User Manual for Nextflow Pipeline Tool

Introduction

Antimicrobial resistance (AMR) is one of the top global health threats, as stated by the World Health Organization (WHO). It leads to resistance against antibiotics, making infections harder to treat. Pseudomonas aeruginosa is a major concern in AMR due to its resistance, which makes treating infections such as bloodstream infections and respiratory system infections more challenging.

To address this issue, there is a need for a variant database that aids in clinical studies, and therapeutic research. This tool is designed to create a variant database using FASTA, FASTQ, and metadata CSV files. It generates VCF, CSV, and database (DB) files as output and provides a web interface for analysis.

Workflow

The tool follows these key processes:

- Indexing
- Trimming and quality check
- Alignment and sorting
- Duplicate marking
- Variant calling and annotation
- Variant filtering
- Database creation using SQLite3

Database Structure

The generated database contains:

1. Variant Table with information such as:

0	Chromosome name
0	Strain name (of the reference)
0	Genome ID
0	Amino acid change
0	Reference allele
0	Alternative allele
0	Gene name etc.,
2. Strain	Table with fields including:
0	Genome ID
0	SRA accession
0	Geographical group
0	Isolation country
0	Genome name (strain name)
0	Taxon ID
0	Resistant phenotype
This tool enables the creation of an instance variant database for all bacterial species and users	
must provide the correct input files corresponding to the correct species to ensure accurate	

Step-by-Step Guide

results.

- 1. Install Docker.
- 2. Use Docker commands to build and run the image.
- 3. First, build the database for SnpEff.
- 4. Use mkdir to create a new folder, making it easier to run and access system data
- 5. Store all files in a single directory, based on the input parameters mentioned below.
- 6. Use the Nextflow command with the required parameters to process the data.

1. System Requirements:

Minimum Requirements:

• CPU: Dual-core processor

• RAM: 6 GB

• Storage: 8 GB free space

• Python: 3.8 or higher

• Nextflow: Latest version

• Tools: Latest version

• Docker: Latest version

• OS: Windows 10/11, Linux (Ubuntu 20.04+)

Recommended Requirements:

• CPU: Quad-core processor

• RAM: 14GB or more

• Storage: 8 GB free space

• Python: 3.10 or higher

• Nextflow: Latest version

• Tools: Latest version

• Docker: Latest version

• OS: Windows 10/11, Linux (Ubuntu 22.04+)

2. Installation:

A. Installing Nextflow Without Docker

Linux (Ubuntu/Debian-based)

- Install dependencies:
 - O sudo apt update && sudo apt install -y openjdk-21-jdk curl
- Install Nextflow:

- O curl -s https://get.nextflow.io | bash
- Move Nextflow to a directory in your system PATH:
 - O sudo mv nextflow /usr/local/bin/
- Verify installation:
 - O nextflow -v

B. Installing Required Tools for the Pipeline Without Docker

1. Installing SnpEff(with and without docker for Windows and Linux)

Build snpEff database

1. Download SnpEff (If Docker is already installed, skip this download step and continue with the remaining steps).

Go to the SnpEff download page.

Download the .zip file and extract it:

 $wget\ https://snpeff.blob.core.windows.net/versions/snpEff_latest_core.zip$

unzip snpEff_latest_core.zip

cd snpEff

Set Up Environment: Ensure SnpEff is executable:

chmod +x snpEff/snpEff.jar

snpeff - version // This ensures that if BWA is already installed.

Generating a .bin file from SnpEff for my_bacteria:(Windows and linux with docker)

- Navigate to the SnpEff directory (inside the container):
 cd /opt/snpeff/snpEff
- Check if your species database is available:

java -jar snpEff.jar databases | grep "my_bacteria"

Download the database for your species:
 java -jar snpEff.jar download my_bacteria

After downloading, SnpEff will automatically generate the required .bin file for my_bacteria.(snpEff/data/my_bacteria/.bin)

If your species directory is missing files except for .bin, you can copy the entire species folder using the following command:

cp -r /opt/snpeff/snpEff/data/Pseudomonas_aeruginosa_pao1 /home/docker/

Generating a .bin file from SnpEff for my bacteria:(Linux without docker)

- Check if your species database is available:
 java -jar snpEff.jar databases | grep "my_bacteria"
- Download the database for your species:
 java -jar snpEff.jar download my_bacteria

After downloading, SnpEff will automatically generate the required .bin file for my_bacteria.(snpEff/data/my_bacteria/.bin)

2. Obtain Genome and Annotation Files

- Download your bacterial genome in FASTA format (.fasta or .fna) from NCBI, Ensembl Bacteria, or another source.
- Obtain the corresponding GFF annotation file (.gff).

