



Universiteit
Leiden
The Netherlands

Metagenomics : beyond the horizon of current implementations and methods

Khachatryan, L.

Citation

Khachatryan, L. (2020, April 28). *Metagenomics : beyond the horizon of current implementations and methods*. Retrieved from <https://hdl.handle.net/1887/87513>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/87513>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/87513> holds various files of this Leiden University dissertation.

Author: Khachatryan, L.

Title: Metagenomics : beyond the horizon of current implementations and methods

Issue Date: 2020-04-28

BacTag - a pipeline for fast and accurate gene and allele typing in bacterial sequencing data

L. Khachatryan¹, M. E. M. Kraakman⁴, A. T. Bernards⁴,
and J. F. J. Laros^{1,2,3}

1 Department of Human Genetics, Leiden University Medical Center, Leiden,
The Netherlands

2 Clinical Genetics, Leiden University Medical Center, Leiden, The Netherlands

3 GenomeScan, Leiden, The Netherlands

4 Department of Medical Microbiology, Leiden University Medical Center, Leiden,
The Netherlands

BMC Genomics, 2019 20:338 doi 10.1186/s12864-019-5723-0

5.1 Background

In order to understand and predict the pathogenic impact and the outbreak potential of a bacterial infection, knowing the species responsible for this infection is not sufficient. Bacterial virulence is often controlled on the sub-species level by the set of specific genes or sometimes even alleles, leading to the necessity of diverse treatment strategies for infections induced by the same bacterial species [301, 302, 303, 304, 305]. For example, antibiotic resistance is one of the most well-known examples where slight variations in a gene can lead to a vast collection of antibiotics resistance profiles within one taxonomic group [306, 307]. Furthermore, different alleles of the same gene can be responsible for distinct adhesion and invasion strategies, reactions to the immune response of the infected organism and toxin production [308, 309]. Besides its relevance for understanding virulence, finding the alleles of specific genes also contributes to a more accurate bacterial classification. One of the most popular methods for subspecies bacterial typing, MultiLocus Sequence Typing (MLST), is based on determination of the alleles of multiple housekeeping genes [310, 311]. Knowing the allele combination allows to identify so called Sequencing Type (ST) of the organism, which is often associated with the important pathogens' attributes such as infection potential [312, 313, 314] or the ability to cause disease in human by transmitting from their animal reservoirs [315, 316, 317]. MLST typing is crucial for the epidemiological studies as it provides fast and accurate identification of geographical dispersal of pathogens and even reveals the migration patterns of the host organism [318, 319].

Despite the importance of the gene and allele typing in the bacterial genomes, there is no "gold standard" method to perform it. For a long time, the presence of particular virulent genes was detected using phenotypic markers such as serotyping [320]. Unfortunately, the set of genetic features that can be revealed using only the phenotype is very limited. Among other restrictions of this group of methods are the inability to grow certain fastidious pathogens in laboratory conditions as well as the extensive delay in cultivation and identification for slowly growing pathogens [321, 322, 323, 166, 324]. In particular cases, the gene and allele identification problem can be solved by using PCR or microarrays with gene- and allele specific primers or probes [325, 326, 327]. These types of methods are much faster and more reliable in comparison to the phenotype-based approaches. However, for the vast majority of genes it is impossible to generate primers or probes that would perform the allele discrimination due to the high similarity among sequences of alleles. Thus, PCR based typing often needs additional analysis, for example, a restriction fragment length polymorphism typing [328, 329] which elaborates the analysis process. PCR-based gene and allele typing most of the time has to be "tailor-made" for the particular group of organisms and the gene of interest. The rapid growth of newly discovered bacteria together with the high mutation rate of some

genes causes the necessity of constant changes in the existing PCR-protocols.

With the improvement of high throughput sequencing techniques and the development of associated bioinformatics software, it became possible to identify the allele variations directly from Whole Shotgun Genome Sequencing (WGS) data by comparing sequencing reads to the reference sequences of the known alleles of the gene of interest in the curated database. Currently, most of the curated and publicly available databases suitable for the gene typing are designed for subspecies classification using the MLST principle. These databases contain variable alleles of housekeeping genes and MLST schemas, associated with those housekeeping genes, for more than 60 bacterial species [330]. There are several tools that perform MLST by aligning assembled WGS data to each sequence in the linked database and reporting the alleles of housekeeping genes with the highest similarity to the provided data [331, 332]. The most recent tools for automated MLST performs the analysis on raw WGS data, as the assembly step is included in its pipeline ([333, 334]). Finally, stringMLST software [335] performs allele identification by comparing the k -mer profiles of raw sequencing data to the k -mer profiles of sequences in the MLST database. This strategy allows to speed up the analysis process drastically, yet the accuracy of the method is lower in comparison with alignment-based ones [336].

Though the WGS-based methods for gene and allele typing potentially requires less effort than any laboratory technique, it has some disadvantages and room for improvement. First of all, the time-consuming separate alignment of WGS data to each sequence in the database can be substituted with a faster algorithm. Furthermore, most of the existing bioinformatics tools for MLST do not provide an option to optimize the analysis settings, which means that the user cannot control, for example, parameters of reads mapping. Finally, it is also not possible to perform the analysis using a database or MLST schema that is not associated with the tool.

In this paper we present BacTag (**B**acterial **T**yping of **a**lles and **g**enes) - a new pipeline, designed to rapidly and accurately detect genes and alleles in sequencing data. Due to the database preprocessing prior to the analysis, BacTag providing a solid and more detailed basis for downstream in comparison with similar tools while retaining the same accuracy. Additionally, our method performs gene and allele detection slightly faster than its current analogs. Our pipeline was successfully tested on both artificial (*E. coli*, *S. pseudintermedius*, *P. gingivalis*, *M. bovis*, *Borrelia spp.* and *Streptomyces spp.*) data and real (*E. coli*, *K. pneumoniae*) clinical WGS samples, by preprocessing the corresponding MLST databases and by performing the subsequent typing. This method is publicly available at <https://git.lumc.nl/l.khachatryan/BacTag>.

5.2 Materials and Methods

5.2.1 Pipeline implementation

The user interface is implemented in Bash, the processing modules are written in GNU Make. Bash allows for user interaction and files maintenance, while GNU Make makes the pipeline suitable for parallel high-performance computing. The pipeline consists of two parts: database preprocessing and sequencing data analysis. Both parts contain modules that include published tools and the scripts from our Python library. The pairwise sequence alignment is performed by the *aln* command from fastools¹. Artificial paired end Illumina FASTQ formatted reads are created by the *make_fastq* local command of sim-reads². Reads are mapped to a reference sequence with BWA *mem* [105]. Alignment sorting and indexing are performed by SAMtools [291]. Potential PCR duplicates are removed using SAMtools *rmdup* command. The SAMtools *mpileup* utility is used to summarize the coverage of mapped reads on a reference sequence at single base pair resolution. Variant calling is performed by the *call* command of BCFtools [337]. To verify whether the called variants for each allele really correspond to the allele sequence, the *vcf-consensus* command of VCFtools [338] is used. Comparison of two VCF files boils down to reporting the number of variants sites that are not equal for both files. Programming languages and software versions used for pipeline construction can be found in Supplementary Table S1. The user may specify parameters for artificial reads generation (by default read length, insert size and coverage are equal 50 nucleotides, 100 nucleotides and 40 respectively), the BWA *mem* and SAMtools *mpileup* utilities for both database preprocessing and sequencing data analysis parts separately. It is also possible to set the ploidy (by default this is one) of the sequencing data, which will be considered during the variants calling in the analysis part of the pipeline.

