

Bioinformatics tools for WGS analysis of infectious diseases

Pimplapas Leekitcharoenphon (Shinny)

Researcher

Research Group for Genomic Epidemiology

DTU Food

Whole-genome sequence
based species ID using
KmerFinder

KmerFinder - bitbucket

<https://bitbucket.org/genomicepidemiology/kmerfinder/src/master/>

Installation

Setting up KmerFinder program

```
# Go to wanted location for resfinder
cd /path/to/some/dir
# Clone and enter the KmerFinder directory
git clone https://bitbucket.org/genomicepidemiology/kmerfinder.git
cd kmerfinder
```

Build Docker image from Dockerfile

```
# Build container
docker build -t kmerfinder .
# Run test
docker run --rm -it \
    --entrypoint=/test/test.sh kmerfinder
```

KmerFinder database - bitbucket

https://bitbucket.org/genomicepidemiology/kmerfinder_db/src/master/

Installation - Clone repository

```
# Go to wanted location for KmerFinder db
cd /path/to/some/dir
# Clone and enter the KmerFinder db directory
git clone https://bitbucket.org/genomicepidemiology/kmerfinder_db.git
cd kmerfinder_db
```

K-mer ?

- A k-mer is a contiguous sequence of k bases
- k is any positive integer
- Sequences with high similarity must share k-mers

sequence **ATGGAAGTCGCGGAATC**

7 mers

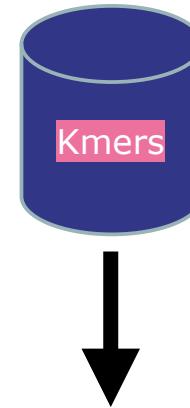
```
ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC
```

Species identification by K-mer

Known species **ATGGAAGTCGCGGAATC**

k-mers

ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

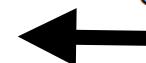


ATGGAAGTCGCGGAATC

Unknown species

k-mers

ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC



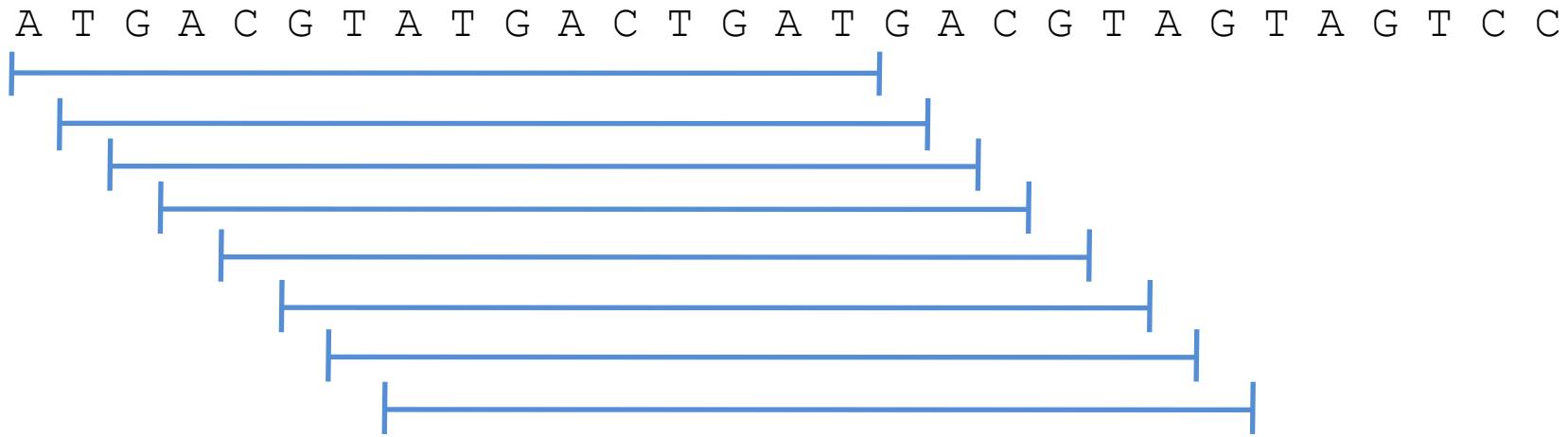
Species

KmerFinder

Using all information in the WGS data

almost

- Genomes is split into 16mers:



- It is necessary to reduce the total amount of 16mers:
 - Only 16mers with particular prefix are (ATGAC) kept

Query bacteria of
unknown species



Reference db bacteria of
known species (template)



Bact1->*E. coli*

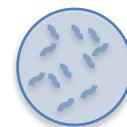


Bact2->*S. enterica*

KMA (k-mer alignment)



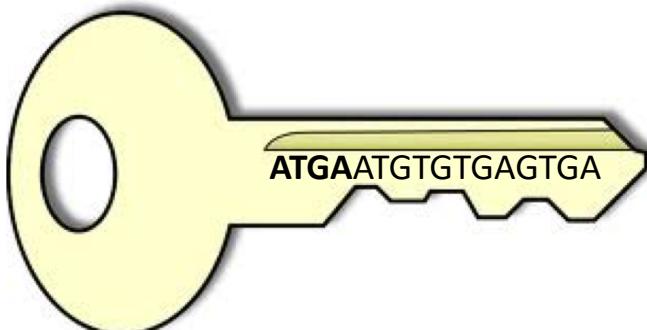
Bact3->*K. pneumoniae*



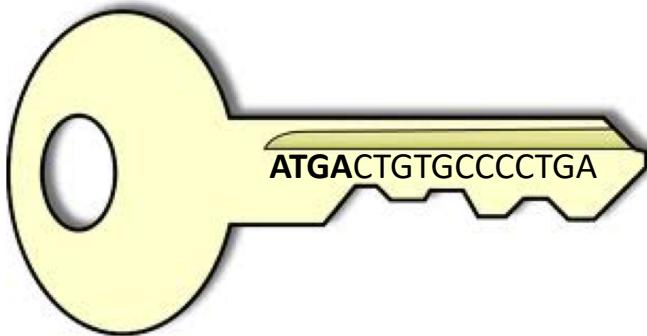
Bact4->*S. aureus*

Prediction: Query bacteria is a *S. aureus*

16mer database



CP001921 (*Acinetobacter baumanii*)
CP000521 (*Acinetobacter baumanii*)
CP002522 (*Acinetobacter baumanii*)



CP001921 (*Acinetobacter baumanii*)
CP002301 (*Buchnera aphidicola*)

Unknown bacteria

16 mers:

ATGAATGTGTGAGTGA

ATGACTGTGCCCTGA

ATGAAAAAAAAAAAAAA



Species	Match	No. of Kmer hits
<i>Acinetobacter baumannii</i>	CP001921	2
<i>Acinetobacter baumannii</i>	CP000521	1
<i>Acinetobacter baumannii</i>	CP002521	1
<i>Buchnera aphidicola</i>	CP002301	1

KmerFinder output – standard scoring method

Center for Genomic Epidemiology

[Home](#)[Services](#)[Instructions](#)[Output](#)

KmerFinder-3.1 Server - Results

KmerFinder 3.1 results:

Template	Num	Score	Expected	Template_length	Query_Coverage	Template_Coverage	Depth
NC_016854.1 Salmonella enterica subsp. enterica serovar Typhimurium str. D23580 complete genome	7004	6094006	30	157485	96.86	99.99	38.70

[EXTENDED OUTPUT](#)

Input Files: *Salmonella-spp-02-03-002_R1_001.trim.fq* *Salmonella-spp-02-03-002_R2_001.trim.fq*

[RESULTS as text \(tab separated\)](#)

Other KmerFinder statistics

$$\text{Query coverage} = \frac{\text{Score (total number of kmers in query sequence that match kmers in template sequence)}}{\text{Total number of kmers in query sequence}}$$
$$\text{Template coverage} = \frac{\text{Score (total number of kmers in template sequence that match kmers in query sequence)}}{\text{Total number of kmers in template sequence (database sequence)}}$$
$$\text{Depth} = \frac{\text{Score (total number of kmers in query sequence that match kmers in template sequence)}}{\text{Total number of kmers in template sequence (database sequence)}}$$

Explanation of the columns in standard and extended output

The following contains a briefly explanation of all columns of the output including the columns in the extended output

Template: shows the accession numbers or name of the template sequences

Assembly: RefSeq assembly accession ID

Num: is the sequence number of accession entry in the KmerFinder database

Score: is the total number of matching Kmers between the query and the template

Expected: is the expected score, i.e. the expected total number of matching Kmers between query and template (randomly selected).

Template length: is the number of Kmers in the template

query_coverage [%]: is the percentage of input query Kmers that match the template.

Coverage [%]: is the template coverage.

