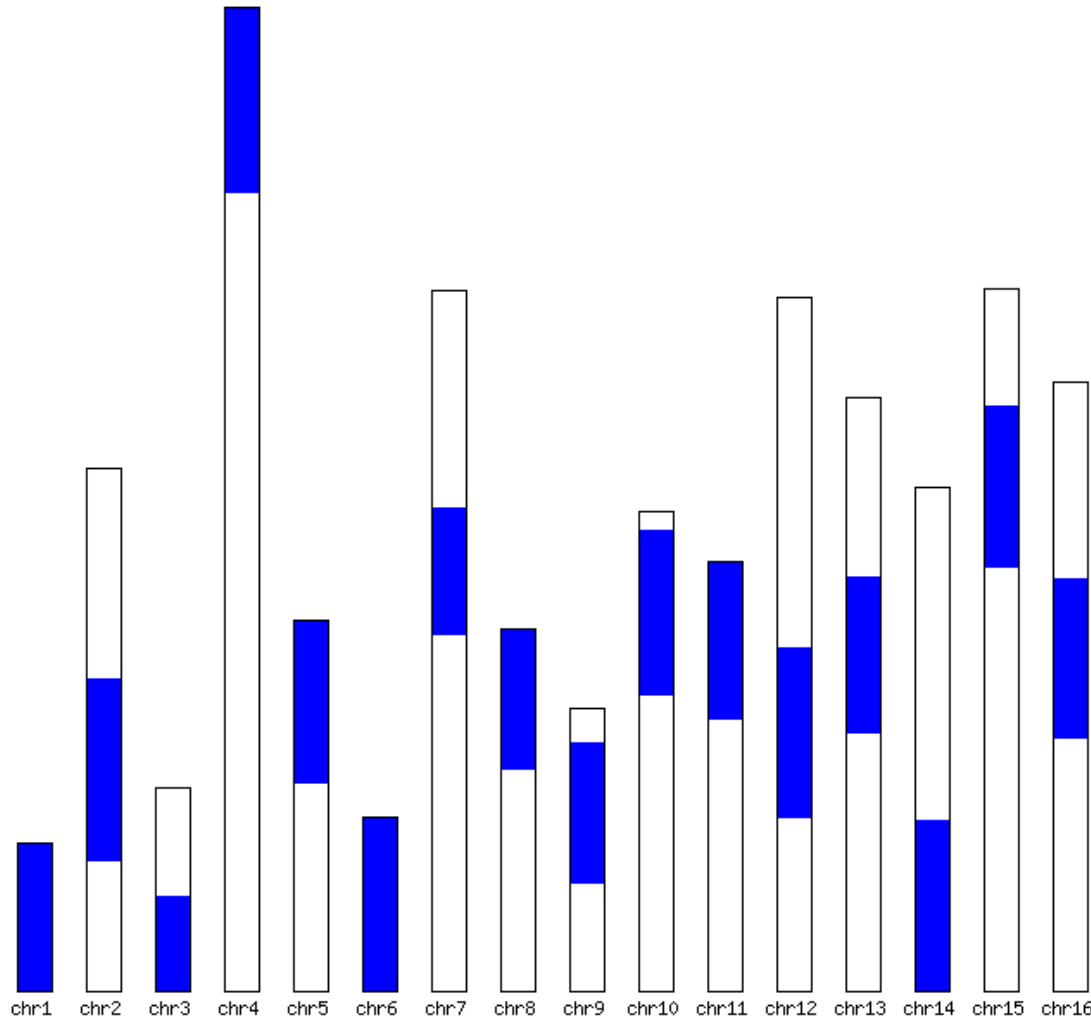
1D<sup>2</sup> is a real improvement in the error rate, unfortunately we get only up to 30% of 1D<sup>2</sup> reads

- The device is able to sequence very long DNA fragments (>100Kb)



~400 high quality reads  
with alignment length  
> 100Kbp  
=> ~4X of yeast genome

Nb bases	2 036 675 349
Nb sequences	137 109
Max length (bp)	461 529
N50 (bp)	50 800
Nb seq. > 50kb	11 695
Nb seq. > 100kb	3 275

