# ACEMFS FUT Minna Bioinformatics Workshop

## Sequence Retrieval & Quality Control using Galaxy

**Itunuoluwa Isewon PhD**
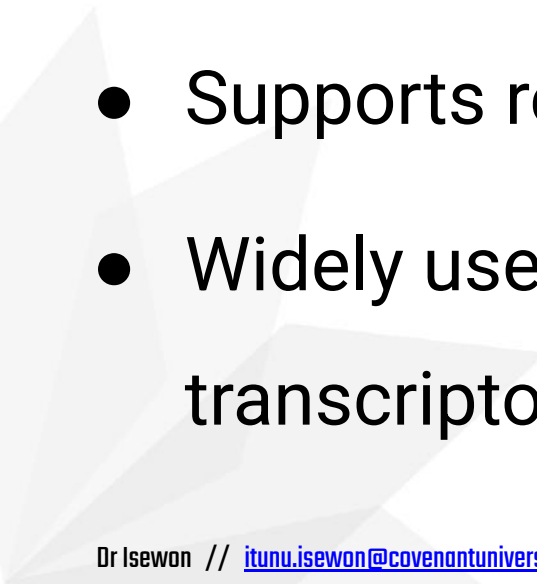Covenant University

# **Workshop Outline**

By the end of this session, participants will be able to:

● Retrieve sequences of mycotoxins & fungal enzymes from databases.

● Upload & organize datasets in **Galaxy**.

● Perform QC using Galaxy tools (**FastQC**)

# Introduction to Galaxy

- Web-based platform for bioinformaticians.

- No coding required.

- Supports reproducible research.

- Widely used for genomics, proteomics and transcriptomics analysis.

# Sequence Retrieval

# Sequence Retrieval

NCBI Exploration:

1.  Go to **NCBI**

2.  In the search bar, paste the **lcc9 gene**.

3.  Click on the dropdown arrow next to the search bar and select **Nucleotide**.

4.  Click the search button.

# Sequence Retrieval

# Sequence Retrieval

5. On the results page click the **first hit**.

6. What is the name of the **organism**?

7. How many **base pairs** does it have?

8. What is its **accession** and **version numbers**?

# Sequence Retrieval

# Sequence Retrieval

9. On the right, click on the dropdown arrow next to Send to:

10. Select file under, Choose **Destination**.

11. Change the format to **FASTA**

12. Click on **Create File** to download the FASTA File.

13. **Repeat** this process for a gene interesting to you.

# Sequence Retrieval

# BLAST

# BLAST

- **BLAST** stands for Basic Local Alignment Search Tool.

- It identifies similarities between biological sequences by comparing nucleotide or protein sequences to a database of sequences.

# BLAST

- **BLASTn (Nucleotide BLAST):** compares one or more nucleotide query sequences to a subject nucleotide sequence or a database of nucleotide sequences.

- **BLASTx (translated nucleotide sequence searched against protein sequences):** compares a nucleotide query sequence that is translated in six reading frames against a database of protein sequences.

- And **many others.**

# BLAST

1. Choose **blastn.**

2. Paste the accession number **EF990898.1** in the Enter Query Sequence box.

3. Select **Nucleotide collection nr/nt** under Database in Choose Search Set section.

4. Under Program Selection choose **Highly similar sequences(megablast).**

5. Set the max target sequences under General Parameters to **50.**

6. Click the **BLAST** button.

| | Descriptions | Graphic Summary | Alignments | Taxonomy |

## Fusarium acuminatum isolate 1A chromosome 5

Sequence ID: CP151264.1   Length: 4969144   Number of Matches: 2

Range 1: 2399923 to 2402423 GenBank   Graphics    ▼ Next Match ▲ Previous Match

| Score | Expect | Identities | Gaps | Strand |
|---|---|---|---|---|
| 1760 bits(953) | 0.0 | 2013/2530(80%) | 51/2530(2%) | Plus/Minus |

```
Query  986     TTTGCTGCCCGCGGCGTTAGCTGCGACGGTGTCTTATGACTTTTCTATTGATTGGGTTCG  1045
               |||| ||||| ||||| || |||||| |||||||||||||||| ||  |||||||||||||
Sbjct  2402423 TTTGATGCCGGCGGCTTTGGCTGCTACGGTGTCTTATGATTTCACTATTGAATGGGTCAG  2402364

Query  1046    AGCAAATCCAGATGGCGCGTTTGAGAGGTCGACGATAGGCATTAATAGAGAGTGGCCGAT  1105
               |||| |||| |||||| ||| ||||||| ||| || ||||| |||| ||| |  ||||||||
Sbjct  2402363 AGCGAATCCTGATGGCGCCTTTGAGAGGCCTACGATTGGCATCAATGGGCGGTGGCCGAT  2402304

Query  1106    ACCGAGGATTGAAGCGAGTATTGGGGATACGGTTTTGGTTTATGTGAGGAATAATTTGGG  1165
                || ||||| || |||| |  || |||||| ||||||  ||| |||||||||||| ||||
Sbjct  2402303 TCCCAGGATCGAGGCGACTGTGGGTGATACGATTTTGGTGAATGCGAGGAATAATCTGGG  2402244

Query  1166    GAATCAGTCTACGAGTTTGCATTTTCATGGGCTTTTCATGAATGGCTCGAATCATATGGA  1225
               |||||||||  ||| || ||||||| ||| ||| ||| |||||||||  || ||||||||
Sbjct  2402243 GAATCAGTCCACGTCGTTGCATTTTCACGGTCTGTTTATGAATGGTTCAAACCATATGGA  2402184

Query  1226    TGGGCCGTCGCAGGTTACGCAGTGCCCTATTCAACCTGGAGAGTCATTTCTCTATAACTT  1285
               ||||||||||| ||||||||| ||  ||||||||| |||| |||||||| ||||||||||
Sbjct  2402183 TGGGCCGTCGCAAGTGACGCAATGTCCAATTCAGCCAGGCGAGTCGTTCCTCTATAACTT  2402124
```