Depth: is the number of matched kmers in the query sequence divided by the total number of Kmers in the template. For read files this estimates the sequencing depth.

tot_query_coverage [%]: is calculated based on the ratio of the score and the number of kmers in the query sequence, where the score includes kmers matched before.

tot_coverage [%]: is calculated based on ratio of the score and the number of unique kmers in the template sequence, where the score includes kmers matched before.

tot_depth: depth value based on all query kmers that can be found in the template sequence .

q_value: is the quantile in a standard Pearson Chi-square test, to test whether the current template is a significant hit.

p_value: is the p-value corresponding to the obtained q_value.

Accession number: accession number of entry ID in fasta file.

Description: additional descriptions available in fasta file, or in the case of organism databases the identifier lines of fasta files.

TAXID: NCBI's TaxID number of the hit

Taxonomy: complete taxonomy of the hit

TAXID Species: NCBI's species TaxID number of the hit (sometimes bacterial strain or substrain TaxIDs can be given above)

Species: Species name

Dataset: 5 unknown genomes

/home/projects/course_23262/course/week03/assembly

/home/projects/course_23262/course/week03/trimmed_reads

Guideline on how to run WGS tool:

DTU Learn > Course Content > Content > Week 3 > Exercises

Output table

<https://docs.google.com/spreadsheets/d/1WmmmtTCFTVp7uOe9vYIRyDilyYHYbR-liPiEyoPkbZo4/edit?usp=sharing>

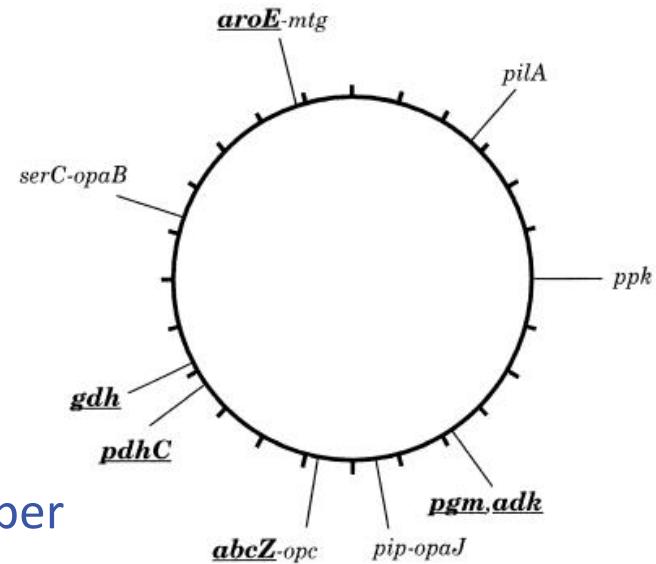
MLST typing

Multilocus Sequence Typing (MLST)

- the golden standard for typing

First developed in 1998 for *Neisseria meningitis*
(Maiden et al. PNAS 1998. 95:3140-3145)

- ❖ The nucleotide sequence of internal regions of app. 7 housekeeping genes are determined by PCR followed by Sanger sequencing
- ❖ Different alleles are each assigned a random number
- ❖ The unique combination of alleles is the sequence type (**ST**)



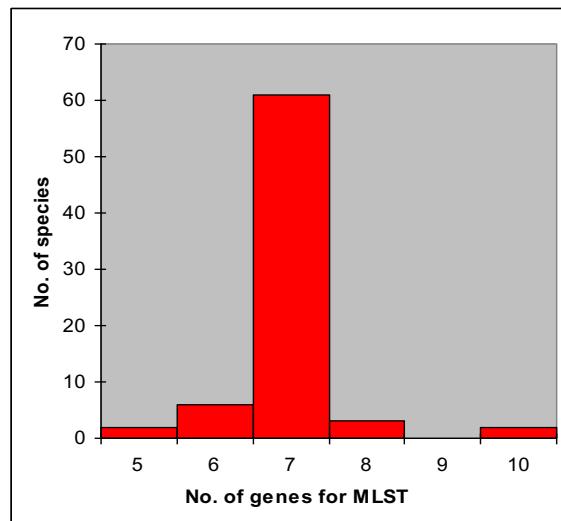
The first 4 ST profiles of *Neisseria meningitis*

Sequence type (ST)	Loci							
	<i>abcZ</i>	<i>adk</i>	<i>aroE</i>	<i>fumC</i>	<i>gdh</i>	<i>pdhC</i>	<i>pgm</i>	
1	1	3	1	1	1	1	3	
2	1	3	4	7	1	1	3	
3	1	3	1	1	1	23	13	
4	1	3	3	1	4	2	3	

A total of 446,824 different sequences (alleles) and 11,113 different Sequence Types exist (13. Feb 2015)

MLST schemes for ~100 species have now been developed and incorporated in the MLST database (called MLST configurations)

For most species, seven genes are used



Original Article

Genomics of an emerging clone of *Salmonella* serovar Typhimurium ST313 from Nigeria and the Democratic Republic of Congo

Pimplapas Leekitcharoenphon^{1,2}, Carsten Friis¹, Ea Zankari¹, Christina A. Svendsen¹, Lance B. Price³, Maral Rahmani¹, Ana Herrero-Fresno¹, Kayode Fashae⁴, Olivier Vandenberg^{5, 6}, Frank M. Aarestrup¹, Rene S. Hendriksen¹

¹Division of Bacterial Genomics and Epidemiology, WHO Collaborating Center for Antimicrobial Resistance in Food borne Pathogens and European Union Reference Laboratory for Antimicrobial Resistance, National Food Institute, Technical University of Denmark, Kemitorvet, Denmark

²Center for Biological Sequence Analysis, Department of System Biology, Technical University of Denmark, Kemitorvet, Denmark

³Division of Pathogen Genomics, Translational Genomics Research Institute (TGen), Flagstaff, USA

⁴Department of Microbiology, University of Ibadan, Ibadan, Nigeria.

⁵Infectious Diseases Epidemiological Unit, Public Health School, Université Libre de Bruxelles, Brussels, Belgium.

⁶Department of Microbiology, Saint-Pierre University Hospital, Brussels, Belgium.

Abstract

Introduction: *Salmonella enterica* serovar Typhimurium ST313 is an invasive and phylogenetically distinct lineage present in sub-Saharan Africa. We report the presence of *S. Typhimurium* ST313 from patients in the Democratic Republic of Congo and Nigeria.

Methodology: Eighteen *S. Typhimurium* ST313 isolates were characterized by antimicrobial susceptibility testing, pulsed-field gel electrophoresis (PFGE), and multilocus sequence typing (MLST). Additionally, six of the isolates were characterized by whole genome sequence typing (WGST). The presence of a putative virulence determinant was examined in 177 *Salmonella* isolates belonging to 57 different serovars.

Results: All *S. Typhimurium* ST313 isolates harbored resistant genes encoded by *bla_{TEM}1b*, *catA1*, *strA/B*, *sul1*, and *dfrA1*. Additionally, *aac(6')Iaa* gene was detected. Phylogenetic analyses revealed close genetic relationships among Congolese and Nigerian isolates from both blood and stool. Comparative genomic analyses identified a putative virulence fragment (ST313-TD) unique to *S. Typhimurium* ST313 and *S. Dublin*.

Conclusion: We showed in a limited number of isolates that *S. Typhimurium* ST313 is a prevalent sequence-type causing gastrointestinal diseases and septicemia in patients from Nigeria and DRC. We found three distinct phylogenetic clusters based on the origin of isolation suggesting some spatial evolution. Comparative genomics showed an interesting putative virulence fragment (ST313-TD) unique to *S. Typhimurium* ST313 and invasive *S. Dublin*.