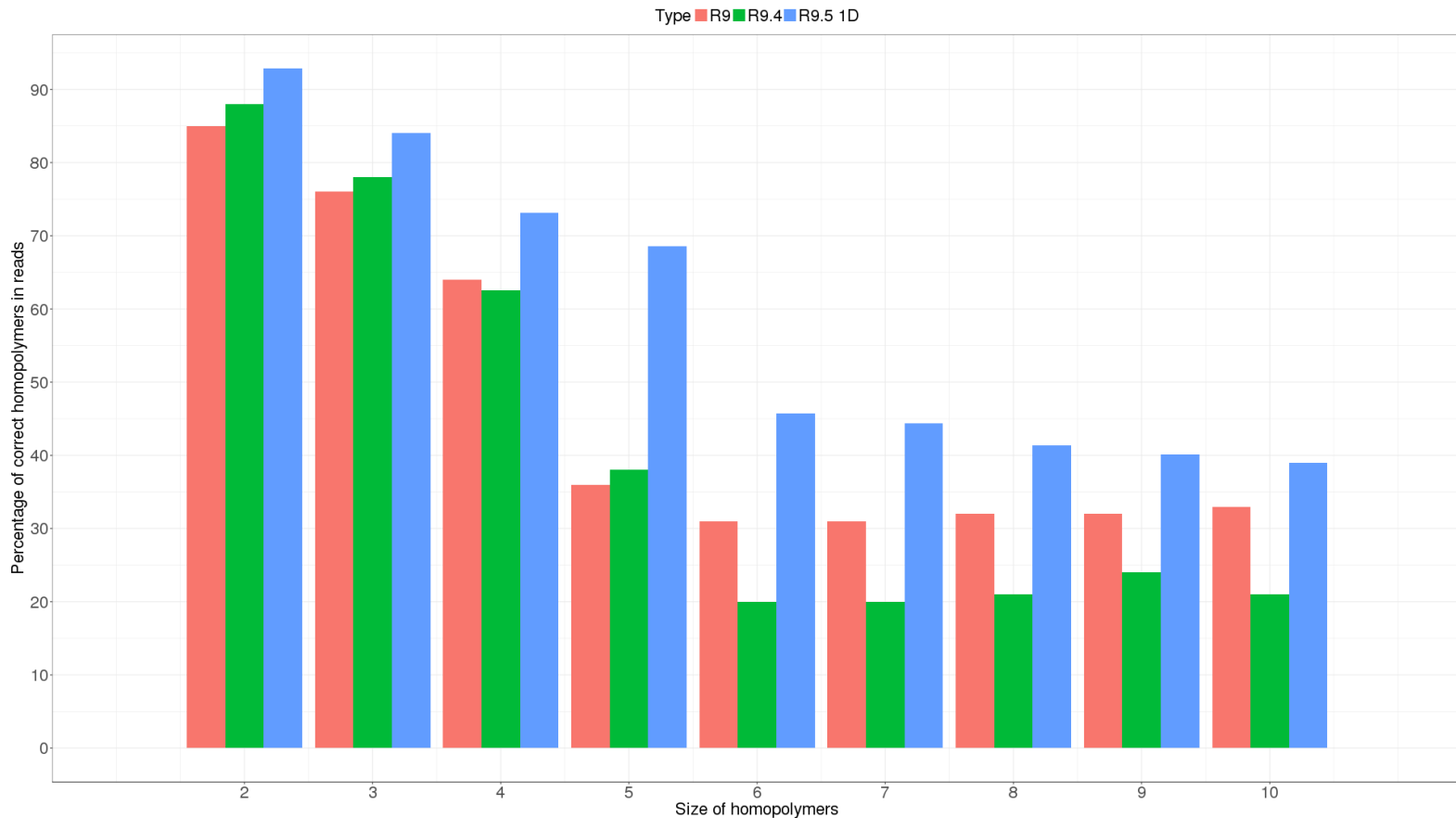


read with the longest alignment for each chromosome

Smallest chromosomes 1 and 6 are obtained in a single nanopore read !



- High error rate in homopolymers is still an issue for de novo sequencing projects, however the R9.5 release (and scrappie) really improve the basecalling of homopolymers.



It is still impossible to generate high quality consensus using nanopore only strategy.

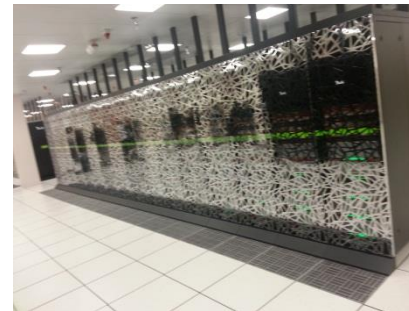


DE LA RECHERCHE À L'INDUSTRIE



[www.cea.fr](http://www.cea.fr)

# *cDNA-Seq and RNA-Seq using the Oxford Nanopore technology*



**Jean-Marc Aury**



[jmaury@genoscope.cns.fr](mailto:jmaury@genoscope.cns.fr)

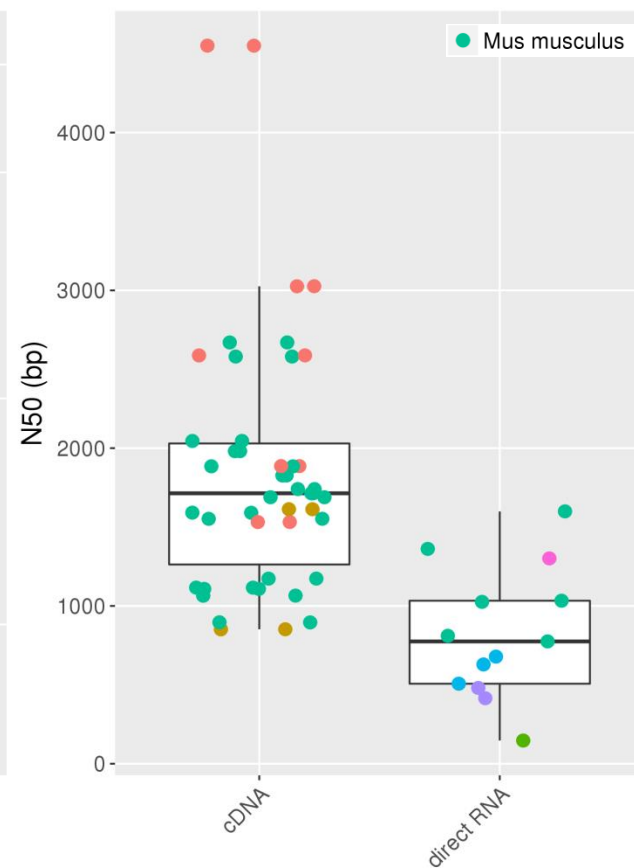
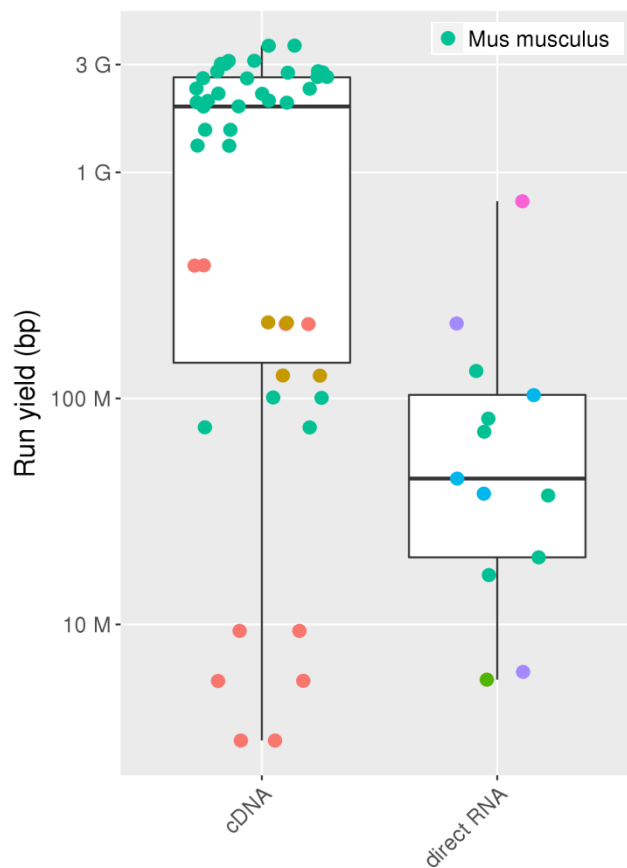
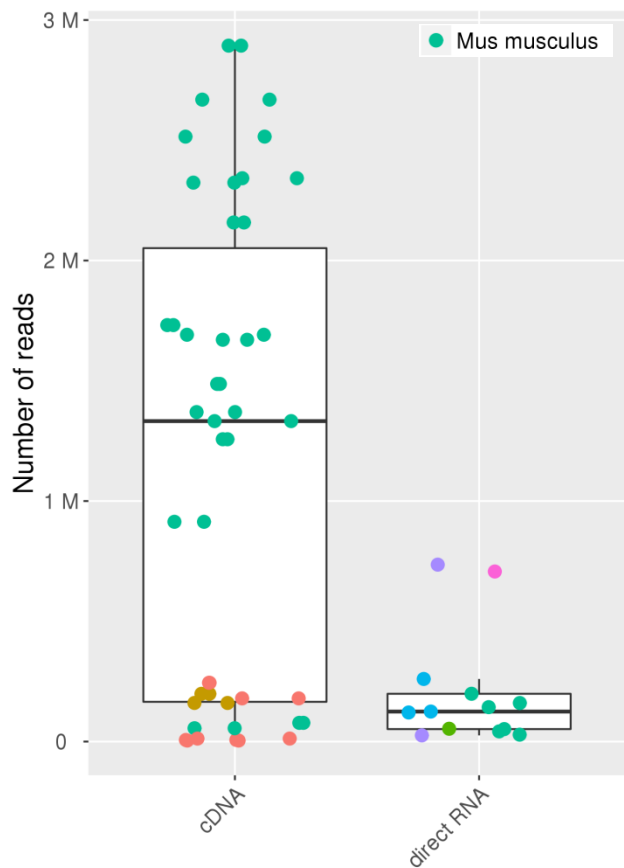


