

1 chewBBACA: A complete suite for gene-by-gene 2 schema creation and strain identification

3
4 Mickael Silva¹, Miguel Machado¹, Diogo N. Silva¹, Mirko Rossi², Jacob Moran-Gilad^{3,4}, Sergio Santos¹, Mario
5 Ramirez¹ and João André Carriço^{1*}

6
7 1 Instituto de Microbiologia, Instituto de Medicina Molecular, Faculdade de Medicina, Universidade de Lisboa,
8 Lisbon, Portugal 2 Department of Food Hygiene and Environmental Health, Faculty of Veterinary Medicine,
9 University of Helsinki, Finland 3 Faculty of Health Sciences, Ben-Gurion University of the Negev, Beer Sheva,
10 Israel 4 Public Health Services, Ministry of Health, Jerusalem, Israel

11
12 *To whom correspondence should be addressed.

13 **Please add author names and institutions into Editorial Manager and they will be added to the final version**
14 **of the article after acceptance. This is to ensure that the journal's double-blind peer review policy is respected.**

15 ABSTRACT

16
17 Gene-by-gene approaches are becoming increasingly popular in bacterial genomic
18 epidemiology and outbreak detection. However, there is a lack of open-source scalable
19 software for schema definition and allele calling for these methodologies. The chewBBACA
20 suite was designed to assist users in the creation and evaluation of novel whole-genome or
21 core-genome gene-by-gene typing schemas and subsequent allele calling in bacterial strains
22 of interest. The software can run in a laptop or in high performance clusters making it useful
23 for both small laboratories and large reference centers. ChewBBACA is available at
24 <https://github.com/B-UMMI/chewBBACA> or as a docker image at
25 <https://hub.docker.com/r/ummidock/chewbbaca/>.

27 DATA SUMMARY

28
29
30 1. Assembled genomes used for the tutorial were downloaded from NCBI in August 2016 by
31 selecting those submitted as *Streptococcus agalactiae* taxon or sub-taxa. All the assemblies
32 have been deposited as a zip file in FigShare
33 (<https://figshare.com/s/9cbe1d422805db54cd52>), where a file with the original ftp link for
34 each NCBI directory is also available.

35
36 2. Code for the chewBBACA suite is available at <https://github.com/B-UMMI/chewBBACA>
37 while the tutorial example is found at https://github.com/B-UMMI/chewBBACA_tutorial.

38
39 **I/We confirm all supporting data, code and protocols have been provided within the article**
40 **or through supplementary data files. ☒**

42 IMPACT STATEMENT

43

44 The chewBBACA software offers a computational solution for the creation, evaluation and
45 use of whole genome (wg) and core genome (cg) multilocus sequence typing (MLST) schemas.
46 It allows researchers to develop wg/cgMLST schemes for any bacterial species from a set of
47 genomes of interest. The alleles identified by chewBBACA correspond to potential coding
48 sequences, possibly offering insights into the correspondence between the genetic variability
49 identified and phenotypic variability. The software performs allele calling in a matter of
50 seconds to minutes per strain in a laptop but is easily scalable for the analysis of large datasets
51 of hundreds of thousands of strains using multiprocessing options. The chewBBACA software
52 thus provides an efficient and freely available open source solution for gene-by-gene
53 methods. Moreover, the ability to perform these tasks locally is desirable when the
54 submission of raw data to a central repository or web services is hindered by data protection
55 policies or ethical or legal concerns.

56

57

58 INTRODUCTION

59

60 Read mapping approaches using Single Nucleotide Polymorphisms (SNP)/Single Nucleotide
61 Variants (SNV) have been widely used for studying bacterial genomes [1]. However, gene-by-
62 gene (GbG) approaches have also been advocated in the context of genomic epidemiology as
63 an expansion of Multilocus Sequence Typing (MLST) [2] allowing portability, scalability, and
64 independence from a defined reference strain. For these reasons, GbG increasingly gains
65 popularity and has been adopted by PulseNet International as the method for bacterial strain
66 discrimination using high throughput sequencing [3]. GbG relies on comparing the draft
67 genome of a strain of interest against a pre-defined schema, typically using a BLAST [4] based
68 approach. This schema can be composed of core loci, which are present in all or the great
69 majority (e.g. 95%) of the analyzed strains (core genome MLST schemas or cgMLST), or
70 including all loci detected in the strains of interest. The latter are referred to as whole genome
71 or pan genome MLST schemas (wgMLST or pgMLST).