3. Use Prokka to Format Files for SnpEff

Format Files for SnpEff:

- Ensure the FASTA file contains a single genome sequence (or properly labeled contigs).
- Verify that the GFF file has accurate annotations, especially for CDS regions.

If annotations are missing or need to be generated, use Prokka:

Note:

- Replace my_bacteria with the actual bacterial name.
- The output files my_bacteria.fasta and my_bacteria.gff will be used in SnpEff.

Run:

prokka --outdir ./prokka_output --prefix my_bacteria my_bacteria.fasta

4. Configure the SnpEff Database

Locate the snpEff.config file in snpEff directory and open the SnpEff directory

then snpEff.config file.

nano snpEff.config

Add a new entry at the end of the file above the Non-standard Databases section:

Note:

- Replace my_bacteria with a unique identifier for your genome.
- Replace MyBacteria with any name.

#---

Non-standard Databases

#---

my_bacteria.genome : MyBacteria

5. Set Up the Genome Directory

Inside the data folder in the SnpEff directory, check a subdirectory for your genome:

ls snpEff/data/my_bacteria

6. Copy the Genome and Annotation Files

Copy the genome FASTA and annotation GFF files into the directory:

cp my_bacteria.fasta snpEff/data/my_bacteria/sequences.fna cp my_bacteria.gff snpEff/data/my_bacteria/genes.gff

Use the Prokka output files

- my_bacteria.fasta sequences.fna
- my_bacteria.gff genes.gff

7. Build the SnpEff Database

Run the following command to build the database:

java -Xmx4g -jar snpEff/snpEff.jar build -gff3 -v my_bacteria

8. Handling Errors

If you encounter an error like this:

>>>WARNING_FILE_NOT_FOUND: Rare Amino Acid analysis: Cannot read protein sequence file 'path/to/snpEff/data/au/protein.fa ', nothing done.

ERROR: CDS check file 'path/to/snpEff/data/au/cds.fa 'not found.

ERROR: Protein check file 'path/to/snpEff/data/au/protein.fa 'not found.

ERROR: Database check failed.

00:00:01 Logging>>>

Note: You need to generate cds.fa and protein.fa files.

9. Generate Missing Files

Extract Coding Sequences (CDS)

Use the GFF3 and FASTA files to generate cds.fa.

Note: Take gff file from prokka.

gffread genes.gff -g sequences.fna -x cds.fa

This extracts CDS sequences and saves them in cds.fa.

Translate CDS to Proteins

Use transeq from the EMBOSS suite to generate protein.fa:

transeq -sequence cds.fa -outseq protein.fa

This generates protein.fa.

Place the Files in the Correct Directory

Move the generated files to the SnpEff directory:

mv cds.fa protein.fa path/to/snpEff/data/au/

10. Rebuild the Database

java -Xmx4g -jar snpEff/snpEff.jar build -gff3 -v au

2. Installing BWA

sudo apt update && sudo apt install -y bwa

Bwa - - help // This ensures that if BWA is already installed.

3. Installing Samtools

sudo apt update && sudo apt install -y samtools

Samtools - - help // This ensures that if Samtools is already installed.

4. Installing VCFtools

sudo apt update && sudo apt install -y vcftools

Vcftools - - help // This ensures that if VCFtools is already installed.

5. Installing BCFtools

sudo apt update && sudo apt install -y bcftools

Bcftools - - help // This ensures that if BCFtools is already installed.

6. Installing Fastp

conda install -c bioconda fastp

Fastp - - help // This ensures that if Fastp is already installed.

7. Installing Sed, Awk, and Cut

sudo apt update && sudo apt install -y sed awk coreutils

Coreutils - - help // This ensures that if coreutils is already installed.