@J\_M\_Aury

December 13, 2017

RNA workshop, Genoscope

A typical cDNA-Seq experiment generates around 2M of reads, in comparison RNA-Seq experiments generate less reads (450bp/s vs 70bp/s)



## Dataset used to perform comparisons



## Brain sample

FC release	R9.4
Nb sequences	1 256 967
Nb bases	2 074 348 139
N50 (bp)	1 885

## Liver sample

FC release	R9.4
Nb sequences	1 369 927
Nb bases	1 956 452 499
N50 (bp)	1 591

cDNA sequencing

## Brain sample

FC release	R9.5
Nb sequences	160 450
Nb bases	81 508 561
N50 (bp)	1 033

## Liver sample

FC release	R9.5
Nb sequences	198 708
Nb bases	131 963 731
N50 (bp)	1 026

Direct RNA

## Brain sample

FC release	HiSeq 4000
Nb sequences	59M
Nb bases	17Gb
N50 (bp)	150

## Liver sample

FC release	HiSeq 4000
Nb sequences	45M
Nb bases	13Gb
N50 (bp)	150

Mapping of reads against RefSeq genes (refseq109) and the mouse genome (GRCm38)

Alignment against GRCm38 using minimap2 (36 cores)

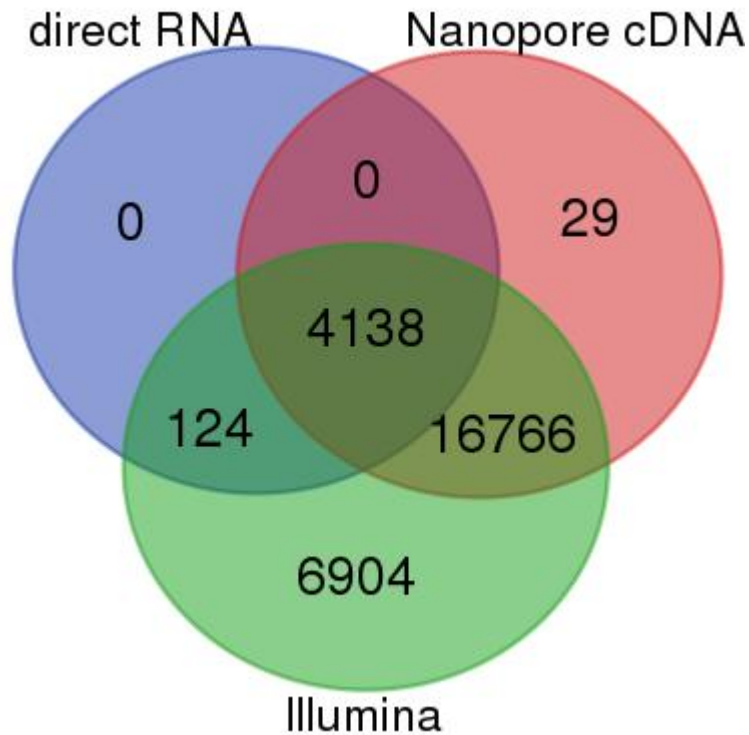
	Number of reads	Mapped reads	Mapped bases (of aligned reads)	Elapsed time (sec)
1D cDNA	1 256 967	90.7%	89.6%	396
RNA direct	160 450	33.8%	82.8%	20

Alignment against RefSeq 105 using bwa-mem (8 cores)