5.2.1.1 Database preprocessing

The database preprocessing workflow is shown in Figure 5.1. We designed the pipeline such that all independent processes are performed in parallel, which reduces the calculation time.

The user provides the database that consists of alleles grouped by genes of interest. Optionally, the user can provide the 5'- and 3'-flanking regions for each gene, otherwise, every allele will be flanked on both sides with a fifty-nucleotide long poly-N sequence. That is done in order to prevent the coverage drop at the end of sequence during the sequencing data mapping. In the first step of the preprocessing stage, the sequences of all alleles belonging to the same gene are aligned in a pairwise manner, yielding the Levenshtein [339] distance for each pair of alleles.

¹ Available from: <https://git.lumc.nl/j.f.j.laros/fastools>. Accessed 27 Oct 2018.

² Available online at <https://git.lumc.nl/j.f.j.laros/sim-reads>

These distances are used to select the allele with the smallest average distance to all other sequences as the gene reference. In the same step the quality of the provided database is checked: it is reported when the same sequence is provided for multiple alleles or when one allele sequence is a subsequence of another. Once the quality report is created, the user can fix the original database when needed. In the next step, artificial Illumina paired end reads are created based on the sequence of each allele.

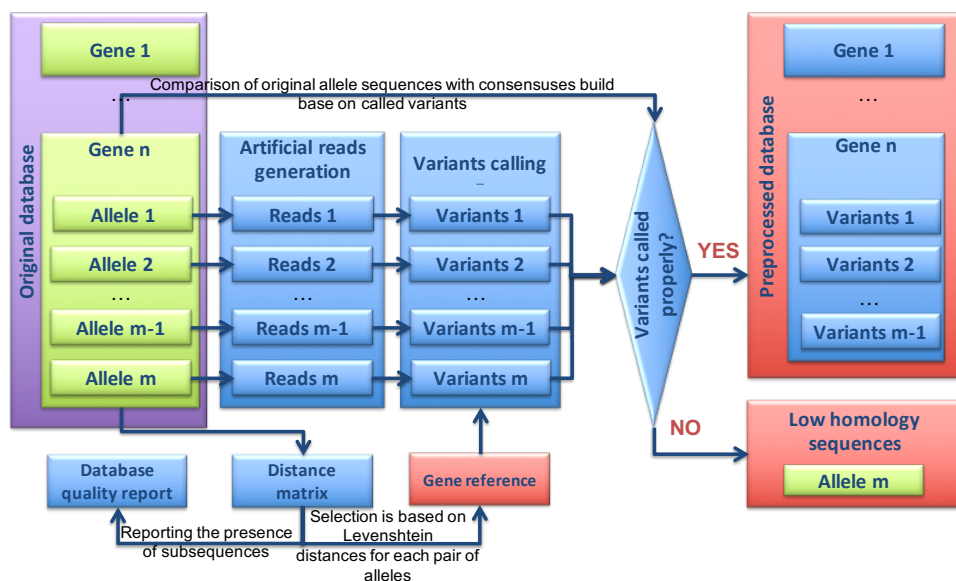


Figure 5.1: Schematic representation of the database preprocessing. All of the processes are illustrated for one gene. Calculations for several genes are done independently in parallel.

Reads are mapped to the selected gene reference, the alignment map file is sorted and indexed, after which the coverage of mapped reads on the reference sequence at a single base pair resolution is summarized and stored in a BCF file, which is used for variants calling. Variants are stored in a VCF file and further subjected to a quality check to verify whether they really correspond to the allele sequence. If variants defining alleles' sequence were not properly called, allele is reported and assigned to the so-called low similarity group of sequences. The low similarity group contains sequences for which the variants were not called correctly during the database preprocessing when using the centroid reference. I.e., for these alleles, the centroid is not an appropriate reference and therefore these sequences should be considered to be references themselves. In the final step the references of all genes are concatenated into one FASTA file, which further serves as the database reference.

5.2.1.2 Sequencing data analysis

The data analysis workflow can be found in Figure 5.2. To initiate the analysis,

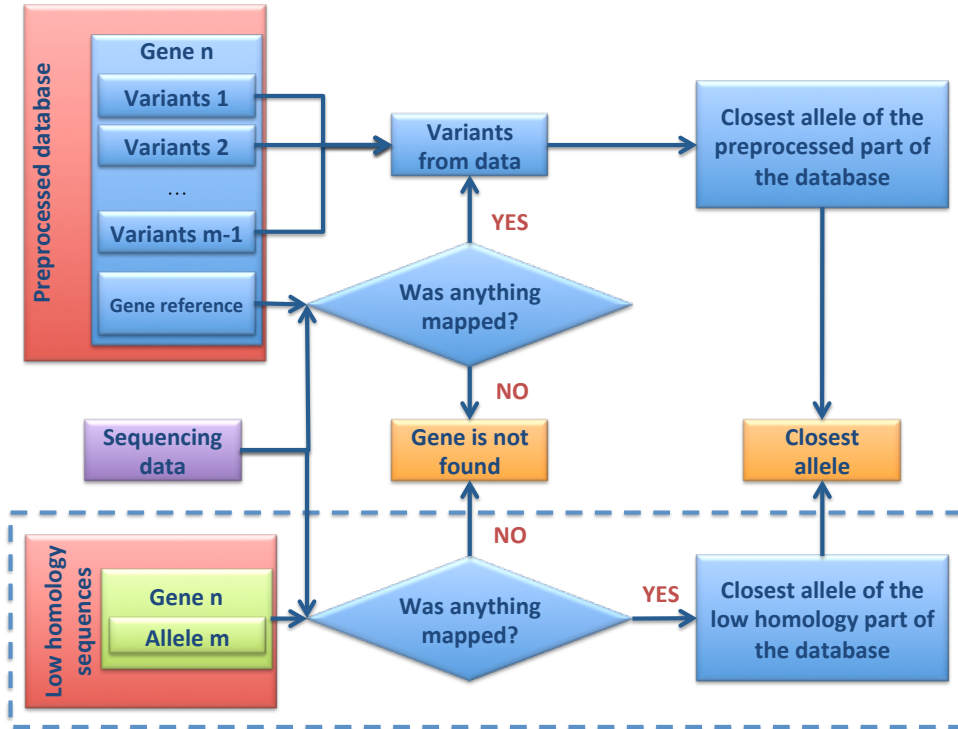


Figure 5.2: Schematic representation of the analysis part of BacTag pipeline. All of the processes are illustrated for one gene. Calculations for multiple genes are done independently in parallel. The analysis of the low homology group of sequences is highlighted by the dashed box and can be manually turned off by the user for the time efficiency.

the user provides two paired FASTQ files. After analysis initialization an output directory is created, which will serve to store the results of the analysis. The user can choose the name of the output directory, otherwise it will have the same name as the basename of the provided FASTQ files. The sequencing data analysis part of the pipeline is comprised of two steps: the main analysis and the analysis of low similarity group of sequences. If no sequences were assigned to the low similarity group during the database preprocessing, only the first step will be performed. The user can manually turn off the second step for time efficiency.

