

# DEVELOPMENT OF A BIOINFORMATICS PIPELINE FOR ROUTINE ANALYSIS OF WHOLE GENOME SEQUENCING DATA OF *ESCHERICHIA COLI* ISOLATES

Bert Bogaerts<sup>1,3</sup>, Stéphanie Nouws<sup>1,3</sup>, Raf Winand<sup>1</sup>, Qiang Fu<sup>1</sup>, Julien Van Braekel<sup>1</sup>, Sarah Denayer<sup>2</sup>, Bavo Verhaegen<sup>2</sup>, Sigrid C. J. De Keersmaecker<sup>1</sup>, Nancy Roosens<sup>1</sup>, Kathleen Marchal<sup>3</sup> and Kevin Vanneste<sup>1</sup>

1) Transversal activities in applied genomics, Sciensano, Brussels (1050), Belgium

2) Foodborne pathogens, Sciensano, Brussels (1050), Belgium

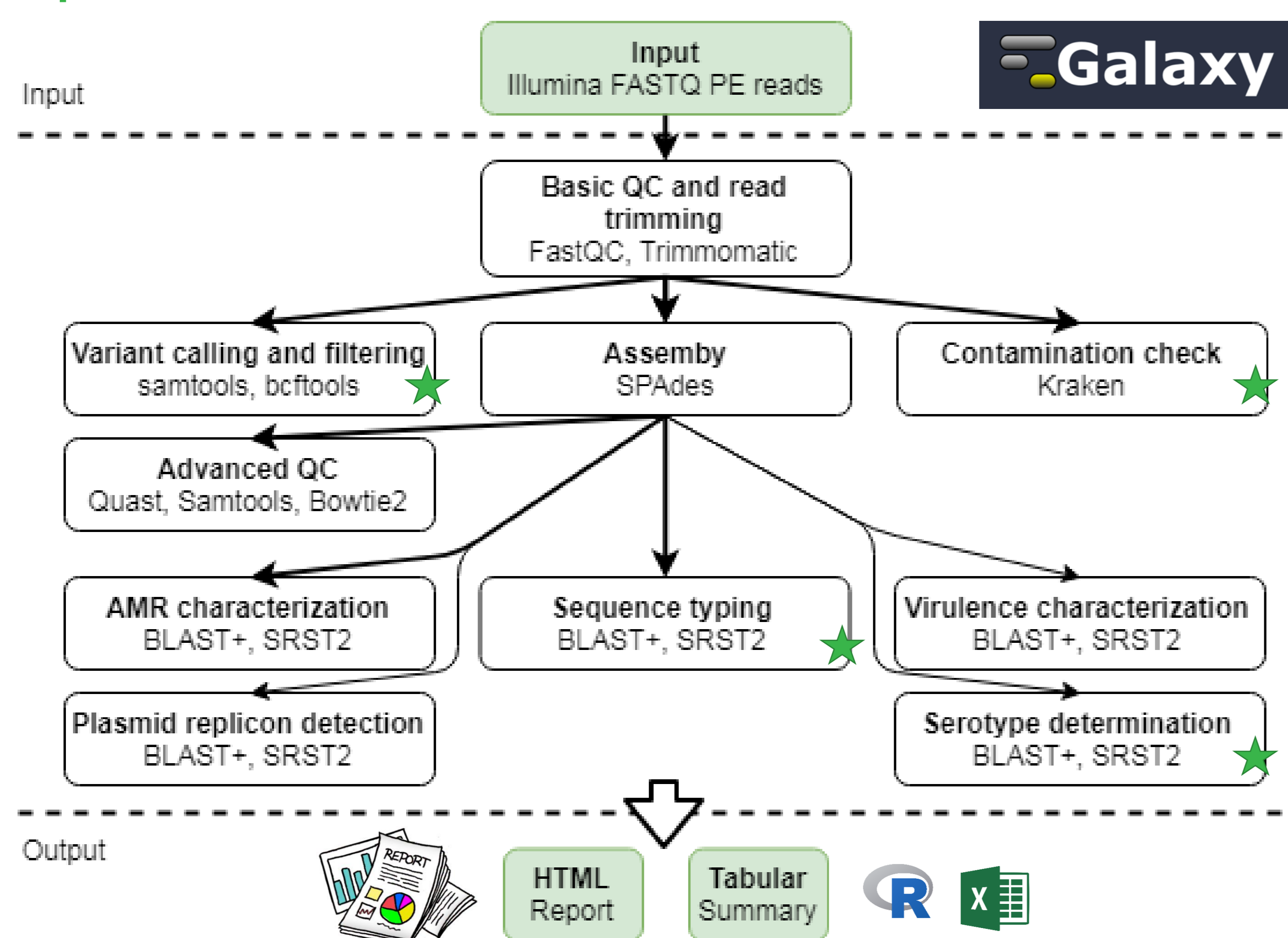
3) UGhent – Department of information Technology, IDLab, imec