# Quality Control

# Quality Control in Galaxy

Obtain **FASTQ** files:

1. Go to **Galaxy**

2. In the tool Bar, click on **Get Data**

3. Choose "**Faster Download and Extract Reads in FASTA/Q format from NCBI SRA**".

4. In the Accession tab, write the accession number of the fastq file: **SRR18453616**.

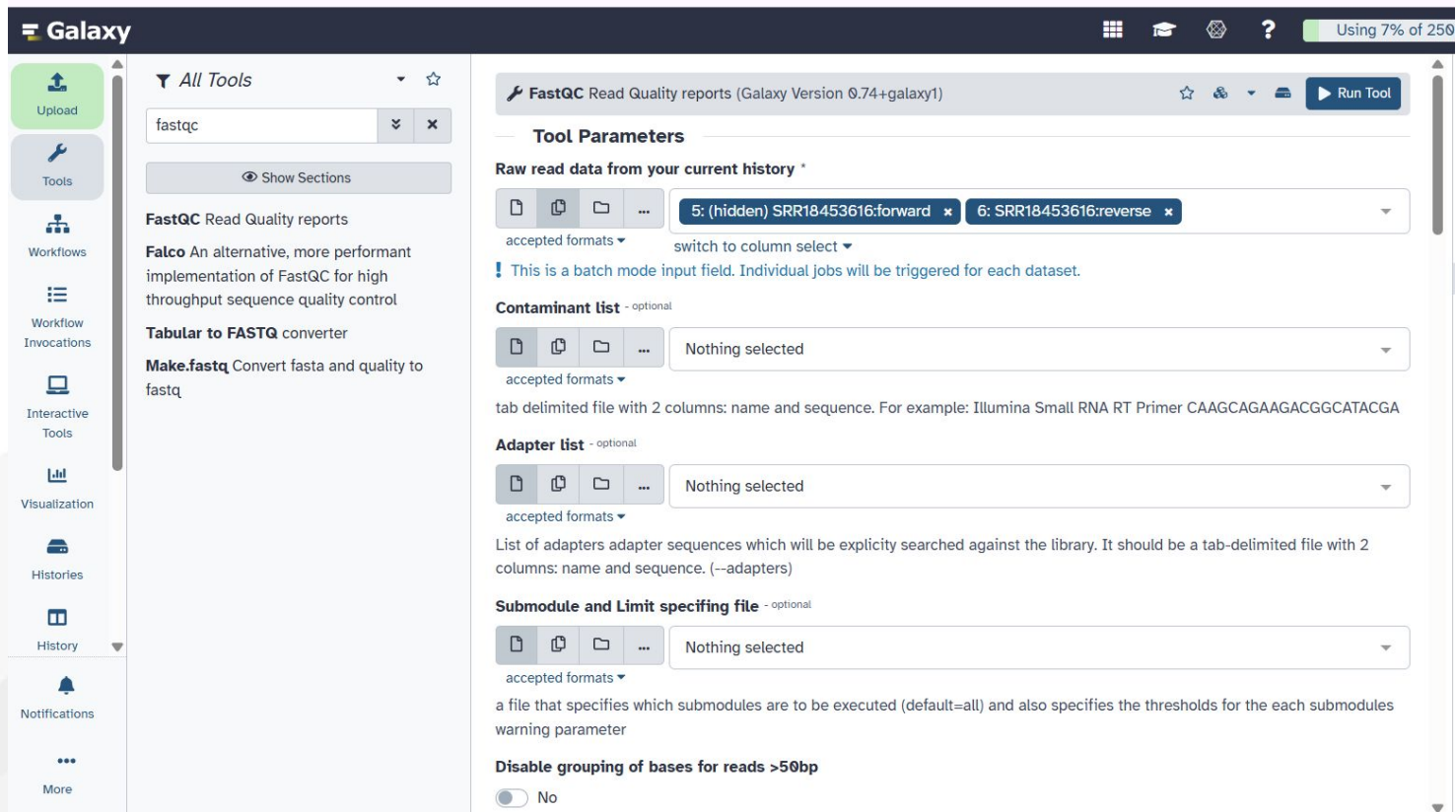# Quality Control in Galaxy

# **Quality Control in Galaxy**

Perform **QC**:

1. On the search bar, type **fastqc**.

2. **Choose** the desired fastq file (paired end) in the raw read tab.