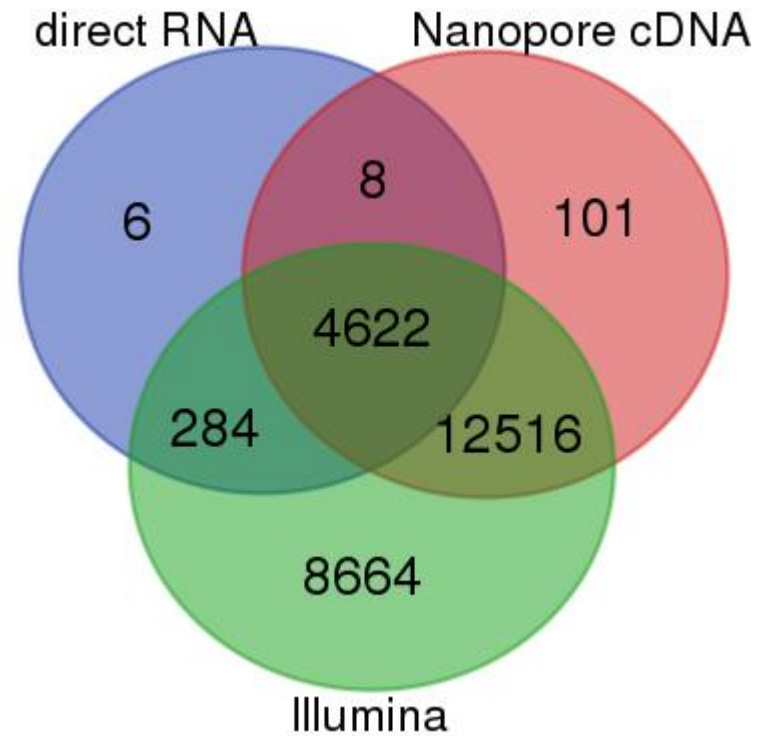
	Number of reads	Mapped reads	Mapped bases (of aligned reads)	Elapsed time (sec)	rRNA	Mitochondrial
1D cDNA	1 256 967	84.7%	64.2%	4 481	21.6%	15.8%
RNA direct	160 450	25.9%	75.2%	65	0.1%	18.5%