## Abstract

The adaptation of whole genome sequencing (WGS) and bioinformatics for routine molecular typing and pathogen characterization in a public health setting remains problematic, which is partly due to the lack of user-friendly and validated data analysis tools that can be used for routine typing in the National Reference Laboratories (NRLs) and peripheral laboratories. In collaboration with the Belgian NRL for *Escherichia coli*, we developed a pipeline for the routine analysis of *E. coli* isolates that was specifically designed to tackle the aforementioned challenges.

## Pipeline architecture



The pipeline uses Illumina WGS data to perform a complete characterization of *Escherichia coli* isolates. The pipeline first generates basic quality reports before trimming the reads. Afterwards, a set of metrics specifically tailored for the pipeline are evaluated to determine if the data quality is sufficient to run the pipeline. The trimmed reads are then used as input for various downstream analyses: contamination check using Kraken, variant calling and filtering against the O157:H7 str. Sakai reference genome, and *de novo* assembly using SPAdes. All of the gene detection based assays can be performed by alignment using **blastn** or using a read mapping based approach with **SRST2**. The pipeline output is provided as either a **HTML report** or a tabular summary file which can be further analysed in other software.

★ Detailed on the rest of the poster

## A) Sequence typing

The pipeline performs sequence typing using the MLST scheme and cgMLST from PubMLST.org. Databases are automatically updated weekly. The allele detection is based on **BLAST+** alignment of the *de-novo* assembly or on read mapping using **SRST2** depending on the pipeline setting. The pipeline output is shown below. The detected sequence type is reported at the top, **color codes** are used to indicate the type of hit (see legend on the right).

MLST (BigSdb, Pasteur);  
MLST, cgMLST, wgMLST (EnteroBase, Warwick)

Classic MLST (Pasteur)

ST	Allele	% Identity	HSP/Locus length	Type	Alignment
563	11	100.00	450/450	DNA	<a href="#">view</a>
	72	100.00	516/516	DNA	<a href="#">view</a>
	134	100.00	468/468	DNA	<a href="#">view</a>
	52	100.00	450/450	DNA	<a href="#">view</a>
	25	100.00	456/456	DNA	<a href="#">view</a>
	145	100.00	561/561	DNA	<a href="#">view</a>
	18	100.00	594/594	DNA	<a href="#">view</a>
	2	100.00	600/600	DNA	<a href="#">view</a>

Hit type	Color
Perfect	Green
Imperfect identity	Yellow
Imperfect short	Orange
Multi-hit	Red
No hit	Grey

[Download \(TSV\)](#)

Last updated: 03-09-2018

## B) Serotype determination

Serotype determination is based on the detection of gene variants that determine the **H-type** (*fliC*, *fliA*, *flkA*, *flnA*, *fimA*) and **O-type** (*wzx*, *wzy*, *wzt*, *wzm*). Detection is based on **BLAST+** alignment of the *de-novo* assembly or on read mapping using **SRST2** depending on the pipeline setting. A set of decision rules is applied to determine the O-type, H-type and the combined serotype. If one of the types cannot be determined it is reported as 'ambiguous'. The database sequences are retrieved from the SerotypeFinder tool provided by DTU, and are automatically updated on a weekly basis.

### SerotypeFinder - O-type

Locus	Length	Coverage	Mismatches	Uncertainty	Depth	Predicted serotype	Accession
wzx_199	1287	98.76	16holes	edge2.0	6.08	O157	AKMA01000036
wzy_201	1185	99.41	7holes	edge0.0	4.62	O157	JH953200

[Download \(TSV\)](#)

Last updated: 03-09-2018

### SerotypeFinder - H-type

Locus	Length	Coverage	Mismatches	Uncertainty	Depth	Predicted serotype	Accession
fliC_15	1758	100.00	0	edge1.0	7.71	H7	AF228487