The main analysis

This part of the pipeline applies to the alleles that were not placed in the low ho-

mology group of sequences during the database preprocessing. analysed reads are mapped to the database reference, obtained after database preprocessing by concatenating all the gene reference sequences. The alignment map file is indexed and sorted and substituted to the removal of potential PCR duplicates. If there are no reads mapped to the gene reference, the gene is reported as not found in the analysed dataset. Otherwise, mapped reads are used to estimate the horizontal coverage of a gene reference at base pair resolution. The obtained BCF coverage summary is used for variant calling, the result of which is stored in VCF format. Variants are compared with variants collected for each gene allele during the preprocessing phase. Once the comparisons are done, the allele with the least difference from the sequencing data will be reported. It is also reported, if heterozygous variants were found in the sample, as that might indicate sequencing or mapping problems as well as the presence of more than one gene allele in the sequencing data. Reports for all genes are concatenated to a single result file, which is placed in the output directory.

Low homology group of sequences analysis

This part of the pipeline works with alleles that were placed in the low homology group of sequences during the database preprocessing. Sequencing reads are subjected to variant calling using each of the alleles from the low homology group as a reference (the same routine with the same parameters as for the main analysis step). If for the particular gene one of the alleles from the low homology group has fewer differences with the sequencing data in comparison to the allele reported during the main analysis, the allele from the low homology group will be reported as present in the sequencing data.

5.2.2 Pipeline testing

All the computational benchmarking was done on chimerashark Blade Server of SHARK computer cluster³ with the maximum of 24 CPUs used at the same time.

5.2.3 Database

5.2.3.1 Genes and alleles

The database preprocessing part of the pipeline was tested using seven curated databases: *E. coli* Achtman MLST⁴ (downloaded January 2018), *K. pneumoniae* Pasteur MLST⁵ (downloaded October 2018), *S. pseudintermedius* MLST⁶ (downloaded

³<https://git.lumc.nl/shark/SHARK/wikis/home>

⁴https://enterobase.warwick.ac.uk/species/ecoli/download_7_gene

⁵<https://bigsd.pasteur.fr/klebsiella/>

⁶<https://pubmlst.org/spseudintermedius/>

February 2019), *P. gingivalis* MLST⁷ (downloaded February 2019), *M. bovis* MLST⁸ (downloaded February 2019), *Borrelia spp.* MLST⁹ ([downloaded February 2019) and *Streptomyces spp.* MLST¹⁰ ([downloaded February 2019). Each database contains sequences of variable regions of housekeeping genes: five for the *Streptomyces spp.* MLST, eight for the *Borrelia spp.* MLST and seven for all the remaining schemas.(see Table 5.1).

MLST schemas were selected for organisms from six different bacterial phyla. These organisms have a GC-content ranging between 29 and 73%. For the database preprocessing the following parameters for BWA mem and SAMtools mpileup tools were selected. Since the database consists of sequences of highly variable regions of housekeeping genes, the alignment mismatch penalty was set to 2 (4 by default) in order to provide the proper alignment for the regions where variants occur in close proximity. The minimum seed length was changed to 15 (19 by default) due to the short length of sequences in the selected database. Penalty for 5'- and 3'-end clipping was set to 100 (5 by default), forcing alignment to detect the variants located at the ends of the variable region. Single end mapped reads (anomalous read pairs, -A) were counted in order to detect variants located at the ends of the variable region. BAQ computation was disabled, as it is oversensitive to regions densely populated with variants. Bases with baseQ/BAQ lower than 13 were not skipped, since the database preprocessing is based on high quality artificial sequencing reads.

5.2.3.2 Flanking regions

The sequences of polymerase chain reaction (PCR) primers commonly applied to amplify each of the housekeeping genes (*E. coli* [340], *K. pneumoniae*¹¹, *S. pseudintermedius*¹², *M. bovis*¹³, *P. gingivalis* [341]) for the selected MLST schemas were used to construct the flanking regions for this study. Each flanking region includes the primer sequence as well as the genomic sequence between the primer and the variable region of interest. The genomic sequence is extracted from the genome of one of the target strains for the corresponding MLST schema (see Table 5.1). In case low-sensitivity PCR primers are used (e.g., for *Borrelia spp.* MLST) or if no PCR primer sequences are available (e.g., for *Streptomyces spp.* MLST), fifty nucleotides before and after the variable regions were used as flanks. Flanking regions have the same orientation as the allele sequences in the database (see section 5.7.3, Additional file 2: Tables S2-S8).

⁷<https://pubmlst.org/pgingivalis/>

⁸<https://pubmlst.org/mbovis/>

⁹<https://pubmlst.org/borrelia/>

¹⁰<https://pubmlst.org/streptomyces/>

¹¹ Available from: http://bigsd.b.pasteur.fr/klebsiella/primers_used.html. Accessed 16 Oct 2018.

¹² Available from: <https://pubmlst.org/spseudintermedius/info/primers.pdf>. Accessed 16 Feb 2019.

¹³ Available from https://pubmlst.org/mbovis/info/M._bovis_MLST_targets_and_primers.pdf. Accessed 16 Feb 2019.

