

# Initial analysis of 160,000 whole genomes from the Emirati Genome Program (EGP)

J. Quilez<sup>1</sup>, D. Sanchez<sup>1</sup>, C. Cremin<sup>1</sup>, H. Wu<sup>1</sup>, H. Sajad<sup>1</sup>, J. Mafofo<sup>1</sup>, G. Katagi<sup>1</sup>, H. Al Mabrazi<sup>1</sup>, V. Kusuma<sup>1</sup>, S. Pejathaya<sup>1</sup>, S. Elavalli<sup>1</sup>, F. Aldhuhoori<sup>1</sup>, I. Eltantawy<sup>1,2</sup>, M. Ibrahim<sup>1,2</sup>, A. Alsuwaidi<sup>1</sup>, K. Thangarajan<sup>1</sup>, A. M. Yousif<sup>3</sup>, M. Olbrich<sup>4</sup>, M. S. Al Ameri<sup>3</sup>, I. Chishty<sup>1</sup>, S. A. Al Marzooqi<sup>3</sup>, W. M. Abdulrahman<sup>3</sup>, A. I. Al Mannaei<sup>3</sup>, M. Mousa<sup>4</sup>, H. Al Safar<sup>4</sup>, I. Iqbal<sup>3</sup>, F. Al Marzooqi<sup>1</sup>, E-K. Albarah<sup>1</sup>, A. Al Awadhi<sup>1</sup>, T. R. Magalhaes<sup>1</sup>

## Introduction

- The EGP has used three sequencing technologies to perform whole-genome sequencing (WGS) of approximately 600,000 participants so far.
- DNA is extracted from either blood or buccal swab (LRS only on blood-derived DNA).
- For the primary and secondary analysis (i.e., from the sequencing raw data to the individual's genetic variants across the genome) we have implemented optimized sequencing platform specific analysis workflows (Fig. 1).
- A key question is: **are the EGP's population-scale genomics datasets of good quality?**

## Methods

- Selected 7 parameters and set cutoffs to define high-quality individual WGS datasets (Table 1).
- Evaluated those parameters in an initial data freeze comprising ~160,000 WGS (Table 2).

## Results

- WGS samples for over 150,000 EGP participants (~95%) passed the high-quality criteria.
- All three sequencing platforms performed similarly well.
- As expected, blood performed better than buccal but the latter still was very good.
- Amongst failed samples, most were buccal and did so because of mapping rate only slightly below the 70% cutoff

**Table 2.** QC summary. Illumina and MGI have lower pass rate because they include buccal swabs while Oxford Nanopore Technologies (ONT) does not.

	Total Samples	Samples passing QC	Pass rate (%)
Illumina	107,739	100,530	93.3
MGI	44,305	43,115	97.3
ONT	7,755	7,602	98.0
<b>All</b>	<b>159,799</b>	<b>151,247</b>	<b>94.6</b>

The Emirati Genome Program (EGP) aims to **sequence the entire Emirati population** (~1 million people) by 2025.

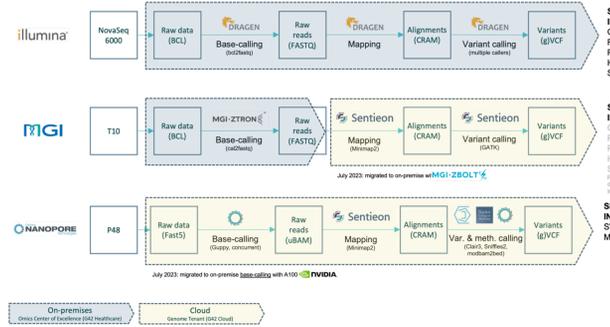
So far, over **600,000 genomes** have been **sequenced** using short- and long-read sequencing.

Implemented **optimized technology-specific analysis workflows** for all three.

**High quality of the genomics datasets**, with about 98% of blood and 90% buccal swab samples passing high-quality criteria.



**Figure 1.** Analysis workflows across sequencing platforms.



**Table 1.** 7 QC metrics and associated cutoffs defining high-quality WGS samples.

QC metric	High-quality cutoff
Yield (Gb)	≥70
% bases >Q30 (>Q10 for ONT)	≥85
Mapping rate (%)	≥70
Mean coverage (X)	≥25X
% genome >10X	≥95
Het/hom ratio	<2.5

Based on UK Biobank, 1000 Genomes, GATK Best Practices and internal benchmarks.

## Acknowledgements

This work would have not been possible without the commitment of all the Emirati citizens who have voluntarily enrolled into the EGP and the support from the Abu Dhabi Department of Health.



jquilez@m42.ae

@m42\_health