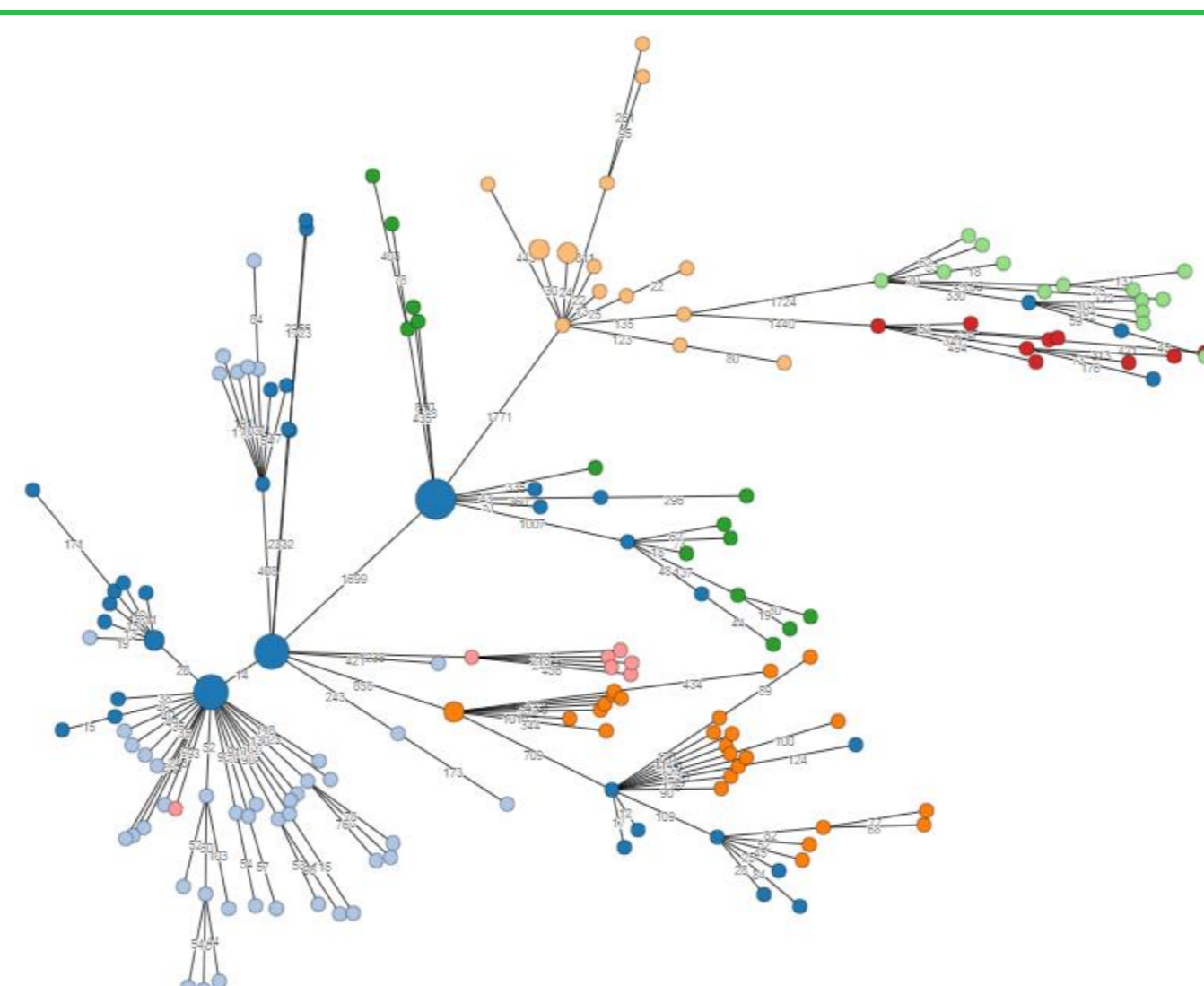
[Download \(TSV\)](#)

Last updated: 03-09-2018

Detected serotype: **O157:H7**

## Validation

We will extensively validate the different assays of the pipeline to demonstrate the employed methods are "fit-for-purpose" and provide high-quality results. Our approach will be similar to one employed in a previous pipeline validation for *Neisseria meningitidis* (Bogaerts et al, <https://doi.org/10.3389/fmicb.2019.00362>). We will evaluate repeatability, reproducibility, accuracy, precision, sensitivity, and specificity of the different bioinformatics assays. The validation strategy will be extended to also include the relatedness between isolates based on the sequence typing output and variant calling assays.



## Conclusion and future work

We present a pipeline for the complete characterization of *Escherichia coli* isolates using (Illumina) WGS data. The pipeline performance is currently being characterized by means of a set of performance metrics and definitions that were specifically adapted towards bioinformatics assays, and which evaluate repeatability, reproducibility, accuracy, sensitivity, precision, and specificity. Preliminary results on a representative set of samples demonstrate high performance, indicating the feasibility of using WGS in routine public health settings to replace classically employed pathogen typing and characterization techniques. Similar pipelines can be developed for other pathogens and case studies, making bioinformatics analyses less complex and more time-efficient for both expert and non-expert users.