Implementation of MLST using WGS data

- Automatic, local download of MLST database once a week
- For each species, we get an ST profile table file and an allele file

ST profile table file

ST	leuA	petC	malF	cysG	holC	nuoL	gltT
1	1	1	1	1	1	1	1
2	1	1	4	1	1	1	1
3	1	1	1	20	1	1	1
4	1	1	1	4	1	1	1
5	2	2	2	2	2	2	2
6	3	3	3	3	3	3	3
7	3	3	3	7	3	3	3
9	3	3	5	5	4	3	4
10	5	4	3	3	6	3	5
11	7	7	7	9	10	8	8
13	7	6	7	9	10	7	8
14	8	8	8	11	12	9	9
16	7	6	8	10	11	8	8
17	1	1	10	12	18	10	1
18	9	1	9	13	14	5	10
19	10	1	10	14	15	11	1
20	1	1	10	12	17	11	11
21	10	1	10	14	15	11	12
26	5	3	3	3	6	3	5
33	11	9	14	15	19	13	10
39	3	3	5	19	4	3	7

Allele file

```
>leuA_1
TACTATCAATGCCATTGGTGAGCGCCCTGGTAACTGGCCCTGGAAGAACTCACTATGGT
GTTGAAAGTACGCAACCGCTTTACAACATTGATACTTCATCCACACATCACGTATCGT
CTCCACCTCCCAGTTACTGCAACGATTGGTGGCATGCCGTGCAACGTAACAAGGCAGT
AGTAGGTGCCAATGCCATTGACATGAATCGGTATCCACAGCACGGTATGCTGCCA
TCGGCGCACCTACGAAATCATGCCCTCAAGAAGTCGGITGGTATGTCGATATGGT
GCTCGGCCGCCATAGCGGCCGTGCTGCCGTGAAACAGCGTCTACCGCAGTGGCTACTT
GCTGGAGGAAGAACATCTAAAACGGTATTTGAAGAACATTCAAACAGCTTTGTGAGAAACA
GCGTTGGTACCGATGTCGACCTGCAAGTACTGATGCAAGATAAACAGTACAGCATGG
CTATCGTTGGCTCAATGACAATCACTGATGTTGTAACCGGCCAACGCACTGGTGG
ATTGTCCGATCCCCAAGGTACGCGTGTGGCGAAACTGCGCAAGGCAACGGCCAGTAGA
TGCACGTGTCGGAGCGCTGGCTGCAAGCAACTGGGTCAAACGGTAACTGGAGTTGGACAGCTATCA
GGTACACAGTGGTATGGGCGCATGCACTGGTGAAGCAAACCT
>leuA_2
TACTATCAATGGCATTGGTGAGCGCGCTGGTAACTGGCGCTGGAAGAACTCACTATGGT
GTTGAAAGTACGCAACCGCTTTACAACATTGATACTTCATCCACACATCACGTATCGT
CTCCACCTCCCAGTTACTGCAACGATTGGTGGCATGCCGTGCAACGTAACAAGGCAGT
AGTAGGTGCCAATGCCATTGACATGAATCGGTATCCATCACGGTATGCTGCCA
TCGGCGCACCTACGAAATCATGCCCTCAAGAAGTCGGITGGTATGTCGACATGGT
ACTCGGCCGCCATAGCGGCCGTGCTGCCGTGAAACAGCGTCTACCGCAGTGGCTACTT
GCTGGAGGAAGAACATCTAAAACGGTATTTGAAGAACATTCAAACAGCTTTGTGAGAAACA
GCGTTGGTACCGATGTCGACCTGCAAGTACTGATGCAAGATAAACAGTACAGCATGG
CTATCGTTGGCTCAATGACAATCACTGATGTTGTAACCGGCCAACGCACTGGTGG
```

Implementation of MLST using WGS data

Method

- The genome is converted to a blast database
- For the specified species, all alleles for all genes are aligned to the database (using blast for assembled genomes/ using KmerFinder (KMA) for raw reads)
- For each gene, the perfectly matching allele is picked (all nucleotides must match across the whole length of the allele)
- If there is no perfectly matching allele, the closest matching allele is outputted along with warnings
- The ST is determined from the combination of alleles

MLST in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/mlst/src/master/>

Installation

Setting up MLST program

```
# Go to wanted location for mlst
cd /path/to/some/dir
# Clone and enter the mlst directory
git clone https://bitbucket.org/genomicepidemiology/mlst.git
cd mlst
```

Build Docker container

```
# Build container
docker build -t mlst .
# Run test
docker run --rm -it \
    --entrypoint=/test/test.sh mlst
```

#Download and install MLST database

```
# Go to the directory where you want to store the mlst database
cd /path/to/some/dir
# Clone database from git repository (develop branch)
git clone https://bitbucket.org/genomicepidemiology/mlst_db.git
cd mlst_db
MLST_DB=$(pwd)
# Install MLST database with executable kma_index program
python3 INSTALL.py kma_index
```

If kma_index has not bin install please install kma_index from the kma repository: <https://bitbucket.org/genomicepidemiology/kma>

MLST in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/mlst/src/master/>

Usage

The program can be invoked with the -h option to get help and more information of the service. Run Docker container

```
# Run mlst container
docker run --rm -it \
    -v $MLST_DB:/database \
    -v $(pwd):/workdir \
    mlst -i [INPUTFILE] -o . -s [SPECIES] [-x]
```

When running the docker file you have to mount 2 directory: 1. mlst_db (MLST database) downloaded from bitbucket 2. An output/input folder from where the input file can be reached and an output files can be saved. Here we mount the current working directory (using \$pwd) and use this as the output directory, the input file should be reachable from this directory as well.

-i INPUTFILE input file (fasta or fastq) relative to pwd -s SPECIES species origin of input file -o OUTDIR output directory relative to pwd
-x extended output. Will create an extented output

<https://cge.cbs.dtu.dk/services/MLST>

Center for Genomic Epidemiology



Username
Password
[New](#) [Reset](#) [Login](#)

Home Services Instructions Output Article abstract

MLST 2.0 (Multi-Locus Sequence Typing)

Software version: [2.0.4 \(2019-05-08\)](#)
Database version: [2.0.0 \(2020-02-10\)](#)

Select MLST configuration
MLST allele sequence and profile data are obtained from [PubMLST.org](#).

Please note that for four organisms, two or three different MLST schemes are available:

- *Acinetobacter baumannii* (*Acinetobacter baumannii* #1 [\[1\]](#), *Acinetobacter baumannii* #2 ([link](#))).
- *Escherichia coli* (*Escherichia coli* #1 [\[4\]](#), *Escherichia coli* #2 [\[5\]](#)).
- *Pasteurella multocida* (*Pasteurella multocida* #1 (*RIRDC*), *Pasteurella multocida* #2 (*multihost*)).
- *Leptospira* (*Leptospira* #1, *Leptospira* #2, *Leptospira* #3).

Select type of data input
Only data from one single isolate should be uploaded. If raw sequencing reads are uploaded KMA will be used for mapping. KMA supports the following sequencing platforms: Illumina, Ion Torrent, Roche 454, SOLiD, Oxford Nanopore, and PacBio.

Please note that "Assembled Genomes/Contigs" should be selected, if you have already assembled your short sequencing reads into one continuous genome or into several contigs. It is indifferent which type of short sequence reads were used to produce the genome/contigs.

Name	Size	Progress	Status

[Upload](#) [Remove](#)

Center for Genomic Epidemiology

[Home](#)[Services](#)[Instructions](#)[Output](#)

MLST-1.7 Server - Typing Results

Sequence Type: *Unknown ST*

Locus	% Identity	HSP Length	Allele Length	Gaps	Allele
acs	99.74	390	390	0	acs_28
aro	100.00	349	498	0	aro_122
gua	100.00	373	373	0	gua_11
mut	100.00	442	442	0	mut_11
nuo	100.00	366	366	0	nuo_4
pps	100.00	370	370	0	pps_12
trp	100.00	443	443	0	trp_3

Please note that one or more loci do not match perfectly to any previously registered MLST allele.
We recommend verifying the results by traditional methods for MLST!

[extended output](#)

MLST Profile: *paeruginosa*

Organism: *Pseudomonas aeruginosa*

Input Files: *18_tag18.contigs.fa*

Extended output

Imperfect match, *acs* loci

```
acs: WARNING, Identity: 99.74%, HSP/Length: 390/390, Gaps: 0, Best Match: acs_28

MLST allele seq: ggcccggtggccaacggcgccaccaccattctgttcgagggcgtgccgaactaccccgac
Hit in genome: ggcccggtggccaacggcgccaccaccattctgttcgagggcgtgccgaactaccccgac

MLST allele seq: gtgaccgcgtggcgaaaatcatcgacaaggcacaaggtaaacatccttacaccgcgcgg
Hit in genome: gtgaccgcgtggcgaaaatcatcgacaaggcacaaggtaaacatccttacaccgcgcgg

MLST allele seq: accgcgatccgcgcgatgtatggccgaaggcaaggcggcggtggccggtgccgacggttcc
Hit in genome: accgcgatccgcgcgatgtatggccgaaggcaaggcggcggtggccggtgccgacggttcc

MLST allele seq: agcctgcgtctgctcggttcgggtggcgagccatcaaccggaaagcctggcagtggtagc
Hit in genome: agcctgcgtctgctcggttcgggtggcgagccatcaaccggaaagcctggcagtggtagc

MLST allele seq: tacgagaccgtcgcccagtcgcgtcccgtatcgacacactgtggcagaccgagacc
Hit in genome: tacgagaccgtcgcccagtcgcgtcccgtatcgacacactgtggcagaccgagacc

MLST allele seq: ggcgcctgcctgatgaccccgctgccggcgcccacgcgatgaagccggctctgcagcc
Hit in genome: ggcgcctgcctgatgaccccgctgccggcgcccacgcgatgaagccggctctgcagcc

MLST allele seq: aagccgttcttcggcgtggtaccggcactg
Hit in genome: aagccgttcttcggcgtggtaccggcactg
```