Number of RefSeq genes seen by each sequencing technology

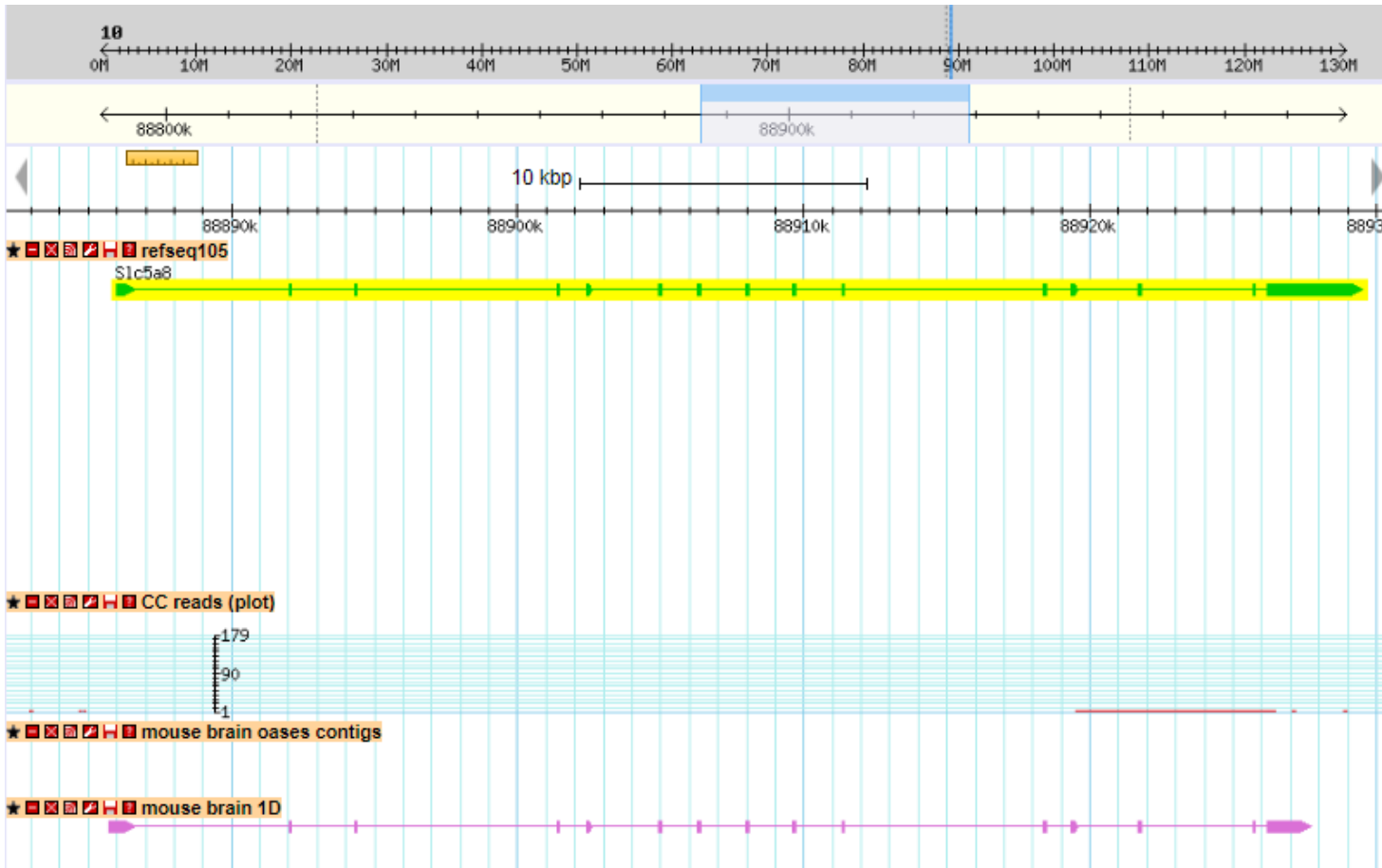
Brain sample



Liver sample



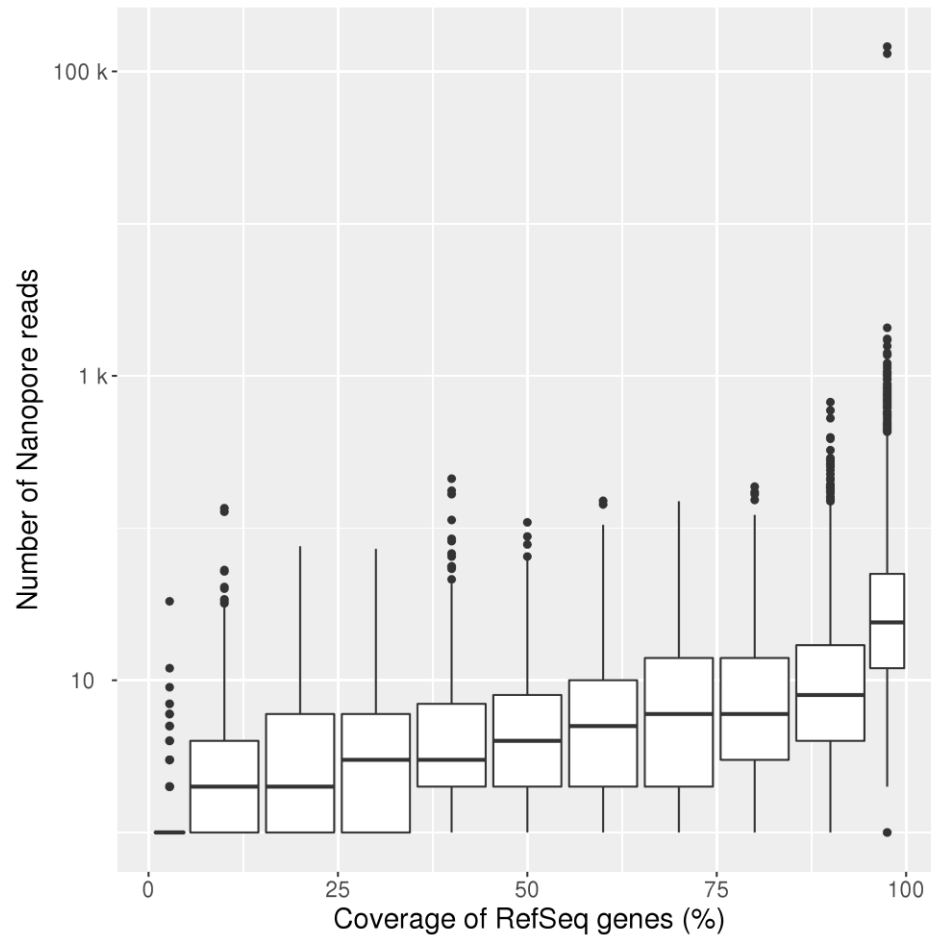
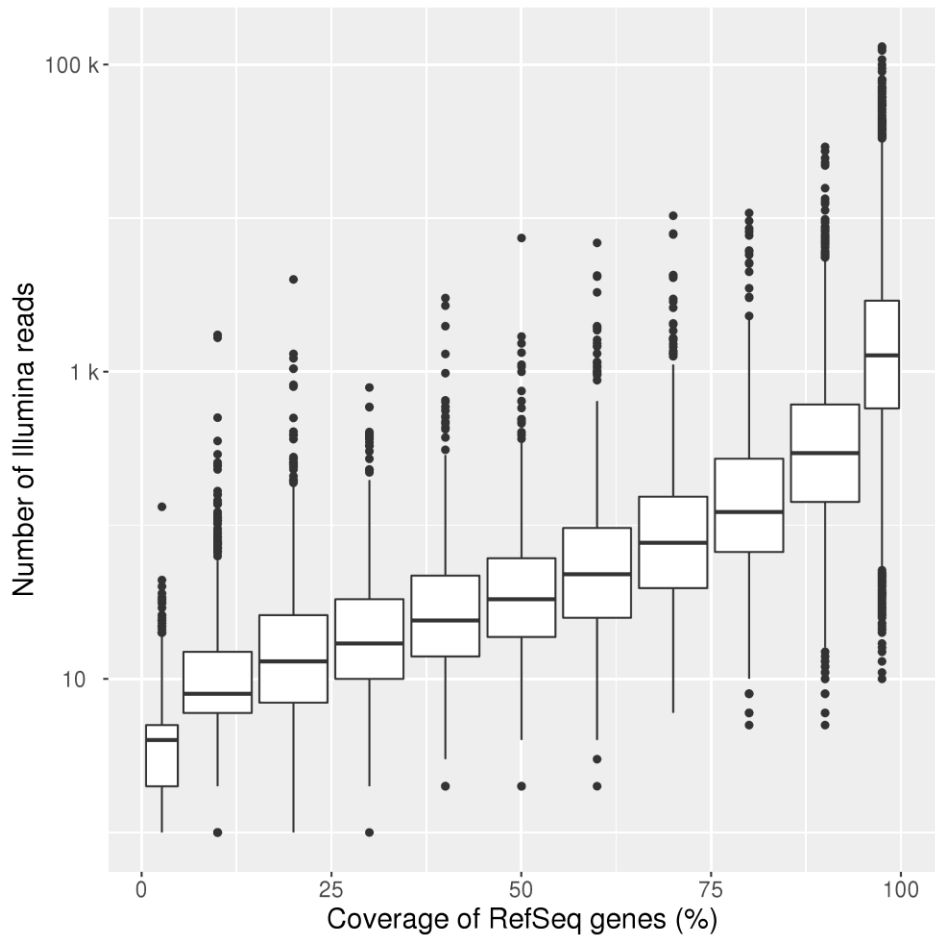
A gene can be covered entirely by a single read