MLST database	Genes including number of alleles per gene	Number of alleles (per gene) in the low similarity group	Strain and reference sequence used for flanking region construction
<i>E. coli</i>	<i>adk</i> (623), <i>fumC</i> (933), <i>gyrB</i> (606), <i>Icd</i> (823), <i>mdh</i> (614), <i>purA</i> (563), <i>recA</i> (512)	<i>fumC</i> (11), <i>gyrB</i> (3), <i>mdh</i> (8)	UMN026, NC_011751.1
<i>K. pneumoniae</i>	<i>gapA</i> (184), <i>infB</i> (141), <i>mdh</i> (245), <i>pgi</i> (221), <i>phoE</i> (365), <i>rpoB</i> (189), <i>tonB</i> (472)	<i>gapA</i> (6), <i>mdh</i> (3), <i>tonB</i> (29)	Kp52.145, FO834906.1
<i>S. pseudintermedius</i>	<i>ack</i> (46), <i>cpn60</i> (96), <i>fdh</i> (26), <i>pta</i> (70), <i>purA</i> (77), <i>sar</i> (38), <i>tuf</i> (24)	-	ED99, NC_017568.1
<i>M. bovis</i>	<i>adh1</i> (15), <i>gltX</i> (17), <i>gpsA</i> (14), <i>gyrB</i> (25), <i>pta2</i> (23), <i>tdk</i> (15), <i>tkl</i> (26)	-	PG45, NC_014760.1
<i>P. gingivalis</i>	<i>ftsQ</i> (40), <i>gpdxJ</i> (37), <i>hagB</i> (37), <i>mcmA</i> (30), <i>pepO</i> (37), <i>pga</i> (27), <i>recA</i> (14)	-	ATCC 33277, NC_010729.1
<i>Borrelia</i> spp.	<i>clpA</i> (296), <i>clpX</i> (258), <i>nifS</i> (230), <i>pepX</i> (261), <i>pyrG</i> (269), <i>recG</i> (285), <i>rplB</i> (250), <i>uvrA</i> (261)	<i>clpA</i> (58), <i>clpX</i> (51), <i>nifS</i> (54), <i>pepX</i> (57), <i>pyrG</i> (51), <i>recG</i> (55), <i>rplB</i> (54), <i>uvrA</i> (45)	<i>B. hermsii</i> DAH, NC_010673.1
<i>Streptomyces</i> spp.	<i>atpD</i> (183), <i>gyrB</i> (179), <i>recA</i> (184), <i>rpoB</i> (183), <i>trpB</i> (200)	<i>atpD</i> (72), <i>gyrB</i> (147), <i>recA</i> (2), <i>rpoB</i> (6), <i>trpB</i> (69)	<i>S. coelicolor</i> A3(2), NC_003888.3

Table 5.1: Preprocessed MLST databases

5.2.3.3 Artificial test data

The sequencing data analysis part of the pipeline was validated by using artificial Illumina reads, based on the complete genomes of 30 different bacterial strains belonging to 13 different bacterial species (see Table 5.2), for which the alleles of housekeeping genes associated with the corresponding MLST schema were previously reported. Paired end FASTQ formatted reads of 100 bp were generated with an insert size of 100. For each genome, an average coverage of 80 was generated in this way.

5.2.3.4 Real test data

The analysis part of the pipeline was tested on 185 paired end Illumina WGS samples belonging to 9 different previously reported sequencing types (ST) of *E. coli* (section 5.7.3, Additional file 3: Table S9) and 98 paired end Illumina WGS samples belonging to 43 different previously reported STs of *K. pneumoniae* (section 5.7.3, Additional file 3: Table S10). Sequencing reads were downloaded from Sequence Read Archive (SRA, [342]). Prior to the analysis, the data quality check and correction (when necessary) was done for each sample using Flexiprep QC pipeline¹⁴.

5.2.3.5 Parameters used for sequencing data analysis

The analysis of both artificial and real samples was done with the same parameters of BWA *mem* as during the database preprocessing. SAMtools *mpileup* parameters were as follow: anomalous read pairs were counted; extended BAQs were calculated for higher sensitivity but lower specificity.

¹⁴Available online at <http://biopet-docs.readthedocs.io/en/latest/pipelines/flexiprep/>

5.3 Results

5.3.1 Building the preprocessed MLST databases

We used BacTag to preprocess seven publicly available MLST databases. During this process we did not detect any duplications or partial sequences for any of the preprocessed databases. When preprocessing *E. coli* Achtman seven genes MLST database, 22 sequences (less than 0.5% of the total number of analysed sequences) belonging to three different genes were assigned to the low similarity group of sequences (see Table 5.1). The run time of the *E. coli* database preprocessing was approximately 2h. The peak memory usage was 150Mb. During the preprocessing of the *K. pneumoniae* database associated with the Pasteur seven genes MLST schema, 38 sequences (2.1% of the total number of analysed sequences) belonging to three different genes were assigned to the low similarity group of sequences. Preprocessing of databases associated with MLST schemas for *S. pseudintermedius*, *M. bovis* and *P. gingivalis* reported no sequences placed in the low similarity group of sequences. For the databases associated with the MLST schemas for *Borrelia* spp. and *Streptomyces* spp. 425 sequences (19.2% of the total number of analysed sequences) and 296 sequences (31.8% of the total number of analysed sequences) were placed in the low similarity group respectively. This large number of low similarity sequences indicates that the alleles in the analysed MLST databases are quite heterogeneous, which can be expected, considering that both aforementioned MLST schemas are genus-specific, not species-specific like other five analysed databases.

Since distance matrix is computed during the preprocessing, the expected CPU time will scale quadratically with the size of the database. We indeed found this behaviour as shown in Figure 5.3.

5.3.2 Testing BacTag on artificial data

We used the preprocessed MLST databases to reveal the presence of the corresponding housekeeping genes and to predict the allele for each of these genes in artificial sequencing data based on complete genomes of 30 different bacterial strains belonging to 15 different species. All housekeeping genes associated with the corresponding MLST schema were identified in each sample. The alleles found by the pipeline matched with the previously reported ones for each but one of the analysed genomes (see Table 5.2). The genome of *P. gingivalis* AJW4 (GenBank accession number NZ_CP011996.1) was previously reported [343] to have the allelic variants *ftsQ*-16, *gpdX*J-9, *hagB*-22, *mcmA*-17, *pepO*-22, *pga*-15 and *recA*-1. However, BacTag analysis revealed the following set of alleles: *ftsQ*-21, *gpdX*J-23, *hagB*-1, *mcmA*-3, *pepO*-20, *pga*-3 and *recA*-7. Manual inspection confirmed that alleles reported by BacTag are correct in case of all aforementioned genes.

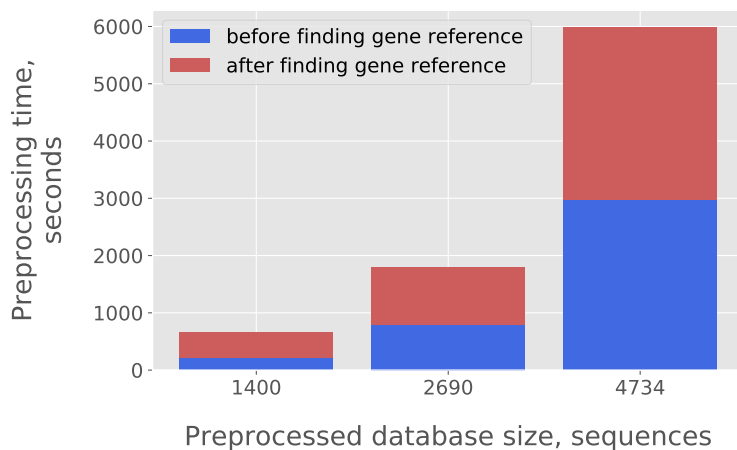


Figure 5.3: The dependence of database preprocessing time from the amount of sequences in the database.