Extended output

Imperfect match, *aro* loci

aro: WARNING, Identity: 100%, HSP/Length: 349/498, Gaps: 0, *aro_122* is the best match for *aro*

MLST allele seq:	atgtcaccgtgccgttcaaggaagaggcctatcgctctggacgaattgagcgagcggg
Hit in genome:	
MLST allele seq:	ccaccccggccggggcggtgaacaccctgatccgcctgccgacggtcgcctgcgcggcg
Hit in genome:	
MLST allele seq:	acaacacccgacggcgccggcttgcgtgcgggacacctgacggcgaacgcccgggtcgagctgc
Hit in genome:	gacctgacggcgaacgcccgggtcgagctgc
MLST allele seq:	gcggcaagcgggttctcctgctcggcgccggcggtgcggtgcgtgggtgctcgaaaccct
Hit in genome:	gcggcaagcgggttctcctgctcggcgccggcggtgcggtgcgtgggtgctcgaaaccct
MLST allele seq:	tcctcggcgagtgcccggcgagttgcgtatcgccaaccgcacggcgccgaaaggccgtgg
Hit in genome:	tcctcggcgagtgcccggcgagttgcgtatcgccaaccgcacggcgccgaaaggccgtgg
MLST allele seq:	acctggccgagcgggttcggccgacacctcgccgcgggtgcacggctgcggtttcggccgaggtcg
Hit in genome:	acctggccgagcgggttcggccgacacctcgccgcgggtgcacggctgcggtttcggccgaggtcg
MLST allele seq:	aaggcccttcgacctgatcgtaacggcacctcgccagttgcggccgacgtgcgcgc
Hit in genome:	aaggcccttcgacctgatcgtaacggcacctcgccagttgcggccgacgtgcgcgc
MLST allele seq:	cgctggcgagacgtgatcgagccccggcgatccgtctgtacgacatgtatgc
Hit in genome:	cgctggcgagacgtgatcgagccccggcgatccgtctgtacgacatgtatgc
MLST allele seq:	aggaaccgactgcctca
Hit in genome:	aggaaccgactgcctca

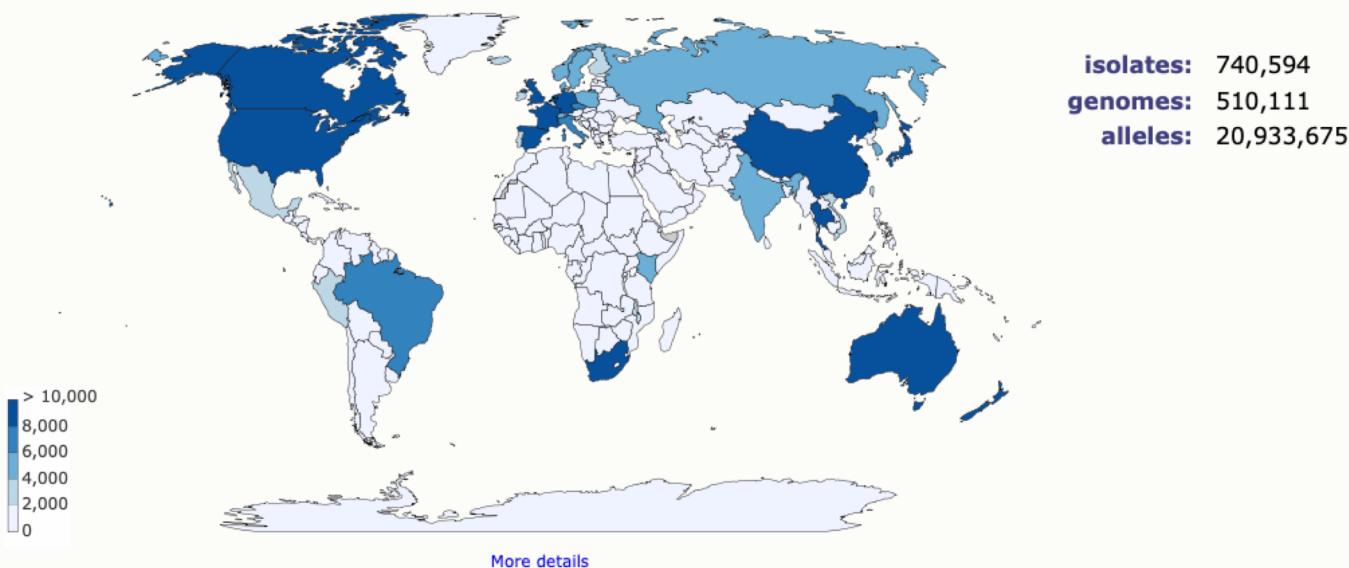
<https://pubmlst.org>

PubMLST

Welcome to PubMLST - Public databases for molecular typing and microbial genome diversity.

- [Log in or register for a PubMLST account](#)
- [Databases](#)
- [Species identification by ribosomal MLST](#)

Source of isolates submitted to PubMLST



This site is hosted at [The Department of Zoology](#), University of Oxford, UK and is funded by The Wellcome Trust.

Is super MLST, e.g. rMLST, the next step?

Microbiology. 2012 Jan 27. [Epub ahead of print]

Ribosomal Multi-Locus Sequence Typing: universal characterisation of bacteria from domain to strain.

Jolley KA, Bliss CM, Bennett JS, Bratcher HB, Brehony CM, Colles FM, Wimalaratna HM, Harrison OB, Sheppard SK, Cody AJ, Maiden MC.

Abstract

No single genealogical reconstruction or typing method currently encompasses all levels of bacterial diversity, from domain to strain. We propose **Ribosomal Multi Locus Sequence Typing (rMLST)**, an approach which **indexes variation of the 53 genes encoding the bacterial ribosome protein subunits (rps genes)**, as a means of integrating microbial genealogy and typing. As with MLST, rMLST employs curated reference sequences to identify gene variants efficiently and rapidly. **The rps loci are ideal targets for a universal characterization scheme as they are: (i) present in all bacteria; (ii) distributed around the chromosome; and (iii) encode proteins which are under stabilising selection for functional conservation.** Collectively, the rps loci exhibit variation that resolves bacteria into groups at all taxonomic and most typing levels providing significantly more resolution than 16S small subunit rRNA gene phylogenies. A web-accessible expandable database, comprising whole genome data from more than 1900 bacterial isolates, including 28 draft genomes assembled de novo from the EBI sequence read archive, has been assembled. The rps gene variation catalogued in this

Ribosomal Multilocus Sequence Typing (rMLST)

Ribosomal Multilocus Sequence Typing (rMLST) is an approach that indexes variation of the 53 genes encoding the bacterial ribosome protein subunits (*rps* genes) as a means of integrating microbial taxonomy and typing.

The *rps* gene variation catalogued in this database permits rapid and computationally non-intensive identification of the phylogenetic position of any bacterial sequence at the domain, phylum, class, order, family, genus, species and strain levels.

rMLST is described in [Jolley et al. 2012 *Microbiology* 158:1005-15](#).

Identify species

[No login required]

genomes: 343,524
alleles: 2,145,412
profiles: 133,582

- Databases
 - [rMLST sequence definitions](#)
 - [Bacterial genomes](#)

***C. jejuni/C. coli* core genome MLST (cgMLST) scheme**

A *C. jejuni* / *C. coli* core genome MLST (cgMLST) scheme v1.0 has been defined from 2,742 Oxfordshire Human Surveillance genomes, for use in the analysis of human disease isolates.

Loci appearing in 95% or more of genomes were detected, using the Genome Comparator function of BIGSdb, from the 1,643 loci defined by the re-annotation of the NCTC11168 genome (Gundogdu et al., 2007). Potential paralogues were identified in 5 subsets of 10 isolates by running Genome Comparator and excluding paralogous loci from distance matrix calculations; firstly, when loci were paralogous in all isolates (default BLAST settings) and secondly when loci were paralogous in any isolates (min 90% sequence similarity BLAST setting). Sequence similarities to loci thus detected were identified in further loci using [GeneDB](#). This resulted in the identification of 22 potential paralogous loci which were removed from the original 95% core list to provide a cgMLST scheme of 1,343 loci.

This scheme is available for selection within the sequence definition and isolate databases. It can also be accessed via the [RESTful API](#).