8. Installing Python3

sudo apt update && sudo apt install -y python3 python3-pip

Python3 - - version // This ensures that if Python3 is already installed.

9. Installing SQLite

sudo apt update && sudo apt install -y sqlite3 || echo "SQLite3 is already installed or an error occurred."

Sqlite3 -- version // This ensures that if SQLite3 is already installed.

3. Installing and Running a Nextflow and Requirement Tools with Docker For Linux:

A. Running the Pipeline Manually

Without Docker in Linux:

• nextflow run main.nf {parameters}.

Example Nextflow command:

(base) virudhagiri@virudhagiri-HP-Laptop-15s-eq2xxx:~/Desktop/ngs\$ nextflow run /home/virudhagiri/Desktop/ngs/exam.nf --ref /home/virudhagiri/Desktop/ngs/snpEff --fna_read "/home/virudhagiri/Desktop/ngs/helo/*.fna" --main_dkvirudhagiri/Desktop/ngs/helo --new db /home/virudhagiri/Desktop/ngs/gui test/csvfy --ref strain 'pao1'

Or

Running the pipeline using GUI gives name, path, files and run nextflow. (if python 3 available).

With Docker in Linux:

- Update the apt package index:
 - O sudo apt-get update
- Install dependencies required for Docker:
 - O sudo apt-get install \
 - O ca-certificates \
 - O curl \
 - O gnupg \
 - O lsb-release \
 - O sudo
- Add Docker's official GPG key:
 - O curl -fsSL https://download.docker.com/linux/ubuntu/gpg | sudo gpg --dearmor -o /usr/share/keyrings/docker-archive-keyring.gpg
- Set up the stable repository for Docker:
 - O echo "deb [arch=amd64 signed-by=/usr/share/keyrings/docker-archive-keyring.gpg] https://download.docker.com/linux/ubuntu \$(lsb_release -cs) stable" | sudo tee /etc/apt/sources.list.d/docker.list > /dev/null
- Update the apt package index again:
 - O sudo apt-get update
- Install Docker Engine:

- O sudo apt-get install docker-ce docker-ce-cli containerd.io
- Start and enable Docker:
 - O sudo systemctl start docker
 - O sudo systemctl enable docker
- Verify that Docker is installed and running:
 - O sudo docker --version
- Build Docker:
 - O sudo docker build -t my_app.
- Run Docker:
 - O sudo docker run -it -v /:/app/ my app /bin/bash
- Inside docker change directory using cd.. then enter the "app" directory.
- @root: /app/ nextflow run main.nf {parameters}.

Or

- The Docker folder must contain the gui.py file, and it should be named exactly as gui.py.
- Navigate to the Docker folder cd /path/to/docker
- Enable X11 access: xhost +
- Run the Docker container:
- sudo docker run -e DISPLAY=\$DISPLAY -v /tmp/.X11-unix:/tmp/.X11-unix
 -v /home/admin:/home/admin -v /home/admin/Desktop/docker:/app my_app
- If any generated file appears locked or cannot be opened, it is likely stored with root ownership. To fix this, change the ownership to your local user.

Check File Ownership:

- ls -l /home/path/to/finalbn_database.db
- If the file is owned by root,then

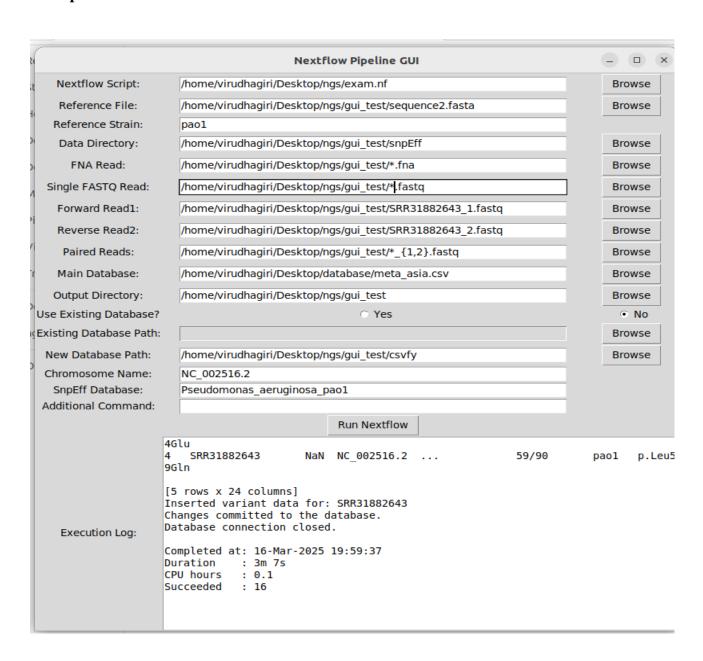
Change Ownership:

