# DEVELOPMENT OF A COMPUTATIONAL PIPELINE FOR NEXT-GENERATION SEQUENCING DATA ANALYSES USING NEXTFLOW AND DOCKER

**OWOLABI, PAUL JESUSANMI**
**(21PBF02261)**
**B.Sc Microbiology, Obafemi Awolowo University, Ile-Ife, Nigeria**

**AUGUST, 2023**

# DEVELOPMENT OF A COMPUTATIONAL PIPELINE FOR NEXT-GENERATION SEQUENCING DATA ANALYSES USING NEXTFLOW AND DOCKER

**BY**

**OWOLABI, PAUL JESUSANMI**
**(21PBF02261)**
**B.Sc Microbiology, Obafemi Awolowo University, Ile-Ife, Nigeria**

**A DISSERTATION SUBMITTED TO THE SCHOOL OF POSTGRADUATE STUDIES, IN PARTIAL FULFILMENT OF THE REQUIREMENTS FOR THE AWARD OF A MASTER OF SCIENCE DEGREE IN BIOINFORMATICS IN THE DEPARTMENT OF COMPUTER AND INFORMATION SCIENCES, COLLEGE OF SCIENCE AND TECHNOLOGY, COVENANT UNIVERSITY, OTA, OGUN STATE, NIGERIA**

**AUGUST, 2023**

# ACCEPTANCE

This is to attest that this dissertation has been accepted in partial fulfilment of the requirements for the award of the degree of Master of Science in Bioinformatics in the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, Nigeria.

**Miss Adefunke F. Oyinloye**
**(Secretary, School of Postgraduate Studies)**                    **Signature and Date**

**Prof. Akan B. Williams**
**(Dean, School of Postgraduate Studies)**                    **Signature and Date**

# DECLARATION

I, **OWOLABI, PAUL JESUSANMI (21PBF02261)** hereby declare that this dissertation titled **"DEVELOPMENT OF A COMPUTATIONAL PIPELINE FOR NEXT-GENERATION SEQUENCING DATA ANALYSES USING NEXTFLOW AND DOCKER"** is a representation of my work and is written and implemented by me under the supervision of Dr. Itunuoluwa M. Isewon of the Department of Computer and Information Sciences, Covenant University, Ota, Nigeria. I attest that this dissertation has in no way been submitted either wholly or partially to any other university or institution of higher learning for the award of a masters' degree. All information cited from published and unpublished literature has been duly referenced.

**OWOLABI, PAUL JESUSANMI**

                                                **Signature and Date**

# CERTIFICATION

This is to certify that this dissertation titled **"DEVELOPMENT OF A COMPUTATIONAL PIPELINE FOR NEXT-GENERATION SEQUENCING DATA ANALYSES USING NEXTFLOW AND DOCKER"**, is an original research carried out by **OWOLABI, PAUL JESUSANMI (21PBF02261)** and meets the requirements and regulations governing the award of Master of Science (M.Sc.) degree in Bioinformatics from the Department of Computer and Information Sciences, College of Science and Technology, Covenant University, Ota, and is approved for its contribution to knowledge and literary presentation.


**Dr. Itunuoluwa M. Isewon**
**(Supervisor)**                                          **Signature and Date**



**Prof. Olufunke O. Oladipupo**
**(Head of Department)**                                  **Signature and Date**



**Prof. Adebukola S. Onashoga**
**(External Examiner)**                                   **Signature and Date**



**Prof. Akan B. Williams**
**(Dean, School of Postgraduate Studies)**                **Signature and Date**

# DEDICATION

I am dedicating this work to the Almighty God, who is my true source, and to those committed to the furtherance of life science research.

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

ADME    Absorption, Distribution, Metabolism, and Excretion

BWA    Burrows-Wheeler Aligner

CNVs    Copy Number Variations

CWL    Common Workflow Language

DNA    Deoxyribonucleic Acid

GATK    Genome Analysis Toolkit

GRCh38    Genome Reference Consortium Human Build 38

GWAS    Genome Wide Association Studies

HTS    High-throughput Sequencing

NCBI    National Centre for Biotechnology Information

NGS    Next generation Sequencing

SNPs    Single Nucleotide Polymorphisms

SNVs    Single Nucleotide Variants

SRA    Sequence Read Archive

SVs    Structural Variants

WES    Whole Exome Sequencing

WGS    Whole Genome Sequencing

# ABSTRACT

Major advances in genomics studies, particularly the introduction of high-throughput sequencing and the evolution of genotyping platforms have led to the emergence of big data in the biological sciences and a growing need to make sense of this data. This has largely fostered the evolution of methods and tools for genomic data analysis (especially of diseased conditions) with the aim of uncovering the genotype-phenotype relationships in such diseased conditions. Due to the growing complexity and volume of next-generation sequencing data available in biological sciences, there is a growing need to developed pipelines that can handle these data while automating most of the steps involved in these analyses. The aim of this study is to develop a computational pipeline for the analysis of next-generation sequencing data using Nextflow and Docker. Since different steps and tools are involved in the analysis of whole genome and whole exome sequencing data, the aim of the study was achieved by developing scripts for selected genome analysis tools, building a computational pipeline for the selected tools and performing unit and integration testing for the pipeline. The pipeline which was built on the framework of the well-established GATK best-practices workflow, integrated the following tools: FastQC, MultiQC, Jellyfish, genomeScope2.0, BWA, GATK and SnpEff. These tools were involved in performing the different steps of the NGS analyses which included quality control check, genome size heterozygosity, alignment or mapping, variant calling and annotation. Nextflow was employed in this pipeline as a workflow management system and Docker was used for containerising all the tools and their software dependencies. The developed pipeline was then tested to verify its utility in NGS data analysis. Pipeline development is very important in genomics research because, it could help improve the quality and reliability of research outcomes and facilitate the sharing and comparison of data across different studies and research groups. Having a pipeline that can effectively be used in quick and simple analysis of genomes will significantly help in uncovering biologically meaningful or clinically significant variants. It is expected that the outcome of this study will significantly impact studies into the genetic basis of human diseases and precision medicine.

**Keywords:** *Next-generation Sequencing, Genomic analysis, Variant calling, Nextflow, Docker, Pipeline*