Resistant gene detection: ResFinder

Source	Accessibility	Approach	Type	Year	Status
Resfinder	Web and standalone	Assembly-based and read based tool	Tool and database	2012	Active update regularly
ARG-ANNOT	Standalone	Assembly-based tools	Tool and database	2014	Archived last update in May 2018
ResfinderFG	Web	Assembly-based tool	Tool and database	2016	Active last update in November 2016
RGI	Web and/or standalone	Assembly-based tools	Tool	2015	Active
ARGs-OAP (v2)	Web and/or standalone	Assembly-based tools	Tool	2016	Active
ARIBA	Standalone	Assembly-based tools	Tool	2017	Active
NCBI-AMRFinder	Standalone	Assembly-based tools	Tool	2018	Active
SRST2	Standalone	Read-based tools	Tool	2014	Active
SEAR	Web and/or standalone (archived)	Read-based tools	Tool	2015	Archived
ShortBRED	Standalone	Read-based tools	Tool	2015	Active
PATRIC	Web	Read-based tools	Tool	2016	Active
SSTAR	Standalone	Read-based tools	Tool	2016	Active
KmerResistance	Web	Read-based tools	Tool	2016	Active
GROOT	Standalone	Read-based tools	Tool	2018	Active
DeepArgs	Web	Read-based tools	Tool	2018	Active

Source	Accessibility	Type	Year	Status
CARD	Web	Database	2013	Active updated monthly
Resfams	Web	Database	2015	Active last update in January 2015
ARDB	Web	Database	2009	Archived last updated in 2009
MEGARes	Web	Database	2016	Active; last update in December 2016
NDARO	Web	Database	2016	Active; started in 2016
Mustard	Web	Database	2018	Active; last update in November 2018
FARME database	Web	Database	2017	Active; last update in 2017
SARG (v2)	Web	Database		Active
Lahey list of β -lactamases	Web	Database	2015	Archived; last update in 2015
BLDB	Web	Database	2018	Active; last update in November 2018
LacED	Web	Database		TEM LacED active: last update in 2017; SHVED archived: last update in April 2010
CBMAR	Web	Database		Last update in September 2014
MUBII-TB-DB	Web	Database		Last update in December 2013
u-CARE	Web	Database		Last update in 2016

ResFinder



- ResFinder is based on curated database, public databases as well as on scientific papers
- The ResFinder is a web-friendly interface and freely accessible tool (It is also a stand-alone tool)
- ResFinder will detect the presence of whole resistance genes, AND chromosomal point mutations causing resistance in the whole genome sequence data (raw reads or assembled genomes)
- ResFinder 4.0 provides in silico antibiograms as reliable as those obtained by phenotypic antimicrobial susceptibility testing
- High concordance (>95%) between phenotypic and predicted antimicrobial susceptibility was observed. Discrepancies were mainly linked to criteria for interpretation of phenotypic tests and suboptimal sequence quality, and not to ResFinder 4.0 performance.

- Bortolaia V. et al. ResFinder 4.0 for predictions of phenotypes from genotypes. *J. Antimicrob Chemother.* 2020 Aug 11;dkaa345

Zankari E, Hasman H, Kaas RS, Seyfarth AM, Agersø Y, Lund O, Larsen MV, Aarestrup FM. Genotyping using whole-genome sequencing is a realistic alternative to surveillance based on phenotypic antimicrobial susceptibility testing. *J. Antimicrob Chemother.* 2013 Apr;68(4):771-7.

ResFinder command line in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/resfinder/src/master/>

ResFinder tool

Setting up ResFinder script and database

```
# Go to wanted location for resfinder
cd /path/to/some/dir

# Clone the latest version and enter the resfinder directory
git clone https://git@bitbucket.org/genomicepidemiology/resfinder.git
cd resfinder
```

Dependencies:

Depending on how you plan to run ResFinder BLAST and KMA can be optional. BLAST is used to analyse assemblies (ie. FASTA files). KMA is used to analyse read data (ie. FASTQ files).

Python modules: Tabulate, BioPython, CGECORE and Python-Git

To install the needed python modules you can use pip

```
pip3 install tabulate biopython cgecore gitpython python-dateutil
```

For more information visit the respective website

```
https://bitbucket.org/astanin/python-tabulate
https://biopython.org
https://bitbucket.org/genomicepidemiology/cge_core_module
https://gitpython.readthedocs.io/en/stable/index.html
```

BLAST (optional)

If you don't want to specify the path of blastn every time you run ResFinder, make sure that blastn is in your PATH.

Blastn can be obtained from:

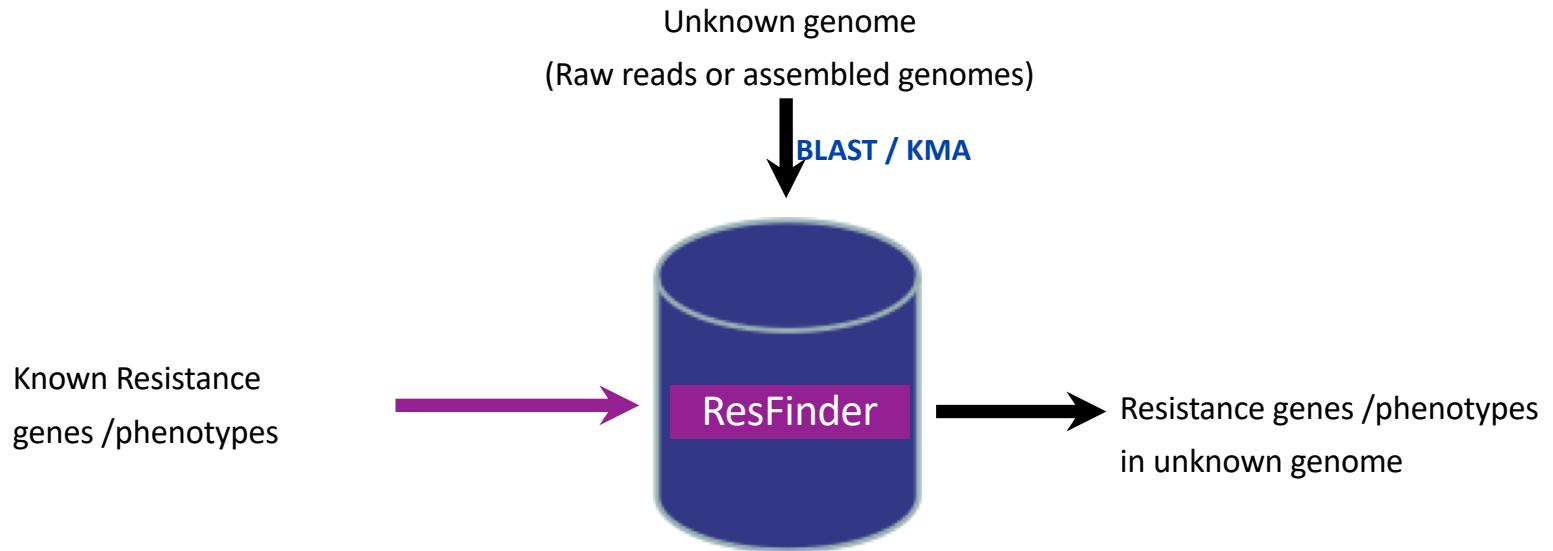
```
ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/
```

ResFinder 4.0 database

ResFinder 4.0 contains 4 databases

- AMR genes (ResFinder) https://bitbucket.org/genomicepidemiology/resfinder_db/src/master/
- Chromosomal point mutation (PointFinder) https://bitbucket.org/genomicepidemiology/pointfinder_db/src/master/
- Translation of genotypes into phenotypes
- Species-specific panels for in silicon antibiograms

ResFinder tool



Center for Genomic Epidemiology



Username

Password

New Reset Login

[Home](#) [Services](#) [Instructions](#) [Output](#) [Overview of genes](#) [Article abstract](#)

ResFinder 4.0

ResFinder identifies acquired genes and/or finds chromosomal mutations mediating antimicrobial resistance in total or partial DNA sequence of bacteria.

The database is curated by:
Frank Møller Aarestrup
(click to contact)

Server updated (20:25 30-Aug-2019 GMT+1): Fixed bug causing point mutations in the database not to be shown on the results page of the website if the required mutations for resistance wasn't found.

Server updated (10:46 28-Aug-2019 GMT+1): Updated databases and made several bug fixes.

Server updated (11:19 28-Jan-2019 GMT+1): Updated antibiotic class definitions and corrected some spelling.

Server updated (13:06 25-Jan-2019 GMT+1): Fixed issue with webserver caused by previous fix.