3. **Leave** all other tabs unchanged.

# Quality Control in Galaxy

# Quality Control in Galaxy

4. Once it runs, two files are generated, a raw data file and a Web Page file.

5. **View** the result by clicking on the web page file produced.

# Quality Control Result

## Summary

✅ Basic Statistics

✅ Per base sequence quality

✅ Per sequence quality scores

❌ Per base sequence content

✅ Per sequence GC content

✅ Per base N content

✅ Sequence Length Distribution

❌ Sequence Duplication Levels

✅ Per base sequence quality

Quality scores across all bases (San

# Quality Control (MultiQC)

**Why:** It helps us to obtain a more intuitive
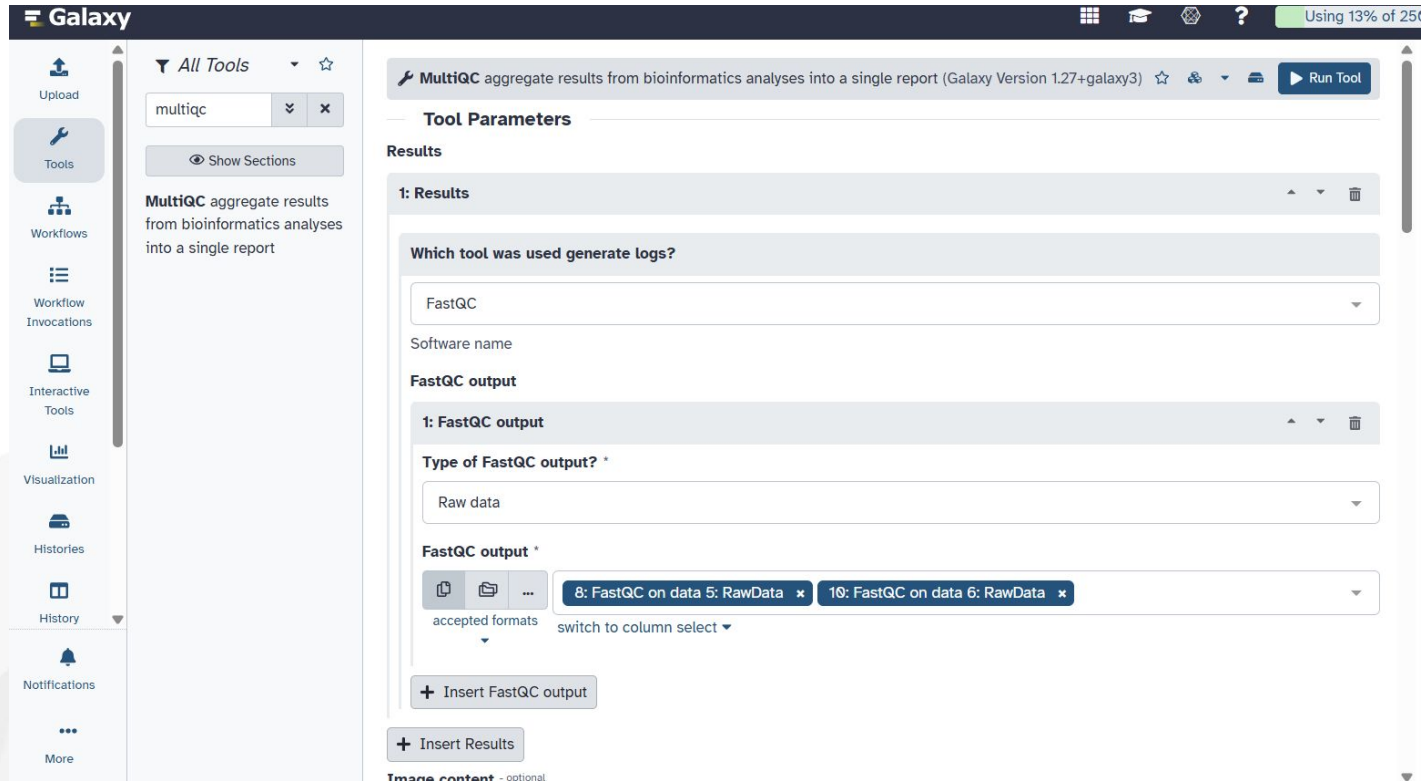
comparison

1.  On the search bar, type **multiqc**

2.  On the "Which tool was used generate logs?" tab,

    choose **Fastqc**

3.  Then click on "**Insert FastQC output**"

# Quality Control (MultiQC)

4. Type of output is raw data

5. Add the raw data files generated earlier

6. **Leave** all other parameters at default

7. **Run** tool

8. **View** the result by clicking on the web page file produced

# Quality Control (MultiQC)

# Quality Control Result

# Why QC is Important

- Ensures data integrity before downstream analysis.

- Detects contamination, errors, or poor-quality sequences.

- Prevents misleading results.

# Why QC is Important

- Ensures data integrity before downstream analysis.

- Detects contamination, errors, or poor-quality sequences.

- Prevents misleading results.