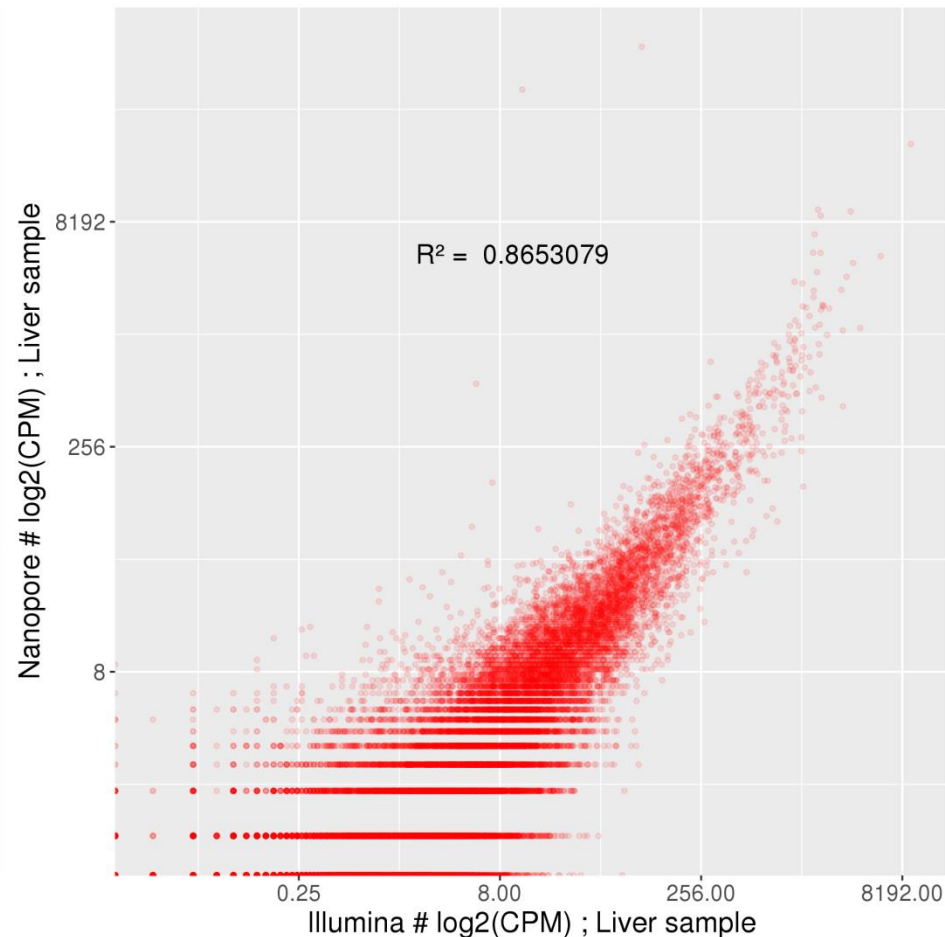
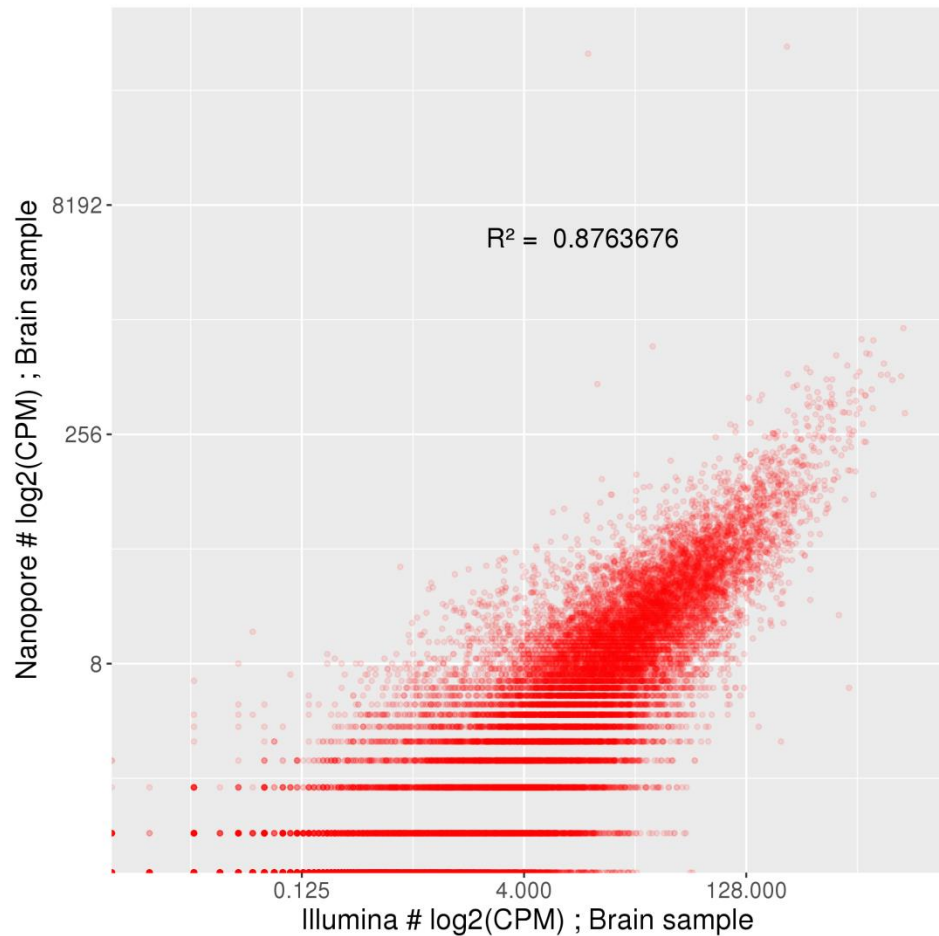
50 illumina reads are aligned and partially cover the gene

The entire gene is covered by a single nanopore read

As expected, less nanopore reads are needed to cover RefSeq genes, when we need at least 500 illumina reads to cover 75% of a given gene, 10 nanopore reads are sufficient