Server updated (11:28 25-Jan-2019 GMT+1): Fixed problem with certain jobs crashing while creating output for the web-tool. Fixed issue where certain antibiotics under certain conditions would be assigned to an incorrect class.

Chromosomal point mutations

Acquired antimicrobial resistance genes

Select species
 *Chromosomal point mutation database exists

Select type of your reads

If you get an "Access forbidden. Error 403": Make sure the start of the web adress is https and not just http. Fix it by clicking [here](#).

Name	Size	Progress	Status
Isolate File			

Upload **Remove**

ResFinder-4.0 Server - Results

Input Files: *strain_1.fasta*

Warning:

One or more resistance genes does not exist in the phenotype database. The Summary table does not take this into account.

escherichia coli	complete			
Antimicrobial	Class	WGS-predicted phenotype	Genetic background	
amikacin	aminoglycoside	No resistance		
tigecycline	tetracycline	No resistance		
tobramycin	aminoglycoside	No resistance		
cefepime	beta-lactam	Resistant	blaCTX-M-15 (blaCTX-M-15_AY044436)	
chloramphenicol	phenicol	No resistance		
piperacillin+tazobactam	beta-lactam	No resistance		
cefoxitin	beta-lactam	No resistance		
ampicillin	beta-lactam	Resistant	blaCTX-M-15 (blaCTX-M-15_AY044436), blaTEM-1B (blaTEM-1B_AY458016)	
ampicillin+clavulanic acid	beta-lactam	No resistance		
cefotaxime	beta-lactam	Resistant	blaCTX-M-15 (blaCTX-M-15_AY044436)	
ciprofloxacin	fluoroquinolone	Resistant	gyrA (p.S83A)	
colistin	polymyxin	No resistance		
sulfamethoxazole	folate pathway antagonist	Resistant	sul2 (sul2_HQ840942), sul1 (sul1_AY115475), sul1 (sul1_AY522923), sul1 (sul1_DQ914960), sul1 (sul1_U12338)	
imipenem	beta-lactam	No resistance		
trimethoprim	folate pathway antagonist	Resistant	dfrA7 (dfrA7_AB161450)	
nalidixic acid	fluoroquinolone	Resistant	gyrA (p.S83A)	
ertapenem	beta-lactam	No resistance		
tetracycline	tetracycline	Resistant	tet(A) (tet(A)_AJ517790)	
fosfomycin	fosfomycin	No resistance		
ceftazidime	beta-lactam	Resistant	blaCTX-M-15 (blaCTX-M-15_AY044436)	
temocillin	beta-lactam	No resistance		
gentamicin	aminoglycoside	No resistance		
meropenem	beta-lactam	No resistance		
azithromycin	macrolide	No resistance		

[Download phenotype table \(txt\)](#)

[Download species specific phenotype table \(txt\)](#)

Fluoroquinolone						
Mutation	Nucleotide change	Amino acid change	Phenotype	PMID	Notes	
gyrA:p.S83A	tcg -> gcg	s -> a	nalidixic acid,ciprofloxacin	8891148, 2168148, 12654733, 12654733		

Aminoglycoside									
Resistance gene	Identity	Alignment Length/Gene Length	Position in reference	Contig or Depth	Position in contig	Phenotype	PMID	Accession no.	Notes
aph(6')-Id	100.0	837/837	1..837	strain_1_contig_11	4362..5198	streptomycin	2653965	M28829	Alternative name strB
aph(3')-lb	100.0	804/804	1..804	strain_1_contig_11	3559..4362	streptomycin	12029529	AF321551	Alternative name strA

Tetracycline									
Resistance gene	Identity	Alignment Length/Gene Length	Position in reference	Contig or Depth	Position in contig	Phenotype	PMID	Accession no.	Notes
tet(A)	100.0	1200/1200	1..1200	strain_1_contig_11	15402..16601	doxycycline,tetra cycline	12654659	AJ517790	

Beta-lactam									
Resistance gene	Identity	Alignment Length/Gene Length	Position in reference	Contig or Depth	Position in contig	Phenotype	PMID	Accession no.	Notes
blaTEM-1B	100.0	861/861	1..861	strain_1_contig_14	84807..85667	amoxicillin,ampici llin,cephalothin,pi peracillin,ticarcillin	15388431	AY458016	Class A
blaCTX-M-15	100.0	876/876	1..876	strain_1_contig_14	81110..81985	amoxicillin,ampici llin,aztreonam,ce fepime,cefotaxim e,ceftazidime,ceft riaxone,piperacilli n,ticarcillin	11470367, 26169409	AY044436	Class A

Folate pathway antagonist									
Resistance gene	Identity	Alignment Length/Gene Length	Position in reference	Contig or Depth	Position in contig	Phenotype	PMID	Accession no.	Notes
sul1	100.0	761/840	1..761	strain_1_contig_9	442478..443238	sulfamethoxazole	unpublished	U12338	

Output



The grey colour

indicates a warning due to a non-perfect match, alignment length is shorter than resistance gene length, % ID <= 100%.

The light green colour

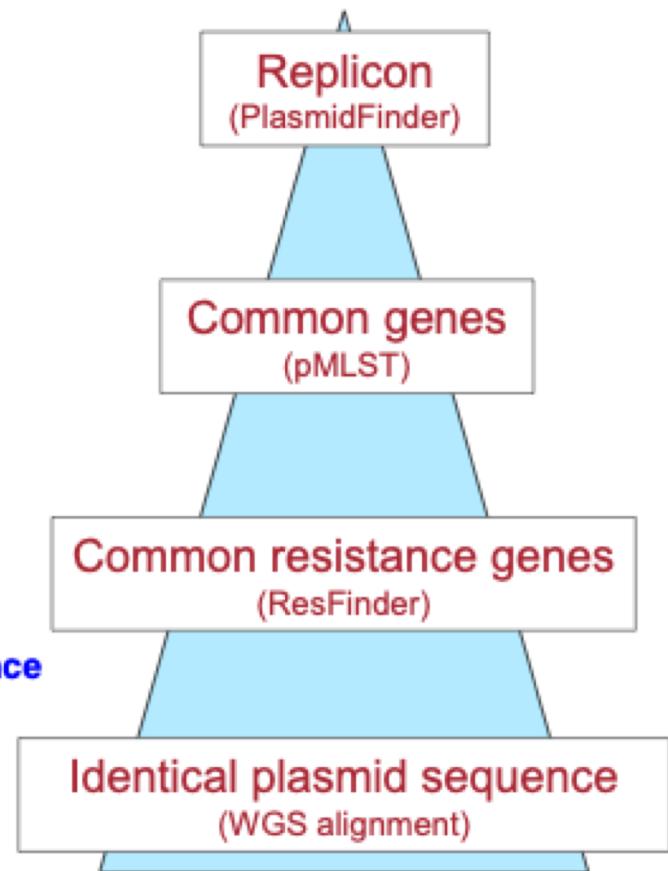
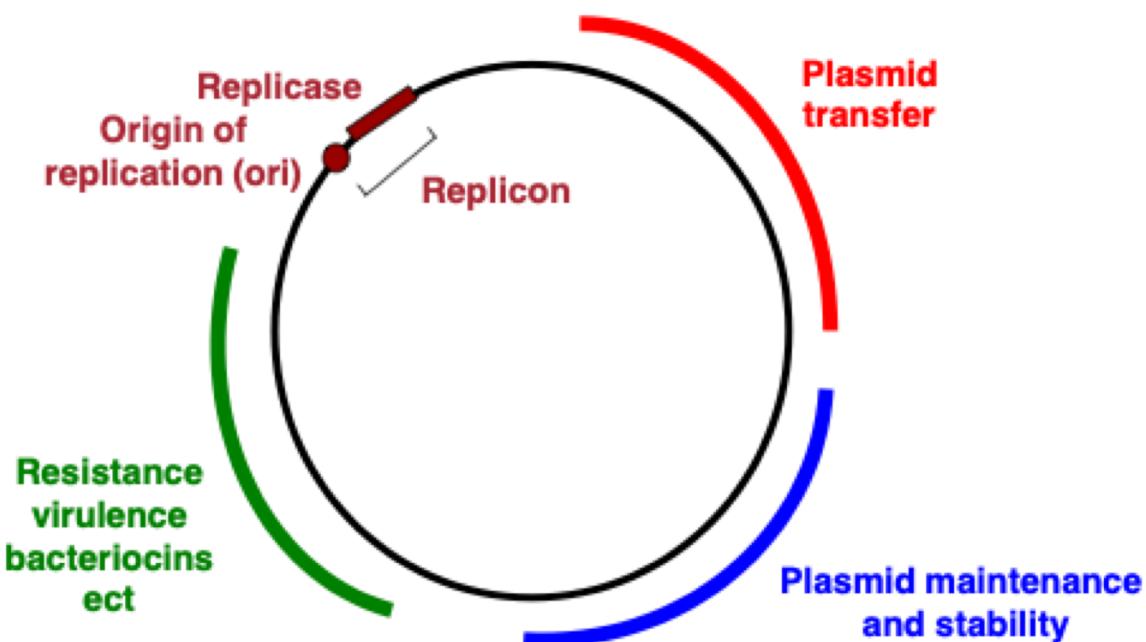
indicates a warning due to a non-perfect match, % ID < 100%, alignment length = resistance gene length.