72 A locus in a schema can be a complete coding sequence (CDS) or a subsequence of it, as in
73 traditional MLST. Defining a locus as a CDS, allows linking the variability found to potential
74 changes in proteins and thus, with phenotype. The definition of locus is currently dependent
75 on the algorithm used for comparing loci and defining the schema, hampering the comparison
76 between different GbG approaches.

77 Only few software are available for GbG allele calling and no tools are available for schema
78 creation and validation. Two commercial platforms offer GbG analyses: Ridom SeqSphere+
79 (<http://ridom.de/seqsphere/>) and Bionumerics ([http://www.applied-](http://www.applied-maths.com/applications/wgmlst)
80 [maths.com/applications/wgmlst](http://www.applied-maths.com/applications/wgmlst)). Since these are proprietary, closed source software, their
81 GbG allele calling algorithms are incompletely described [5-6], although Ridom schemas have
82 been publicly released (<http://www.cgmlst.org/>).

83 BIGSdb was the first open-source freely available platform allowing cgMLST analysis [7] and
84 is currently the basis of the PubMLST website (<https://pubmlst.org/>). More recently,
85 Enterobase provides comprehensive cgMLST and wgMLST schemas and an allele calling
86 engine for three major foodborne bacterial pathogens (<https://enterobase.warwick.ac.uk/>).

87 A limitation of Enterobase is the requirement to submit reads to the website or to public
88 repositories (NCBI SRA/EBI ENA), since currently no stand-alone versions of their allele calling
89 algorithm are available. At present, the only published open-source stand-alone GbG allele

90 calling algorithm is Genome Profiler [8] which, however, uses a single CPU core making it
91 unsuitable for large scale analyses.
92 We developed chewBBACA to be a complete stand-alone pipeline for GbG analyses, including
93 constructing and validating novel cg/wgMLST schemas and performing CDS allele calling
94 suitable for large scale studies.

95

96 THEORY AND IMPLEMENTATION

97 chewBBACA is composed of three main modules: Schema Creation, Allele Calling and Schema
98 Evaluation. These modules can be used together in order to define and evaluate new
99 wg/cgMLST schemas for species of interest. A general workflow of such process is presented
100 in Fig. 1.

101

102 wg/cgMLST Schema creation

103 The First module is the Schema Creation, which allows the definition of wg/cgMLST schemas
104 from user provided complete genomes or draft assemblies, focusing on excluding paralogous
105 loci, detection of contaminated/poor quality assemblies and supporting user decisions
106 towards the identification of the most appropriate schema through interactive graphic data
107 analysis. This module uses an iterative approach for CDS comparison in the selection of loci
108 that is more computationally efficient than the Markov Clustering Step typically used in
109 software such as OrthoMCL [9] or CD-hit [10]. In order to create a wgMLST schema, the user
110 provides a set of genomes in FASTA format. The algorithm first defines the CDSs of each
111 genome using Prodigal [11]. In the next step, all the CDSs in the genomes are compared in a
112 pairwise fashion, resulting in a single file containing all unique CDSs identified in the genomes.
113 This comparison is a two-step process. Firstly, all the CDSs having identical sequence of other
114 CDSs but being smaller in length are removed and the larger CDS is kept. At the same time,
115 the algorithm also removes all CDSs with a length less than indicated in the “-l” parameter. In
116 the second step, the remaining CDSs are clustered in unique loci by performing an all-against-
117 all BLASTP search and calculating the Blast Score Ratio (BSR) [12]. CDSs with a BSR pairwise
118 comparison equal or greater than 0.6 are considered alleles of the same locus and the larger
119 allele (in bp) is kept in the list. This procedure defines the schema as a set of CDSs, each
120 representing the largest single allele of distinct loci. The Allele Calling module is then used to
121 populate the schema with alleles using the same genomes used for its creation. This step
122 allows the identification and exclusion of possibly paralogous loci. The Allele Calling algorithm
123 detects if a CDS in the genome under analysis matches more than one locus in the schema,
124 indicating that those loci can be paralogous. The Allele Calling module outputs a list of such
125 loci to be removed from the wgMLST schema or to be further investigated. From the created
126 wgMLST schema, cgMLST schemas can be defined by selecting the loci that are present in a
127 predetermined percentage of the analysed strains, typically 95%-99%.