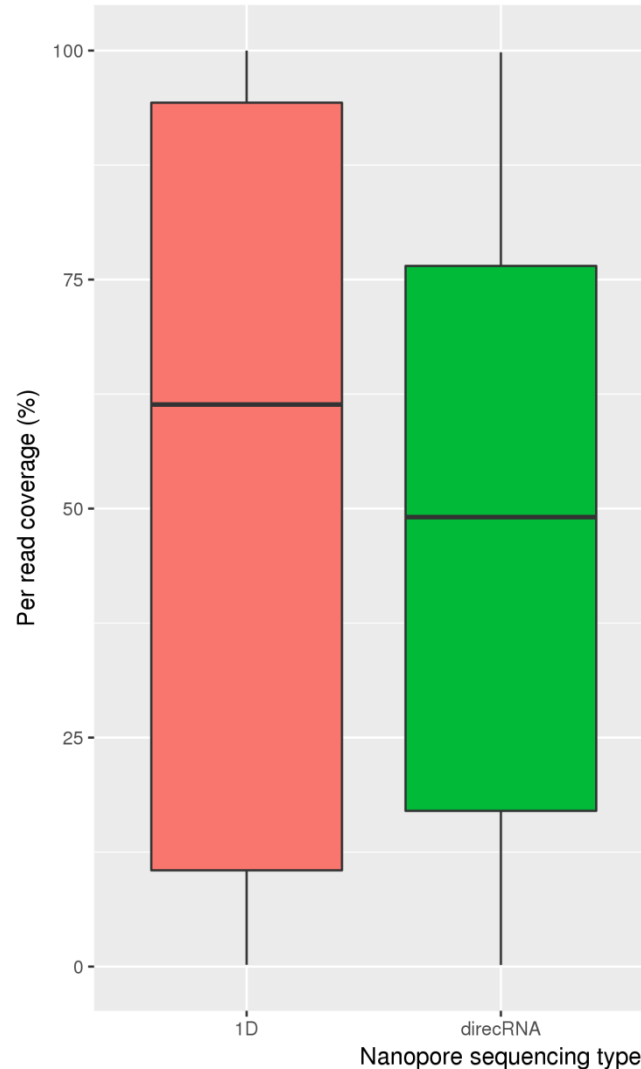
Expression levels (brain and liver samples) are correlated between Illumina and Nanopore experiments





# Are all reads full-length RNA ?

A small proportion of reads are full-length RNA, in average a cDNA and RNA read cover 55% and 47% respectively of a RefSeq gene



We tested the TeloPrime amplification kit from Lexogen



Based on Lexogen's unique Cap-Dependent Linker Ligation (CDLL) and long reverse transcription (long RT) technology, it is highly selective for full-length RNA molecules that are both capped and polyadenylated.

2 sequencing runs from brain and liver samples

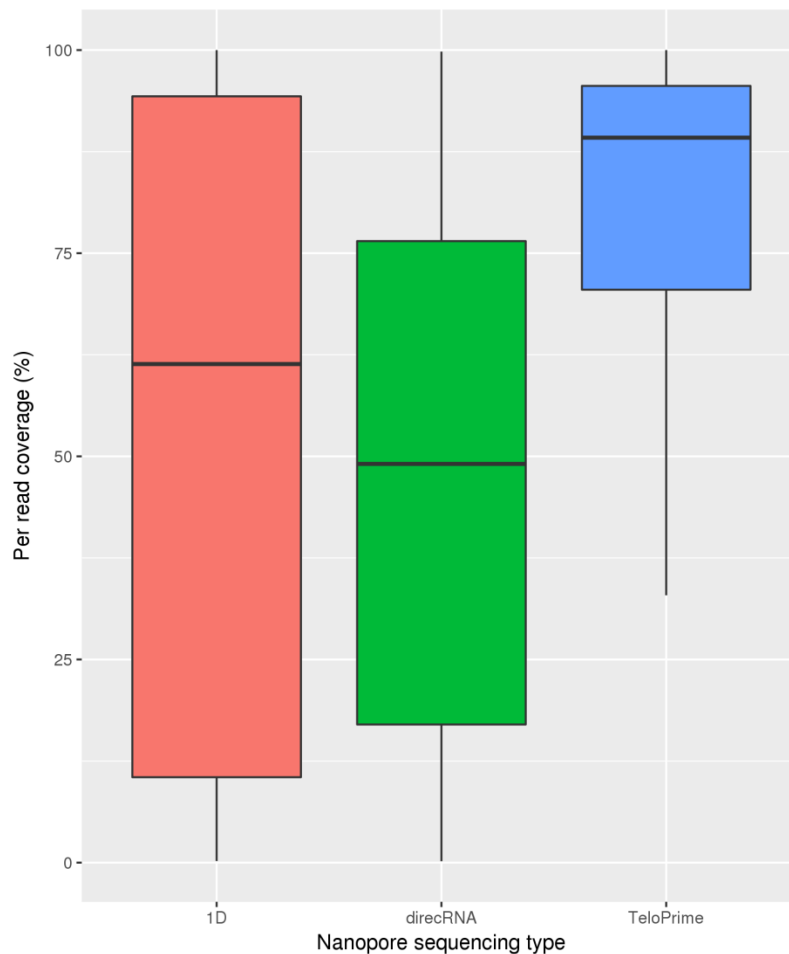
### Brain sample

FC release	R9.5
Nb sequences	2 668 975
Nb bases	2 641 896 941
N50 (bp)	1 116

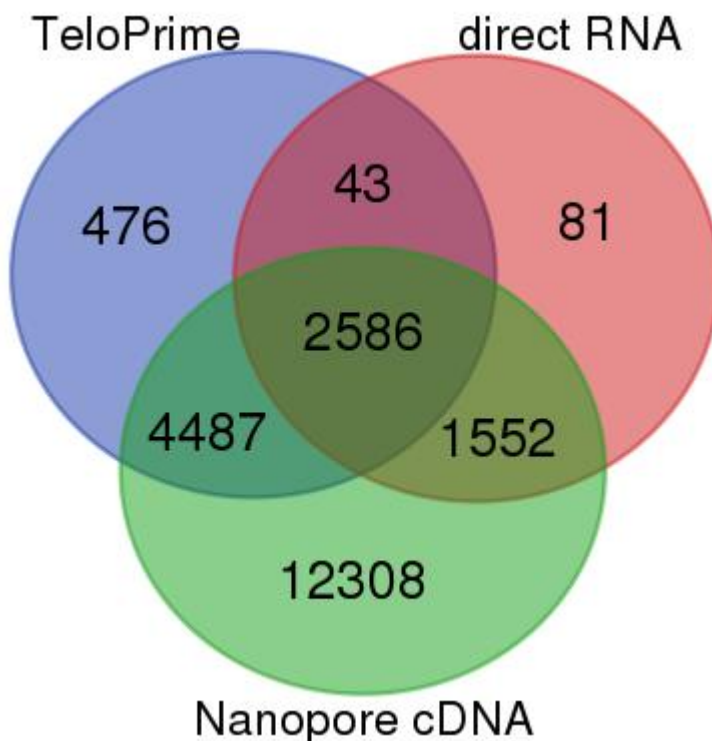
### Liver sample

FC release	R9.5
Nb sequences	1 691 454
Nb bases	1 312 184 503
N50 (bp)	896

TeloPrime reads better cover RefSeq genes, compared to cDNA and RNA sequencing. in average a TeloPrime read cover 80% of a RefSeq gene