PlasmidFinder

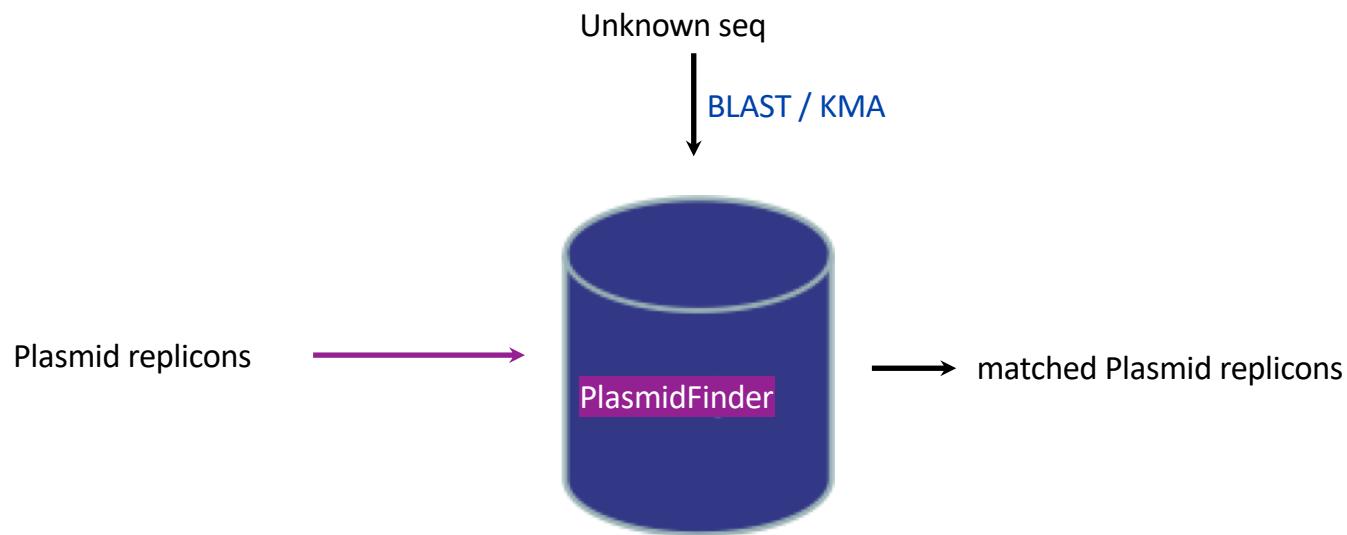
Background

- Plasmids are double-stranded circular or linear DNA molecules. They can replicate and transfer between different bacterial species and clones
- Most of the known plasmids have been identified because they confer phenotypes that are subject to positive selection on bacterial host such as the presence of antimicrobial resistance or virulence genes
- It is important not only to study the molecular epidemiology of different bacterial clones but also to study and understand the molecular epidemiology of transferable plasmids
- For this specific purpose, plasmid typing systems are needed

Hieratical typing strategy



PlasmidFinder



<https://cge.cbs.dtu.dk/services/PlasmidFinder/>

Center for Genomic Epidemiology

Username
Password

Home

Services

Instructions

Output

Article abstract

PlasmidFinder 1.3

PlasmidFinder identifies plasmids in total or partial sequenced isolates of bacteria.

[Use this link to find the link for downloading the database](#)

The input sequence must be in one-letter nucleotide code and the filename must not contain spaces. E.g. correct filename: "test_VTEC:H4.fasta", wrong filename: "test VTEC Ha.fasta"

[Test sequence](#)

The database is curated by:
Henrik Hasman
(click to contact)

View the [version history](#) of this server.

Select database

Select multiple items, with Ctrl-Click (or Cmd-Click on Mac)

Plasmid - Enterobacteriaceae
Plasmid - Gram-positive (under construction)

Select threshold for %ID

95 %

Select type of your reads

Assembled Genome/Contigs*

 Isolate File

Name

Size

Progress

Status

 Upload

 Remove

Center for Genomic Epidemiology

Home

Services

Instructions

Output

PlasmidFinder-1.3 Server - Results

PlasmidFinder Results

SETTINGS:

Selected %ID threshold: **95.00**

PlasmidFinder - Enterobacteriaceae						
Plasmid	%Identity	Query/HSP length	Contig	Position In contig	Note	Accession number
<i>IncP</i>	99.44	535 / 534	strain_1_contig_11	13583..14117	alpha	L27758
<i>IncFII(pRSB107)</i>	97.70	261 / 261	strain_1_contig_18	23064..23324	pRSB107	AJ851089
<i>IncFIB(AP001918)</i>	96.92	682 / 682	strain_1_contig_18	27994..28675		AP001918
<i>IncI1</i>	97.89	142 / 142	strain_1_contig_14	72210..72351	Alpha	AP005147
<i>IncQ1</i>	100.00	450 / 450	strain_1_contig_11	895..1344		HE654726

[extended output](#)

[Results as text](#)

[Results tab separated](#)

[Hit in genome sequences](#)

[Plasmid sequences](#)

Input Files: ***strain_1.fasta.txt***

VirulenceFinder



Bacterial pathogenicity and virulence

- **Pathogenicity.** This is the potential capacity of certain species of microbes to cause an infectious process.
- **Virulence.** signifies the degree of pathogenicity of the given strain. Virulence, therefore, is an index of the qualitative individual nature of the pathogenic microorganism.

E. coli is a **pathogenic** bacterium. However, not all *E. coli* are equally **virulent**. Some have specific virulence factors, which make them more **virulent** than others.

VirulenceFinder in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/virulencemanager/src>

Installation

Setting up VirulenceFinder program

Warning: Due to bugs in BioPython 1.74, if you are not using the Docker container, do not use that version if not using Python 3.7.

```
# Go to wanted location for virulencemanager
cd /path/to/some/dir
# Clone and enter the virulencemanager directory
git clone https://bitbucket.org/genomicepidemiology/virulencemanager.git
cd virulencemanager
```

Build Docker container

```
# Build container
docker build -t virulencemanager .
# Run test
docker run --rm -it \
    --entrypoint=/test/test.sh virulencemanager
```

#Download and install VirulenceFinder database

```
# Go to the directory where you want to store the virulencemanager database
cd /path/to/some/dir
# Clone database from git repository (develop branch)
git clone https://bitbucket.org/genomicepidemiology/virulencemanager_db.git
cd virulencemanager_db
VIRULENCE_DB=$(pwd)
# Install VirulenceFinder database with executable kma_index program
python3 INSTALL.py kma_index
```

If kma_index has no bin install please install kma_index from the kma repository: <https://bitbucket.org/genomicepidemiology/kma>

VirulenceFinder database

https://bitbucket.org/genomicepidemiology/virulencefinder_db/src/master/

Genomic Epidemiology / Databases

virulencefinder_db

Clone

	master	Files	Filter files	
Name		Size	Last commit	Message
 .gitignore		11 B	2018-07-19	Add INSTALL script, gitignore kma binary files
 INSTALL.py		1.87 KB	2018-07-24	Add non_interactive argument
 config		587 B	2017-03-31	Update of config
 listeria.fsa		134.12 KB	2019-09-30	Added new genes to listeria from MyDbFinder db
 notes.txt		20.48 KB	2020-05-29	Updated Ecoli with ExPEC specific genes
 s.aureus_exoenzyme.fsa		38.34 KB	2021-05-06	Change file endings to unix
 s.aureus_hostimm.fsa		49.27 KB	2021-05-06	Change file endings to unix
 s.aureus_toxin.fsa		185.35 KB	2021-05-06	Change file endings to unix
 stx.fsa		177.34 KB	2021-05-06	Change file endings to unix
 virulence_ecoli.fsa		5.85 MB	2021-05-06	Change file endings to unix
 virulence_ent.fsa		189.83 KB	2018-10-01	Fix ccf10:1:AE016830.1 entry

<https://cge.cbs.dtu.dk/services/VirulenceFinder/>

The screenshot shows the homepage of the Center for Genomic Epidemiology. The header features a red background with the text "Center for Genomic Epidemiology". Below the header is a dark grey navigation bar with four links: "Home", "Services", "Publications", and "Contact". The main content area has a white background. At the top of this area, the text "VirulenceFinder 2.0" is displayed. Below it is a horizontal navigation bar with five items: "Service" (which is highlighted in a light grey box), "Instructions", "Output", "Article abstract", "Citations", and "Version history". Underneath this bar, there are two lines of text: "Software version: 2.0.3 (2020-05-21)" and "Database version: (2020-05-29)". To the right of these lines is a yellow callout box containing the text "The database is curated by: Flemming Scheutz, SSI (click to contact)". Further down the page, there are four sections with dropdown menus: "Select species" (with options Listeria, S. aureus, Escherichia coli, Enterococcus), "Select threshold for %ID" (set to 90 %), "Select minimum length" (set to 60 %), and "Select type of your reads" (set to Assembled or Draft Genome/Contigs* (fasta)).