128

129 Allele calling Algorithm

130 The Allele Calling algorithm is based on CDSs identified by Prodigal [11] with similarity
131 determined using a BLASTP BSR approach, allowing the detection of alleles with divergent
132 DNA sequences but similar encoded proteins. This allows the identification of alleles that
133 would be considered absent loci with BLASTN, while retaining the full diversity found at the
134 DNA sequence level. The algorithm is defined as presented in Fig. 2. A BLASTP database is
135 created, containing all the translated CDSs identified by Prodigal in the query genome. A 100%
136 DNA identity comparison is performed on all the genome of interest CDSs against each locus

137 allele database. If an exact match is found, an allele identification is attributed to the CDS (and
138 tagged as *EXC – Exact Match*, in the statistics output). If not, a BLAST BSR approach is used to
139 identify the allele. To improve computational efficiency, chewBBACA performs the similarity
140 search on each locus in the schema separately, parallelising the jobs using the specified
141 number of CPUs. For each locus, a short list containing the most divergent alleles is queried
142 against the BLASTP database. The BSR is calculated for each hit and based on these results
143 and a size validation step, the locus is either considered not found (tagged as *LNF – Locus Not*
144 *Found*) or a new allele of the locus is inferred. The size validation step excludes alleles larger
145 than or smaller than 20% of the locus allele length mode (Defined as *ASM – Alleles Smaller*
146 *than Mode* or *ALM – Alleles Larger than Mode*) (Fig.3a). Furthermore, the identification of loci
147 as duplicated in the genome of interest is also reported. Such matches are identified as *Non-*
148 *Informative Paralogous Hits (NIPH)*, if at least two CDSs have best matches with alleles of the
149 same locus but presenting less than 100% identity, or *NIPHEM – NIPH Exact Match* if 100%
150 identity to existing alleles is detected (Fig. 3b). Furthermore, the algorithm detects whether
151 the CDS match is close to the 5' or 3' ends of a contig and a larger allele that contains the
152 matched sequence would exceed the contig length. Such sequences are tagged as *Possible*
153 *Locus On the Tip (PLOT)* (Fig. 3c). Finally, the Allele Calling module identifies possible
154 paralogous (as described above) checking if there are CDS matching alleles in two or more
155 different loci (Fig. 3d). After running each genome, the loci database is updated with the
156 newly found alleles and, whenever required, a locus short-list is also updated with a new
157 divergent allele.

158

159 **Schema Evaluation**

160 This module allows the assessment of the suitability of including each locus in a schema
161 through a suite of functions to graphically explore and evaluate the type and extent of allelic
162 variation detected in each of the chosen loci. This module also creates multiple sequence
163 alignments of the alleles of each locus using MAFFT [13] and constructs neighbour-joining
164 trees using ClustalW2 [14], allowing the exploration of the potential consequences of the
165 variability at each locus. This module can be used to analyse any existing cg/wgMLST schema,
166 including those created by other methodologies, since the analysis input is a set of FASTA
167 files, one per locus, with all identified alleles.

