# Toolkit to explore and analyze nanopore sequencing data on a Hadoop framework

Asunción Gallego<sup>1</sup>, Joaquin Tarraga<sup>1</sup>, Vicente Arnau<sup>2</sup>, Ignacio Medina<sup>3</sup> and Joaquin Dopazo<sup>1,4,5</sup> <sup>1</sup> Computational Genomics Department,Centro de Investigación Príncipe Felipe(CIPF) <sup>2</sup> Departamento de informática,ETSE,Universidad de Valencia <sup>3</sup> HPC Service, University Information Services, University of Cambridge <sup>4</sup> Bioinformatics of Rare Diseases(BIER),CIBER de Enfermedades Raras(CIBERER) <sup>5</sup> Functional Genomics Node,(INB) at CIPF

# INTRODUCTION

The use of nanopore technologies is expected to spread in the future because they are portable and can sequence long fragments of DNA molecules without prior amplification. The first nanopore sequencer available, the MinION<sup>™</sup> from Oxford Nanopore Technologies, is a USB-connected, portable device that allows real-time DNA analysis. In addition, other new instruments are expected to be released soon, which promise to outperform the current short-read technologies in terms of throughput. Despite the flood of data expected from this technology, the data analysis solutions currently available are only designed to manage small projects and are not scalable.

#### **RESULTS**

Here we present HPG Pore[1], a toolkit for exploring and analysing nanopore sequencing data. HPG Pore can run on both individual computers and in the Hadoop distributed computing framework, which allows easy scale-up to manage the large amounts of data expected to result from extensive use of nanopore technologies in the future.

## NANOPORE SEQUENCE ANALYSIS



Annex Finn Promeingallegotests/AnthLomanuato	1. col_H01495_3311_3_H100_H100_H100_H1045   Cour Test  Alt_Col_H01495_3311_3_H100_H1020_H1041 Hun5  Cour Test  Cour Test
Income Files     Accord Files     Accord Files     Accord Files     Accord Files     Accord Files     Accord Files	antualor, E. col., HIGL 655, 3331, 3, UNI 00, 514-20, 514-41 (Fauth Sector Sect
Lomanuder, L. coll, #01055, 2011, 3, cf     *	Control New Provide Annual Street and 20,000 disease Called Learning Inter- Canter Called





## HPG PORE'S COMMANDS

Command HPG Pore	Local	Hadoop	Description	
stats	yes	yes	Compute statistics and generates plots and histograms	
events	yes	yes	Extract raw data of the electronic signal measured for a given MinION read	
signal	yes	yes	Plot the electronic signal measured over time	
Fasta	yes	yes	Extract the sequences in format FASTA	
Fastq	yes	yes	Extract the sequences in format FASTQ	
import	no	yes	Store the individual FAST5 files into a single Hadoop MapFile in the Hadoop Distributed File System(HDFS)	
export	no	yes	Store the MinION reads to FAST5 local files from the Hadoop MapFile	
fast5names	no	yes	Extract the filenames of the imported FAST5 files	

#### **HPG PORE'S RUNTIMES**



Runtimes (left panel) and increase in speed (right panel) as the number of nodes increase in the Hadoop system in two different scenarios: FAST5 file containing 32,000 (blue line), 100,000 (red line), 300,000 (green line) and 1 million (dark blue line) sequences. Dotted line in the lower panel represents the ideal speed-up according to the number of nodes used. Speed-ups have been calculated using 3 nodes as the starting point given that the 1 million reads could not be calculated for 1 only one node.



## FEATURES, RUNTIMES AND SCALABILITY ON THREE NANOPORE TOOLS HPG PORE, PORE AND PORETOOLS

Feature	HPG Pore	PoRe	Poretools
Extract Fastq	$\checkmark$	✓	$\checkmark$
Extract Fasta	$\checkmark$	$\checkmark$	$\checkmark$
Organise fast5 into run folders	-	$\checkmark$	-
Create tar files of runs	-	-	✓
Organise the results into run folders	$\checkmark$	-	-
Plot yield	$\checkmark$	$\checkmark$	~
Plot signal	✓	$\checkmark$	✓
Extract run stats	~	$\checkmark$	~
Read length histogram	$\checkmark$	$\checkmark$	✓
Read length (max., avg., min)	$\checkmark$	$\checkmark$	~
Mean read quality	✓	-	-
Nucleotides content: count and %	✓	-	✓
%GC	$\checkmark$	-	-
Plot Frequency- %GC	$\checkmark$	-	-
Plot per base sequence content	✓	-	-
Read quality histogram	$\checkmark$	-	-
Reads per channel histogram	$\checkmark$	$\checkmark$	$\checkmark$
Nucleotides per channel histogram	$\checkmark$	$\checkmark$	-



Our study shows that runtimes in poRe, Poretools and HPG Pore (running locally) are approximately linearly dependent on the number of sequences in the FAST5 file, with a trend towards an increased slope for high numbers of sequences. HPG Pore runs the fastest, followed by Poretools, while poRe presents remarkably slower execution times.

# CONCLUSIONS

HPG Pore allows for virtually unlimited sequencing data scalability, thus guaranteeing its continued management in near future scenarios. HPG Pore is available in GitHub at http://github.com/opencb/hpg-pore.

#### REFERENCES

[1] Tarraga J, Gallego A, Arnau V, Medina I, Dopazo J. 2016 HPG pore: an efficient and scalable framework for nanopore sequencing data. BMC Bioinformatics. 17(1):107