Even with a higher number of reads, TeloPrime reads spread over a limited number of genes (~8k vs ~21k using 1D protocol)



### Nanopore cDNA

FC release	R9.4
Nb sequences	1 256 967
Nb bases	2 074 348 139
N50 (bp)	1 885

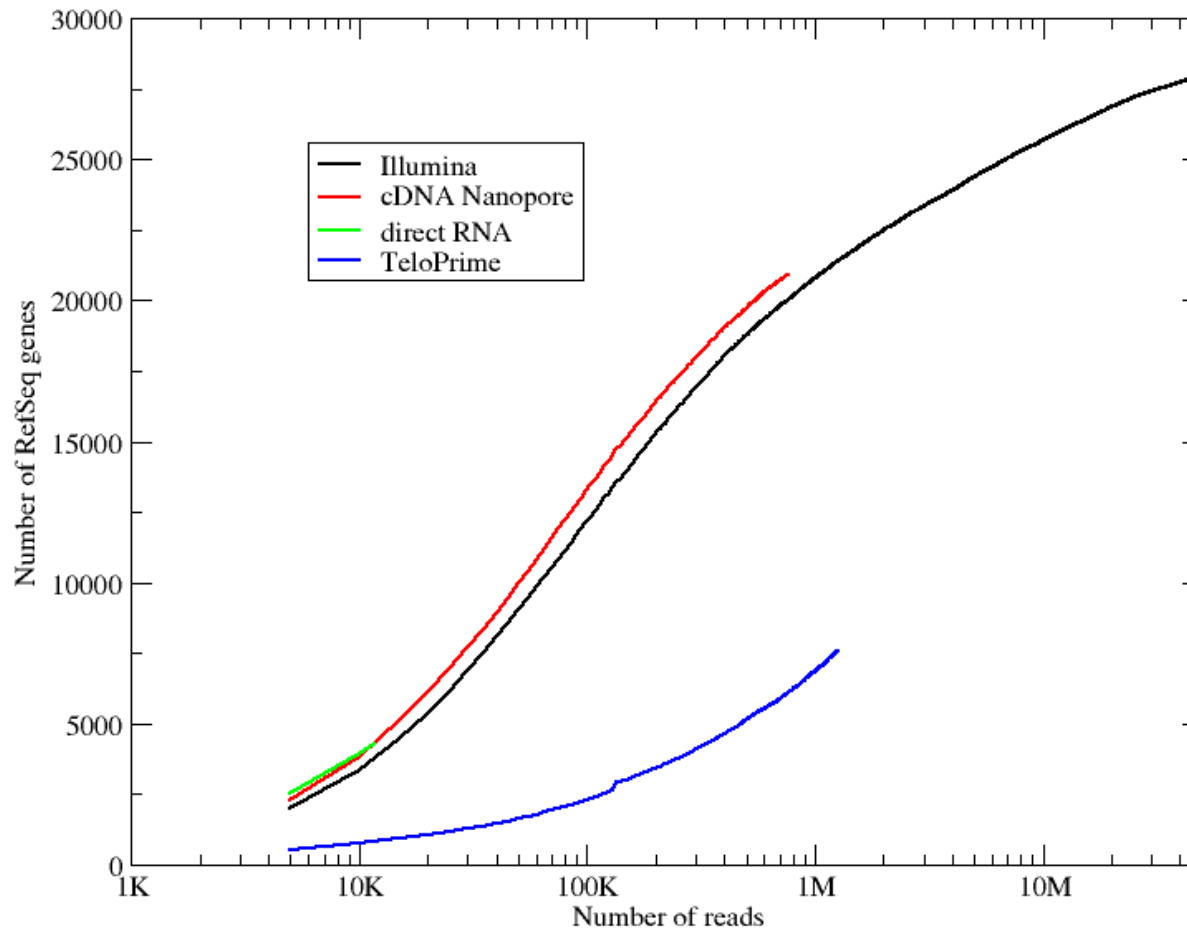
### Direct RNA

FC release	R9.5
Nb sequences	160 450
Nb bases	81 508 561
N50 (bp)	1 033

### TeloPrime

FC release	R9.5
Nb sequences	2 668 975
Nb bases	2 641 896 941
N50 (bp)	1 116

We need to sequence at a higher depth with the TeloPrime amplification kit to be able to catch a high proportion of RefSeq genes



- Today the throughput of the MinION device is sufficient for profiling eukaryotic gene expression, gene prediction can take advantage of long reads to avoid transcriptome assembly
- The potential of the device to sequence long reads is impressive, sequencing of large eukaryotic genomes is now possible even with the MinION device
- Error rate is acceptable for de novo sequencing projects (a high proportion of reads with less than 10% of errors), still an issue with homopolymers
- Need to improve the “wetlab part” to increase the proportion of full-length reads, TeloPrime kit seems to bring a real improvement



R&DBioSeq Team

[www.genoscope.cns.fr/rdbioseq](http://www.genoscope.cns.fr/rdbioseq)



[jmaury@genoscope.cns.fr](mailto:jmaury@genoscope.cns.fr)



@J\_M\_Aury

- Genoscope labs
  - Bioinformatic : Corinne Da Silva, Stefan Engelen, Benjamin Istace and Marion Dubarry
  - Nanopore Sequencing : Corinne Cruaud, Odette Beluche, Emilie Payen, Thomas Guérin and Arnaud Lemainque
- Members of the ASTER project
- Funding agencies : CEA, Genoscope, France Génomique and ANR