VirulenceFinder 2.0

Service [Instructions](#) [Output](#) [Article abstract](#) [Citations](#) [Version history](#)

Software version: 2.0.3 (2020-05-21)

Database version: (2020-05-29)

The database is curated by:

Flemming Scheutz, SSI

(click to contact)

Select species

Listeria
S. aureus
Escherichia coli
Enterococcus

Select threshold for %ID

90 %



Select minimum length

60 %



Select type of your reads

Only data from one single isolate should be uploaded. If raw sequencing reads are uploaded KMA will be used for mapping. KMA supports the following sequencing platforms: Illumina, Ion Torrent, Roche 454, SOLiD, Oxford Nanopore, and PacBio.

Assembled or Draft Genome/Contigs* (fasta)



VirulenceFinder-1.2 Server - Results

SETTINGS:

Selected %ID threshold: **98.00**

Virulence - <i>E. coli</i>						
Virulence factor	%Identity	Query/HSP length	Contig	Position In contig	Protein function	Accession number
<i>mcmA</i>	99.64	279 / 279	NODE_17_length_48340_cov_62.616714	40909..41187	Microcin M part of colicin H	AJ515251
<i>lpfA</i>	100.00	573 / 573	NODE_4_length_115337_cov_62.053581	84857..85429	Long polar fimbriae	KC207123
<i>iss</i>	99.71	342 / 342	NODE_195_length_89121_cov_54.610832	87701..88042	Increased serum survival	CU928160
<i>prfB</i>	100.00	882 / 882	NODE_75_length_157387_cov_57.585850	94324..95205	P-related fimbriae regulatory gene	CP002970

[extended output](#)

[Results as text](#)

[Results tab separated](#)

[Hit in genome sequences](#)

[Virulence gene sequences](#)

Input Files: *EC19_2011_70_34_3-illumina_pe_velvet1.1.04_kmer63_cov57_cut0.fna*

- **Non-pathogenic:** Few or no obvious virulence factors (you already have information regarding which this could be).
- **UPEC:** Adhesion factors needed to avoid being flushed away from the urinary tract. Also siderophore proteins such as the one encoded by the iroN gene for iron chelation in urine can be relevant due to the iron-limiting environment in the bladder. Presence Pap (P) fimbriae (papG adhesion) can be a sign of increased virulence, as these fimbriae are associated with progression of a urinary tract infection into pyelonephritis (Dan: Nyrebændelse).
- **VTEC, STEC (EHEC):** A hall-mark virulence factor is the Shiga toxin. This is encoded by the stx2A and the stx2B genes.
- **EAEC:** Many different virulence factors (especially aggregative adherence fimbriae (AAFs) located on a 100-kb pAA plasmid, mycolycins such as those encoded by the pic gene and toxins such as those encoded by the pet and the astA gene) are believed to be responsible for the EAEC phenotype. However, it has recently been found that the regulator encoded by the aggR gene also located on the pAA plasmid is coordinating the virulence factors. Therefore, detection of the aggR gene is a good marker for EAEC.
- **ETEC:** This pathotype is known for its production of Heat-Stable Toxin (ST) or Heat-Labile Toxin (LT). The former can be encoded by the sta1 or stb genes and the latter by the elt or ltcA genes.

EPEC	Enteropathogenic E.coli	pEAF	ehly					
ETEC	Enterotoxigenic E.coli	elt	ltcA	sta1	stb	cfa	elt	est
EAEC	Enteroaggregative E.coli	aaf	aggR	aaiC	aatA			
VTEC, STEC (EHEC)	Shiga toxin-producing E.coli	stx2A	stx2B					
EIEC	Enteroinvasive E.coli	pinV						

SalmonellaTypeFinder

Accept only FASTQ

SalmonellaTypeFinder in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/salmonellatypewriter/src/master/>

Installation using Docker

Build time: Docker image takes approximately 1 hour to build.

Setting up SalmonellaTypeFinder

```
# Go to wanted location for SalmonellaTypeFinder
cd /path/to/some/dir
# Clone and enter the salmonellatypewriter directory
git clone https://bitbucket.org/genomicepidemiology/salmonellatypewriter.git
cd salmonellatypewriter
```

The installation can either be with MLST typing capability or without. Choose without, if you intend to provide SalmonellaTypeFinder with an MLST type each time you run it, or if you have the CGE MLST tool installed, and intend to provide the path to it when invoking SalmonellaTypeFinder. Choose with, if you sometimes or always need SalmonellaTypeFinder to find the MLST type.

Build Docker container without MLST typing capability

```
# Rename Docker files
mv Dockerfile Dockerfile_complete
mv Dockerfile_light Dockerfile
# Build container
docker build -t salmonellatypewriter .
```

Build Docker container with MLST typing capability Dependencies: KMA. If kma and kma_index has not bin installed please install kma_index from the kma repository: <https://bitbucket.org/genomicepidemiology/kma>

```
# Build container
docker build -t salmonellatypewriter .
# Go to wanted location for MLST database
# Note: The location needs to be able to be attached to your docker container.
cd /path/to/db_dir
# Clone database from git repository
git clone https://bitbucket.org/genomicepidemiology/mlst_db.git
# Install MLST database with executable kma_index program
cd mlst_db
python3 INSTALL.py kma_index non_interactive
```

SalmonellaTypeFinder database

https://bitbucket.org/genomicepidemiology/salmonellatypefinder_db/src/master/

Genomic Epidemiology / Databases

salmonellatypefinder_db

Clone

master ▾ Files ▾ Filter files

📁 /

Name	Size	Last commit	Message
.DS_Store	6 KB	2017-11-30	all databases added
README.md	118 B	2019-02-20	README.md edited online with Bitbucket
db.json	130.4 KB	2017-11-30	all databases added
scheme_species_map.tab	4.23 KB	2017-11-30	all databases added
senterica.fsa	2.1 MB	2017-11-30	all databases added
senterica.txt.clean	90.14 KB	2017-11-30	all databases added

SerotypeFinder

For serotyping of *E.coli*

SerotypeFinder in CGE bitbucket

<https://bitbucket.org/genomicepidemiology/serotypefinder/src/master/>

Installation

Setting up SerotypeFinder program

Warning: Due to bugs in the BioPython 1.74, if not using the Docker container, do not use that version if not using Python 3.7.

```
# Go to wanted location for serotypefinder
cd /path/to/some/dir
# Clone and enter the serotypefinder directory
git clone https://bitbucket.org/genomicepidemiology/serotypefinder.git
cd serotypefinder
```

Build Docker container

```
# Build container
docker build -t serotypefinder .
```

#Download and install SerotypeFinder database

```
# Go to the directory where you want to store the serotypefinder database
cd /path/to/some/dir
# Clone database from git repository (develop branch)
git clone https://bitbucket.org/genomicepidemiology/serotypefinder_db.git
cd serotypefinder_db
STFINDER_DB=$(pwd)
# Install SerotypeFinder database with executable kma_index program
python3 INSTALL.py kma_index
```

If kma_index has no bin install please install kma_index from the kma repository: <https://bitbucket.org/genomicepidemiology/kma>

SerotypeFinder database

https://bitbucket.org/genomicepidemiology/serotypefinder_db/src/master/

Genomic Epidemiology / Databases

serotypefinder_db

Clone

master ▾

Files ▾

Filter files



/

Name

Size

Last commit

Message

.gitignore

11 B

2019-02-27

Add .gitignore file

H_type.fsa

163.4 KB

2017-01-23

Initial commit

INSTALL.py

1.67 KB

2018-08-03

Added database install script for serotypefinder

O_type.fsa

551.5 KB

2020-09-24

deleted KY115227 and fixed EU549863

config

282 B

2017-01-23

Initial commit

notes.txt

0 B

2017-01-23

notes added