5.3.3 Testing BacTag on real *E. coli* and *K. pneumoniae* data

We tested BacTag on 185 *E. coli* and 97 *K. pneumoniae* clinical publicly accessible WGS datasets, with each test yielding either one of nine *E. coli* or one of 44 *K. pneumoniae* sequencing types (STs). *E. coli* samples were analysed using the preprocessed *E. coli* Achtman seven genes MLST database, while *K. pneumoniae* samples were analysed using the preprocessed *K. pneumoniae* Pasteur seven genes MLST database. Each sample was shown to contain all expected seven housekeeping genes; alleles of those genes identified using our method corresponded to the expected ones for all but one sample (see Table 5.3). This sample was checked additionally using web-based tools for the MLST [333, 334]. Results of this independent check were completely identical to the ones obtained by our pipeline and suggest that the sample belongs to *E. coli* ST95 instead of ST73. Furthermore, according to the original publication [344], MLST was never done for this and 21 other samples analysed during the same study in order to confirm their sequencing type. Thus, we conclude that in Ref. [344] one of the samples was incorrectly assigned to *E. coli* ST73.

Our pipeline reported the presence of multiple variants at the same position for eight *E. coli* samples belonging to three different STs and 55 samples of *K. pneumoniae* belonging to 24 different STs (see Table 5.3)). This might suggest the presence of contamination in the sequenced DNA samples or the existence of pseudogenes in the genome of the sampled organisms.

Species and strain	GeneBank AC	Identified alleles
<i>E. coli</i> 042	FN554766.1	<i>adk</i> -18, <i>fumC</i> -22, <i>gyrB</i> -20, <i>lcd</i> -23, <i>mdh</i> -5, <i>purA</i> -15, <i>recA</i> -4
<i>E. coli</i> E2348/69	FM180568.1	<i>adk</i> -15, <i>fumC</i> -15, <i>gyrB</i> -11, <i>lcd</i> -15, <i>mdh</i> -18, <i>purA</i> -11, <i>recA</i> -11
<i>E. coli</i> E24377A	CP000800.1	<i>adk</i> -6, <i>fumC</i> -213, <i>gyrB</i> -33, <i>lcd</i> -1, <i>mdh</i> -24, <i>purA</i> -8, <i>recA</i> -7
<i>E. coli</i> IHE3034	NC_017628.1	<i>adk</i> -37, <i>fumC</i> -38, <i>gyrB</i> -19, <i>lcd</i> -37, <i>mdh</i> -17, <i>purA</i> -11, <i>recA</i> -26
<i>E. coli</i> IMT5155	CP005930.1	<i>adk</i> -55, <i>fumC</i> -38, <i>gyrB</i> -19, <i>lcd</i> -37, <i>mdh</i> -17, <i>purA</i> -11, <i>recA</i> -26
<i>E. coli</i> RS218	NZ_CP007149.1	<i>adk</i> -37, <i>fumC</i> -38, <i>gyrB</i> -19, <i>lcd</i> -37, <i>mdh</i> -17, <i>purA</i> -11, <i>recA</i> -26
<i>E. coli</i> UMN026	NC_011751.1	<i>adk</i> -21, <i>fumC</i> -35, <i>gyrB</i> -115, <i>lcd</i> -6, <i>mdh</i> -5, <i>purA</i> -5, <i>recA</i> -4
<i>S. pseudintermedius</i> NA45	NZ_CP016072.1	<i>ack</i> -2, <i>cpn60</i> -10, <i>fdh</i> -2, <i>pta</i> -1, <i>purA</i> -5, <i>sar</i> -1, <i>tuf</i> -2
<i>S. pseudintermedius</i> ED99	NC_017568.1	<i>ack</i> -3, <i>cpn60</i> -9, <i>fdh</i> -2, <i>pta</i> -1, <i>purA</i> -1, <i>sar</i> -1, <i>tuf</i> -1
<i>S. pseudintermedius</i> HKU10-03	NC_014925.1	<i>ack</i> -2, <i>cpn60</i> -55, <i>fdh</i> -3, <i>pta</i> -42, <i>purA</i> -14, <i>sar</i> -2, <i>tuf</i> -1
<i>M. bovis</i> Ningxia-1	NZ_CP023663.1	<i>adh</i> 1-4, <i>gltX</i> -3, <i>gpsA</i> -2, <i>gyr</i> -3, <i>pta</i> 2-17, <i>tdk</i> -3, <i>tk</i> t-4
<i>M. bovis</i> HB0801	NC_018077.1	<i>adh</i> 1-4, <i>gltX</i> -3, <i>gpsA</i> -2, <i>gyr</i> -3, <i>pta</i> 2-5, <i>tdk</i> -3, <i>tk</i> t-4
<i>M. bovis</i> NM2012	NZ_CP011348.1	<i>adh</i> 1-4, <i>gltX</i> -3, <i>gpsA</i> -2, <i>gyr</i> -3, <i>pta</i> 2-5, <i>tdk</i> -3, <i>tk</i> t-4
<i>M. bovis</i> CQ-W70	NC_015725.1	<i>adh</i> 1-4, <i>gltX</i> -5, <i>gpsA</i> -2, <i>gyr</i> -3, <i>pta</i> 2-5, <i>tdk</i> -3, <i>tk</i> t-4
<i>M. bovis</i> PG45	NC_014760.1	<i>adh</i> 1-3, <i>gltX</i> -2, <i>gpsA</i> -4, <i>gyr</i> -2, <i>pta</i> 2-1, <i>tdk</i> -3, <i>tk</i> t-2
<i>M. bovis</i> 08M	NZ_CP019639.1	<i>adh</i> 1-4, <i>gltX</i> -3, <i>gpsA</i> -2, <i>gyr</i> -3, <i>pta</i> 2-5, <i>tdk</i> -3, <i>tk</i> t-4
<i>P. gingivalis</i> ATCC 33277	NC_010729.1	<i>ftsQ</i> -5, <i>gpdxJ</i> -9, <i>hagB</i> -1, <i>mcmA</i> -1, <i>pepO</i> -1, <i>pga</i> -5, <i>recA</i> -5
<i>P. gingivalis</i> AJW4	NZ_CP011996.1	<i>ftsQ</i> -21, <i>gpdxJ</i> -23, <i>hagB</i> -1, <i>mcmA</i> -3, <i>pepO</i> -20, <i>pga</i> -3, <i>recA</i> -7
<i>P. gingivalis</i> A7A1-28	CP013131.1	<i>ftsQ</i> -1, <i>gpdxJ</i> -12, <i>hagB</i> -1, <i>mcmA</i> -1, <i>pepO</i> -1, <i>pga</i> -1, <i>recA</i> -1