168

169 A more complete description of each module and their functionalities is available at
170 <https://github.com/B-UMMI/chewBBACA/wiki>

171

172 **Benchmark**

173 The performance of chewBBACA's allele calling algorithm was evaluated for *Streptococcus*
174 *agalactiae* assemblies (~2Mb genome) using a cgMLST schema of 1,264 loci. Benchmarks
175 were performed on a high-performance cluster (HPC) with Intel® Xeon™ E5-2630 v4 @
176 2.20GHz CPUs, up to 256Gb RAM and an SSD distributed storage in RAID1; a laptop with Intel®
177 Core™ i5-7200U @ 2.50GHz x 4 CPUs, 8Gb RAM and a NVMe SSD storage; and a laptop with
178 Intel® Core™ i7-3630QM @ 2.40GHz x 8 CPUs, 8Gb RAM and a SATA2 HDD storage. Allele
179 calling was conducted for 100 *S. agalactiae* assemblies in the HPC cluster using 2 to 40 CPUs,
180 and for a subset of 50 assemblies in both laptops using 2 and 4 CPUs (Fig 4). Each CPU data
181 point was run 5 times. In terms of CPU performance, the time it took to run each sample
182 decreased almost linearly up to 16 CPUs, at which point disk possibly I/O storage access
183 becomes the bottleneck and no increase in performance is observed. At peak performance,

184 allele calling takes approximately 4 seconds for each sample. Since this algorithm is I/O
185 intensive, a substantial increase in performance can be observed in a laptop with SSD storage,
186 in comparison to another with HDD storage. In any case, allele calling in a laptop using 4 CPUs
187 can be performed for each sample in approximately 9 seconds with SSD storage and 15
188 seconds with HDD storage.

189

190 **Usage Example**

191 A tutorial providing a complete usage example, demonstrating the creation of a schema for
192 *Streptococcus agalactiae*, from publicly available complete genomes and assemblies available
193 at NCBI/ENA is provided at https://github.com/B-UMMI/chewBBACA_tutorial.

194

195 **CONCLUSION**

196

197 The chewBBACA suite was developed to allow performing GbG analyses in high-end Unix
198 based laptops but chewBBACA can also be easily run in HPC, facilitating its adoption into large-
199 scale automated analysis pipelines. The good performance of the software in current laptops
200 and in HPCs, allows a flexible implementation in small laboratories or large reference centres.
201 The allele calling engine of chewBBACA uses FASTA files with draft assemblies or complete
202 genomes as input and returns as output an allelic profile matrix and a set of FASTA files
203 containing the full allelic diversity of each locus. Currently available cg/wgMLST schemas, can
204 be adapted to run using chewBBACA's allele calling engine. The chewBBACA suite is the first
205 to provide schema creation tools and to enforce CDS allele calling, which can be important to
206 evaluate phenotypic diversity including the identification of the potential mechanisms
207 underlying the success of particular clones. Since there is an urgent need for bioinformatics
208 solutions that will allow the development of nomenclature-based schemas [15], future work
209 will focus on centralised repositories for schemas and allele definitions that can be
210 synchronised with local allele calling outputs to facilitate the development of common
211 schemas and nomenclatures for cg/wgMLST allowing a more widespread application of GbG
212 methodologies in public health.

213

214

215 **AUTHOR STATEMENTS**

216

217 **Funding information**

218 Mickael Silva and Miguel Machado have been supported by INNUENDO project
219 (<https://www.innuendoweb.org>) co-funded by the European Food Safety Authority (EFSA),
220 grant agreement GP/EFSA/AFSCO/2015/01/CT2 ("New approaches in identifying and
221 characterizing microbial and chemical hazards"). The conclusions, findings, and opinions
222 expressed in this review paper reflect only the view of the authors and not the official position
223 of the European Food Safety Authority (EFSA).