- sudo chown sastra:sastra /home/path/to/finalbn_database.db
- Now, the file should be accessible on your system.
- sastra local system user.

Note:

- -e DISPLAY=\$DISPLAY Passes the display environment variable for GUI.
- -v /tmp/.X11-unix:/tmp/.X11-unix Shares the X11 socket for graphical output.
- -v /home/admin:/home/admin Mounts your home directory inside the container.
- -v /home/admin/Desktop/docker:/app Mounts your project directory inside the container which contains the gui.py file.
- my_app The name of your Docker image.

Example GUI:



B. Input Parameters

: nextfow run path/to/main.nf

- -ref Provide the reference genome file for a species (.fasta or .fna).(eg., path/to/sequence.fna or fasta)
- -datadir Specify the path to the SnpEff directory. (eg., path/to/snpEff)
- -fna_read Provide assembled genome files (.fna) or the path containing .fna files.
- eg., Give "path/to/*.fna"
- **-read** Provide raw sequencing reads (.fastq) or the path containing .fastq files.(Single end fastq)
- eg., Give "path/to/*.fastq"
- -read1 Provide the forward raw read file (.fastq) with _1 as a delimiter
 (e.g., path/to/SRR12333_1.fastq).
- -read2 Provide the reverse raw read file (.fastq) with _2 as a delimiter (e.g., path/to/SRR12333_2.fastq).
- -reads Provide the path containing multiple .fastq files.(Paired end fastq)
 eg., Give "path/to/* {1,2}.fastq "
- -main_db Provide a CSV metadata file containing details of assembled genome .fna files or raw reads.(eg., path/to/meta.csv)
- -- ref_strain Provide the name of the reference strain as it appears in the FASTA file.(eg., 'pao1')
- -output_dir Specify the path where output CSV and VCF files will be stored. This should be a new directory with nothing inside it. (e.g., path/to/any directory)
- -chr_name Provide the name of the chromosome as it appears in the FASTA file.(eg., 'NC_257237')
- -snpEff_db Provide the name of the SnpEff database directory.(eg., 'psedo_au')
- -exist_db Provide the path and filename of an existing database file (if available).(eg., path/to/finaldata.db)

-new_db - If no existing database is available, provide this parameter along with the path where the new database file should be stored. Importantly, the specified directory should not contain any existing database file. (e.g., path/to/any directory)

C.Input Type

Ref: A reference FASTA or FNA file.

Fna read:

- Can be a single file (e.g., path/to/287.1738.fna).
- Can be multiple files stored in a directory (path/to/*.fna).
- If using fna_read, the file extension must be .fna; otherwise, do not use this parameter.

Read:

- Can be a single file (e.g., path/to/SRR89484.fastq or path/to/ERR83938.fastq).
- Can be multiple files stored in a directory (path/to/*.fastq).
- If using read, the file extension must be .fastq, and the filename must start with SRR or ERR.
- If the filename does not start with SRR or ERR, rename it accordingly or do not use this parameter.

Read 1 and Read 2:

- Can be single files:
 - O path/to/SRR89484_1.fastq and path/to/SRR89484_2.fastq.
 - O path/to/ERR83938_1.fastq and path/to/ERR83938_2.fastq.
- If using read, the file extension must be .fastq, and the filename must start with SRR or ERR.
- 1 and 2 must be added before the file extension.
- If the filename does not follow this format, rename it accordingly or do not use this parameter.