Table 5.2: To be continued on the next page

Species and strain	GeneBank AC	Identified alleles
<i>Borrelia hermsii</i> DAH	NC_010673.1	<i>clpA</i> -68, <i>clpX</i> -165, <i>nifS</i> -149, <i>pepX</i> -171, <i>pyrG</i> -179, <i>recG</i> -188, <i>rplB</i> -157, <i>worA</i> -175
<i>Borrelia turicatae</i> 91E135	NC_008710.1	<i>clpA</i> -71, <i>clpX</i> -166, <i>nifS</i> -150, <i>pepX</i> -172, <i>pyrG</i> -180, <i>recG</i> -189, <i>rplB</i> -158, <i>worA</i> -176
<i>Borrelia anserina</i> BA2	CP005829	<i>clpA</i> -212, <i>clpX</i> -179, <i>nifS</i> -161, <i>pepX</i> -186, <i>pyrG</i> -196, <i>recG</i> -204, <i>rplB</i> -170, <i>worA</i> -188
<i>Borrelia recurrentis</i> A1	NC_011244	<i>clpA</i> -213, <i>clpX</i> -164, <i>nifS</i> -162, <i>pepX</i> -187, <i>pyrG</i> -197, <i>recG</i> -205, <i>rplB</i> -156, <i>worA</i> -189
<i>Borrelia parkeri</i> SLO	CP005851	<i>clpA</i> -214, <i>clpX</i> -180, <i>nifS</i> -163, <i>pepX</i> -188, <i>pyrG</i> -198, <i>recG</i> -206, <i>rplB</i> -171, <i>worA</i> -190
<i>Borrelia coriaceae</i> Co53	CP005745	<i>clpA</i> -215, <i>clpX</i> -181, <i>nifS</i> -164, <i>pepX</i> -189, <i>pyrG</i> -199, <i>recG</i> -207, <i>rplB</i> -172, <i>worA</i> -191
<i>Borrelia crocidurae</i> Achema	CP003426	<i>clpA</i> -216, <i>clpX</i> -164, <i>nifS</i> -165, <i>pepX</i> -190, <i>pyrG</i> -200, <i>recG</i> -208, <i>rplB</i> -173, <i>worA</i> -192
<i>Streptomyces coelicolor</i> A3(2)	NC_003888.3	<i>atpD</i> -127, <i>gyrB</i> -124, <i>recA</i> -131, <i>rpoB</i> -126, <i>trpB</i> -142
<i>Streptomyces fulvisimilis</i> DSM 40593	CP005080.1	<i>atpD</i> -133, <i>gyrB</i> -130, <i>recA</i> -13, <i>rpoB</i> -36, <i>trpB</i> -147
<i>Streptomyces griseus</i> NBRC 13350	NC_010572.1	<i>atpD</i> -6, <i>gyrB</i> -8, <i>recA</i> -8, <i>rpoB</i> -8, <i>trpB</i> -8
<i>Streptomyces albidoflavus</i> J1074	NC_020990.1	<i>atpD</i> -36, <i>gyrB</i> -5, <i>recA</i> -5, <i>rpoB</i> -36, <i>trpB</i> -39

Table 5.2: Testing the pipeline on artificial WGS data

5.3.4 Comparing BacTag with web-based tools for *E. coli* Achtman MLST

We measured the time required for the analysis, using 30 samples belonging to the ST131 with the dataset size varying from 0.2 to 3. Gb. We performed the MLST typing in two modes: with and without analysis of the low similarity sequences group. As can be seen in Figure 5.4a and b, the processing time of BacTag depended on the sequencing sample size and the analysis mode. The larger the input sequencing data is, the more time is required for typing regardless of the analysis mode. Performing the typing including the analysis of low similarity group (mode 2) increases the processing time. Including low similarity sequences into the analysis did not affect the final output, for all samples tested during this research.

SRA Run AC	Reported ST	Expected ST	Genes with multiple variants at the same position
ERR966604	95	73	-
SRR2767732	16	16	<i>Icd</i>
SRR2767734	21	21	<i>Icd, mdh</i>
SRR2970643	131	131	<i>fumC</i>
SRR2970737	131	131	<i>adk, fumC, gyrB, mdh, recA, purA</i>
SRR2970742	131	131	<i>fumC</i>
SRR2970753	131	131	<i>fumC</i>
SRR2970774	131	131	<i>fumC</i>
SRR2970775	131	131	<i>fumC</i>
SRR5973405	1164	1164	<i>phoE</i>
SRR5973308	1180	1180	<i>phoE</i>
SRR5973303	13	13	<i>phoE</i>
SRR5973253	133	133	<i>phoE</i>
SRR5973334	133	133	<i>phoE</i>
SRR5973324	1373	1373	<i>phoE</i>
SRR5973251	1426	1426	<i>gapA, phoE</i>
SRR5973269	147	147	<i>gapA</i>
SRR5973320	1876	1876	<i>phoE</i>
SRR5973351	188	188	<i>gapA</i>
SRR5973329	20	20	<i>phoE</i>
SRR5973408	2267	2276	<i>phoE</i>
SRR5973397	25	25	<i>phoE</i>
SRR5973248	258	258	<i>gapA</i>
SRR5973283	258	258	<i>gapA</i>
SRR5973279	258	258	<i>gapA</i>
SRR5973271	258	258	<i>gapA</i>
SRR5973336	258	258	<i>gapA</i>
SRR5973319	258	258	<i>gapA</i>
SRR5973317	258	258	<i>gapA</i>
SRR5973294	258	258	<i>gapA</i>
SRR5973291	258	258	<i>gapA</i>
SRR5973289	258	258	<i>gapA</i>
SRR5973400	258	258	<i>gapA</i>
SRR5973382	258	258	<i>gapA</i>
SRR5973381	258	258	<i>gapA</i>
SRR5973287	258	258	<i>gapA</i>
SRR5973240	307	307	<i>phoE</i>

Table 5.3: To be continued on the next page

SRA Run AC	Reported ST	Expected ST	Genes with multiple variants at the same position
SRR597324	307	307	<i>phoE</i>
SRR5973282	307	307	<i>phoE</i>
SRR5973280	307	307	<i>phoE</i>
SRR5973339	307	307	<i>phoE</i>
SRR5973322	307	307	<i>phoE</i>
SRR5973288	307	307	<i>phoE</i>
SRR5973396	307	307	<i>phoE</i>
SRR5973380	307	307	<i>phoE</i>
SRR5973379	307	307	<i>phoE</i>
SRR5973376	307	307	<i>phoE</i>
SRR5973373	307	307	<i>phoE</i>
SRR5973361	307	307	<i>phoE</i>
SRR5973355	307	307	<i>phoE</i>
SRR5973284	23	23	<i>phoE</i>
SRR5973332	35	35	<i>phoE</i>
SRR5973389	35	35	<i>phoE</i>
SRR5973368	35	35	<i>phoE</i>
SRR5973393	405	405	<i>phoE</i>
SRR5973311	412	412	<i>phoE</i>
SRR5973371	429	429	<i>tonB</i>
SRR5973327	466	466	<i>phoE</i>
SRR5973407	466	466	<i>phoE</i>
SRR5973239	492	492	<i>phoE</i>
SRR5973301	502	502	<i>phoE</i>
SRR5973348	753	753	<i>phoE</i>
SRR5973362	8	8	<i>phoE</i>

Table 5.3: Results of pipeline testing on real *E. coli* and *K. pneumoniae* data. Only samples with results different from expected are shown.