## C) Contamination check

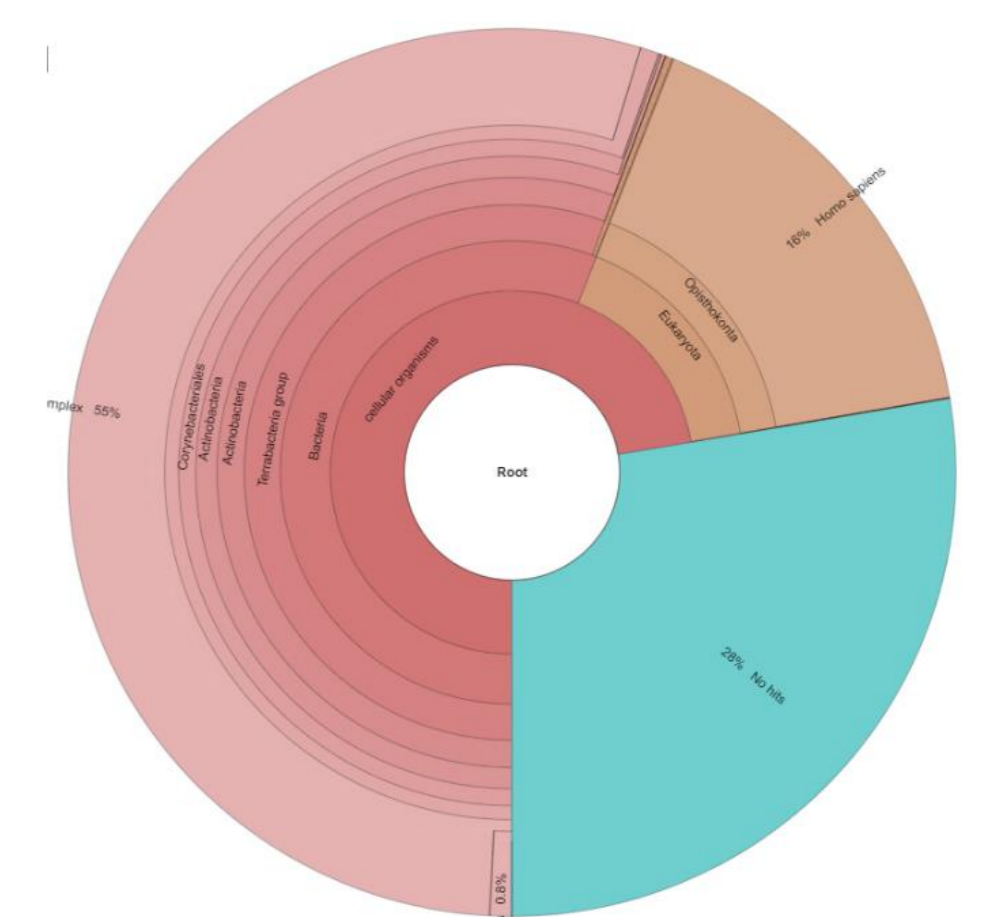
The pipeline uses **Kraken** to check for possible contaminants in the input sample. Warnings and errors are raised for unexpected species that have over 1% and 5% K-mer abundance, respectively. Results are visualized in an interactive **Krona** report.

Species	Percentage
Expected	
<i>Escherichia coli</i>	40.23
Contaminants	
None found	

Percentage	Level	Name
25.58	Unclassified	unclassified
74.42		root
74.40		cellular organisms
74.39	Domain	Bacteria
74.33	Phylum	Proteobacteria
74.27	Class	Gammaproteobacteria
74.18	Order	Enterobacteriales
73.63	Family	Enterobacteriaceae
60.45	Genus	<i>Escherichia</i>
40.23	Species	<i>Escherichia coli</i>

[Krona Report](#)

Last updated: 24-01-2018



## D) Variant calling & filtering

**Bowtie2** is used to map the trimmed reads to the O157:H7 str. Sakai reference genome (NC\_002695.2). Afterwards, *samtools mpileup* followed by *bcftools call* is used to call variants.

Reference: *Escherichia coli* O157:H7 (NC\_002695.2)

### Read mapping

Mapping rate
73.99%



### Filtering

Filter	Variants passed
Depth	64772/67523
SNP quality	63828/64772
Mapping quality	45225/63828
Distance	39216/45225
Z-score	38986/39216

### Output files

	Number of variants	VCF file
Unfiltered	67523	<a href="#">Download</a>
Filtered (All positions)	38986	<a href="#">Download</a>