Reads:

- Can be a single paired-end file (e.g., path/to/SRR89484_1 and _2.fastq or path/to/ERR83938_1 and _2.fastq .fastq).
- Can be multiple files stored in a directory (path/to/*_{1,2}.fastq).
- The files must not be clipped (e.g., each file should contain both forward and reverse reads).
- If using read, the file extension must be .fastq, and the filename must start with SRR or ERR.
- If the filename d
- paooes not follow this format, rename it accordingly or do not use this parameter.

Main db:

- The file extension must be .csv.
- The file must contain the following fields(columns):
 - O antibiotic, genome_id, genome_name (strain name), sra_accession(with or without data), isolation_country, geographical_group, taxon_id, and resistant_phenotype.
- Other fields are optional.

4. Visualization

A. Web Interface (If python3 is available and not working in docker environment)

- 1. Loads your generated database file (SQLite) to host a local server
- 2. Runs using Streamlit with the command:

streamlit run path/to/app.py path/to/database.db

- 3. Redirects to a web interface where you can provide filtering options for variant data.
- 4. Displays Variant table and interactive plots for analysis..

Without Docker Installation:(Linux)

- pip install streamlit
- echo 'export PATH=\$HOME/.local/bin:\$PATH' >> ~/.bashrc
- source ~/.bashrc

- streamlit --version
- streamlit run path/to/app.py path/to/database.db

With Docker Run:(Linux)

sudo docker run -it -v /:/home/ -p 8501:8501 my_app streamlit run /home/path/to/stream_test_sqlite/Home.py /home/path/to/db_index_few_check_all_last1.db --server.port=8501 --server.address=0.0.0.0

Note:

- **docker run -it** Runs the container in interactive mode with a terminal.
- -v /:/app/ Mounts the entire host filesystem (/) to /app/ inside the container
- -p 8501:8501 Maps port 8501 from the container to the host for Streamlit access.
- my_app Specifies the Docker image to run (my_app).
- streamlit run /app/path/to/Home.py /app/path/to/db_index_few_check_all_last1.db Runs Streamlit with the given script and database.(ensure both files are inside /app/, which is the mounted directory.)
- --server.port=8501 Runs Streamlit on port 8501.
- --server.address=0.0.0.0 Makes Streamlit accessible externally.

This will create and host a link to the web interface.

With Docker Run:(Windows)

docker run -it -v C:\:/home/ -p 8501:8501 my_app streamlit run /home/path/to/Home.py /home/path/to/db_index_few_check_all_last1.db --server.port=8501 --server.address=0.0.0.0 Note:

- docker run -it Runs the container in interactive mode with a terminal.
- -v C:\:/home/ Mounts the entire C: drive from Windows to /home/ inside the container.

•

- -p **8501:8501** Maps port 8501 from the container to the host, allowing external access to Streamlit.
- my_app Specifies the Docker image to run (my_app).
- streamlit run /home/path/to/Home.py /home/path/to/db_index_few_check_all_last1.db Runs Streamlit with the given script and database.(ensure both files are inside /home/, which is the mounted directory.)
- --server.port=8501 Runs Streamlit on port 8501.
- --server.address=0.0.0.0 Makes Streamlit accessible from any device on the network.

This will create and host a link to the web interface.

Note: Minimum 16 GB RAM require.(depends on database size)

B. Sqlite Browser Software

 You can use this software or any SQLite-based software to visualize the database, search, and analyze data using SQL queries.

5. Troubleshooting

Common Issues

- GUI problem: Check whether Xvxcsrc is running or not. If not, click XLaunch and run it..(windows)
- Check whether the parameter is provided correctly, including case sensitivity as uppercase letters will also cause an error.
- Nextflow command not found: Ensure Nextflow is in your system PATH or docker environment PATH.

- snpEff command error: The snpEff configuration file is not found, or the data directory is missing.
 - O If you encounter this issue, follow the snpEff manual (as mentioned in the Installation section) to properly recreate the snpEff database, ensure the correct path to the configuration file, and export the path using .bashrc. Additionally, set the necessary permissions using chmod 777 on the snpEff folder.
- Java not found: Install OpenJDK using sudo apt install openjdk-21-jdk.