The same 30 samples were submitted for analysis to web-based tools for MLST typing: MLST1.8 [333] and Enterobase [334]. These methods perform the assembly of submitted WGS data and use the obtained contigs for the BLAST-based comparison with sequences in the MLST database. For both tools, the results of the WGS assembly can be downloaded after the analysis is finished, MLST 1.8 also provides information about BLAST alignments for the best matching alleles as an output. The analysis of the 30 samples with MLST 1.8 took from 299 to 569 (median 454) minutes per job, the processing time did not correlate with the input data size (Figure 5.4c). MLST 1.8 failed to perform the assembly (and thus to finish the MLST) for two samples.

Long processing time can be explained by high load of the tool server. However, that cannot be checked as it is only possible to track the time in between job submission to the server and the time when job is finished. It is unfortunately not possible to assess when the actual calculations for the particular sample started. Another tool, Enterobase, failed to perform the analysis of one sample (due to the problems with assembly) and did not define the correct ST for one other sample. However, Enterobase shows when each part of the analysing pipeline is being launched, which allowed us to determine the time required for the analysis of each sample and compare it to our tool (Figure 5.5). The processing time for Enterobase was comparable to our tool and also seems to be dependent on the size of the submitted WGS data (Figure 5.4d).

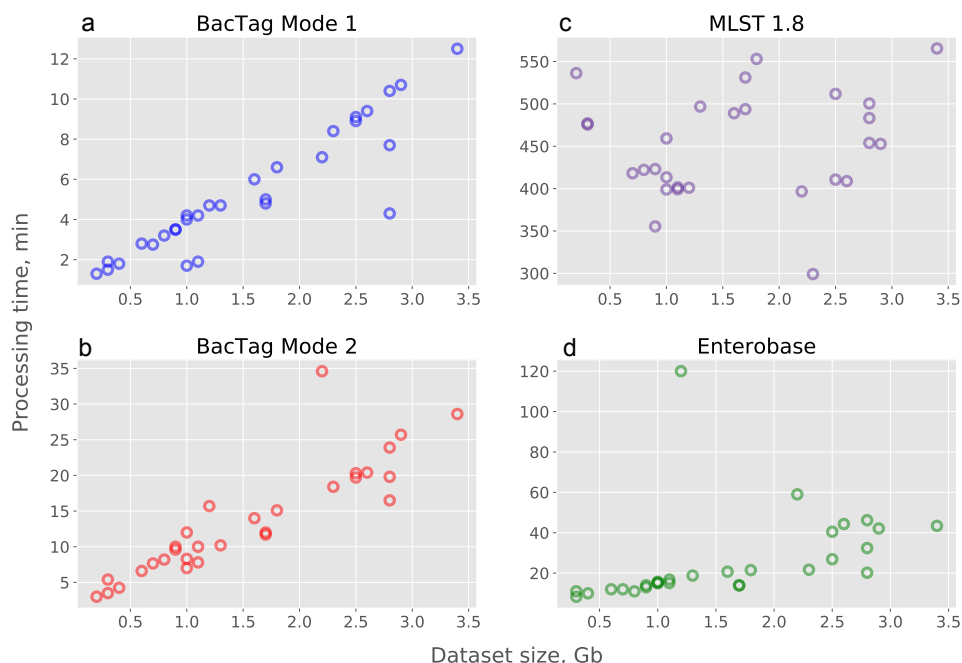


Figure 5.4: Time required for the analysis of 30 samples belonging to the ST131 by two modes of BacTag (a and b), MLST 1.8 (c) and Enterobase (d)

5.4 Discussion

In this paper we described BacTag - a new pipeline designed to perform fast and accurate gene and allele detection directly using WGS data. Our method was shown

to work faster and more accurate than most popular current bioinformatics tools due to the absence of the necessity to compare sequencing data with each sequence in the database. Instead, we preprocess the reference database once prior to the analysis in order to store all the mismatches between different alleles of the same gene. Under the assumption that all alleles of the same gene are highly similar, it is easy to check whether the gene of interest is present in the sequencing data by mapping the reads to the most "average" gene allele. Variants detected after such mapping can be compared with the information obtained during the database preprocessing in order to retrieve the allele of the detected gene. Since the database preprocessing needs to be done only once, this approach significantly reduces the time required for the analysis of multiple samples. Additionally, the possibility of parallel computation allows to speed up the database preprocessing significantly since all of the independent computations can be done in parallel.

Most of the existing tools for automatic gene and allele detection are based on fixed and rarely updated databases. The possibility to choose the database that will be preprocessed as well as to check the quality of that database is another essential feature of BacTag. It is important to note that the pipeline allows the user to set the parameters for the database preprocessing and sequencing data analysis. The same database, preprocessed with different parameters, allows the user to control in which case the variants for some alleles are not properly called. Thus, the user can determine the optimal parameters to detect as many of the alleles of interest as possible and apply this knowledge to the experimental design. On the other hand, preprocessing the database with the parameters of already existing sequencing data provides an estimate of the alleles that likely will not be properly detected.

While the current tools for gene allele identification require assembly of the WGS data prior to the comparison with the reference database, we chose to work directly with raw sequencing data. This was done in order to preserve the information about positions with multiple reported variants, which would be lost in case of bacterial genome assembly. That information is crucial for the detection of possible sample contaminations, presence of pseudogenes and, potentially, for extending our pipeline to metagenomic datasets. Furthermore, BacTag can work with sequencing data that for some reasons cannot be assembled.

Two main limitations of the pipeline need to be addressed. First, our approach assumes that a considerable part of the same gene alleles is highly similar. The more alleles of the same gene that do not fulfill this requirement, the slower the pipeline will work: sequences for which the pipeline will not be able to call the proper variants will be checked by direct read mapping. Second, the pipeline also does not provide proper analysis results if several alleles of the same gene are present in the sequencing data (this can be caused, among other reason, by the mixed-strain infection of the same subject, see [345, 346, 347]). More detailed evaluation of the horizontal coverage of the detected genes as well as the additional analysis of the positions with multiple variants reported could potentially help to resolve this

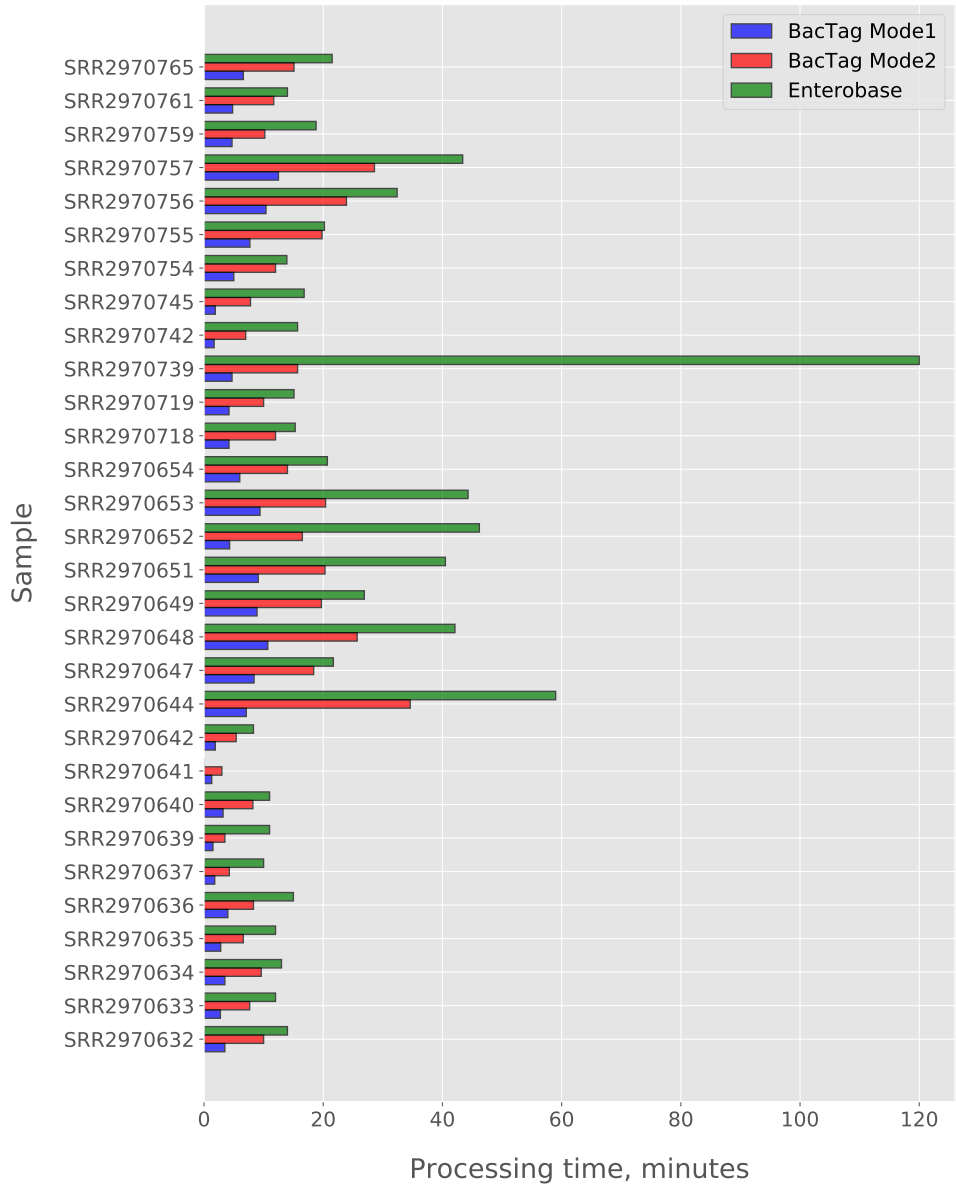


Figure 5.5: Comparing of the processing time required for the Achtman seven genes MLST analysis of 30 WGS *E. coli* samples.

problem and extend the approach in order to perform the analysis on complicated metagenomic datasets.

5.5 Conclusions

We have introduced BacTag - a new pipeline for fast and accurate gene and allele recognition based on database preprocessing and parallel computing. In contrast to the majority of already existing methods, BacTag avoids the comparison of sequencing data to each allele sequence present in the database due to the database preprocessing. While the database preprocessing provides analysis time reduction, it also provides important information about database quality. Amongst other advantages of our method are the possibility to cope with any user-provided database, and the absence of the assembly step that potentially may help extend our approach to metagenomics datasets. We believe that our approach can be useful for a wide range of projects, including bacterial subspecies classification, clinical diagnostics of bacterial infections, and epidemiological studies.

5.6 Abbreviations

- MLST - multi-locus sequence typing;
- NGS - next-generation sequencing;
- ST - sequence type;
- WGS - whole-genome shotgun sequencing.

5.7 Author Statements

5.7.1 Acknowledgements

The authors would like to thank Martijn Vermaat, Sander Bollen and Peter van 't Hof for the helpful discussions and suggestions. We also would like to thank Louk Rademaker for the feedback on this manuscript.

5.7.2 Funding information

This work is part of the research program "Forensic Science" which is funded by grant number 727.011.002 of the Netherlands Organisation for Scientific Research (NWO). The funding body had no direct influence on the design of the study, collection of samples, analysis or interpretation of the data.

5.7.3 Availability of data and materials

- All the data analysed in this study are included in this manuscript and its Additional files 1¹⁵, 2¹⁶ and 3¹⁷.
- Results of the analysis done in this study are available via Figshare:
<https://doi.org/10.6084/m9.figshare.c.4041512.v1>
- BacTag is publicly available via
<https://git.lumc.nl/l.khachatryan/BacTag>.

¹⁵Available online https://static-content.springer.com/esm/art%3A10.1186%2Fs12864-019-5723-0/MediaObjects/12864_2019_5723_MOESM1_ESM.pdf

¹⁶Available online https://static-content.springer.com/esm/art%3A10.1186%2Fs12864-019-5723-0/MediaObjects/12864_2019_5723_MOESM2_ESM.pdf

¹⁷Available online https://static-content.springer.com/esm/art%3A10.1186%2Fs12864-019-5723-0/MediaObjects/12864_2019_5723_MOESM3_ESM.pdf

5.7.4 Authors' contributions

LK conception, pipeline design, acquisition of data, analysis and interpretation of data, manuscript drafting; MEMK conception, acquisition of data, manuscript editing; ATB conception, general supervision; JFJL pipeline design, manuscript editing, general supervision. All authors have read and approved this manuscript.

5.7.5 Ethics approval and consent to participate

Since in this research no human material or clinical records of patients or volunteers were used, this research is out of scope for a medical ethical committee. This was verified by the Leiden University Medical Center Medical Ethical Committee.

5.7.6 Competing interests

The authors declare that they have no competing interests.