224 This work was partially supported by the following projects: ONEIDA project (LISBOA-01-0145-
225 FEDER-016417) co-funded by FEEI - "Fundos Europeus Estruturais e de Investimento" from
226 "Programa Operacional Regional Lisboa 2020" and by national funds from FCT - "Fundação
227 para a Ciência e a Tecnologia" and BacGenTrack (TUBITAK/0004/2014) [FCT/ Scientific and
228 Technological Research Council of Turkey (Türkiye Bilimsel ve Teknolojik Araştırma Kurumu,
229 TÜBİTAK)].

230

231 **Acknowledgements**

232 The authors would like to thank Eduardo Taboada and Peter Kruczkiewicz from National
233 Microbiology Laboratory at Lethbridge, Public Health Agency of Canada, Lethbridge, Alberta,
234 Canada., for insightful discussions.

235

236 **Conflicts of interest**

237 None declared.

238

239

240 **ABBREVIATIONS**

241

242 SNP, Single Nucleotide Polymorphisms; SNV, Single Nucleotide Variants; GbG, gene-by-gene;
243 cgMLST, core genome MLST; wgMLST, whole genome MLST; pgMLST, pan genome MLST;
244 CDS, coding sequence; BSR, Blast Score Ratio;

245

246

247 **REFERENCES**

248

249 1. Lynch, T. et al. A Primer on Infectious Disease Bacterial Genomics. *Clinical Microbiology*
250 *Reviews*, 2016; 29(4), pp.881–913.

251 2. Maiden, M.C. et al. Multilocus sequence typing: a portable approach to the identification
252 of clones within populations of pathogenic microorganisms. *Proceedings of the National*
253 *Academy of Sciences of the United States of America*, 1998; 95(6), pp.3140–3145.

254 3. Nadon, C. et al. PulseNet International: Vision for the implementation of whole genome
255 sequencing (WGS) for global food-borne disease surveillance. *Euro surveillance: bulletin*
256 *européen sur les maladies transmissibles = European communicable disease bulletin*, 2017;
257 22(23), pp.13–24.

258 4. Altschul, S.F. et al. Basic local alignment search tool. *Journal of molecular biology*, 1990;
259 215(3), pp.403–410.

260 5. Moura, A. et al. Whole genome-based population biology and epidemiological surveillance
261 of *Listeria monocytogenes*. *Nature Microbiology*, 2016; pp.1–10.

262 6. Ruppitsch, W. et al. Defining and Evaluating a Core Genome Multilocus Sequence Typing
263 Scheme for Whole-Genome Sequence-Based Typing of *Listeria monocytogenes*. D. J.
264 Diekema, ed., 2015; 53(9), pp.2869–2876.

265 7. Jolley, K.A. & Maiden, M.C.J. BIGSdb: Scalable analysis of bacterial genome variation at the
266 population level. *BMC Bioinformatics*, 2010; 11, p.595.

267 8. Zhang, J. et al. Refinement of whole-genome multilocus sequence typing analysis by
268 addressing gene paralogy., 2015; 53(5), pp.1765–1767.

269 9. Li, L., Stoeckert, C.J. & Roos, D.S. OrthoMCL: identification of ortholog groups for eukaryotic
270 genomes. *Genome Research*, 2003; 13(9), pp.2178–2189.

271 10. Fu, L. et al. CD-HIT: accelerated for clustering the next-generation sequencing data.
272 *Bioinformatics*, 28(23), 2012; pp.3150–3152.

273 11. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site
274 identification. *BMC Bioinformatics*, 2010; 11(1), p.119.

- 275 12. Rasko, D.A., Myers, G. & Ravel, J. Visualization of comparative genomic analyses by BLAST
276 score ratio. BMC Bioinformatics, 2005; 6(1).
- 277 13. Katoh, K. et al. MAFFT: a novel method for rapid multiple sequence alignment based on
278 fast Fourier transform. Nucleic Acids Research, 2002; 30(14), pp.3059–3066.
- 279 14. Larkin, M.A. et al. Clustal W and Clustal X version 2.0. Bioinformatics, 2007; 23(21),
280 pp.2947–2948.
- 281 15. Moran-Gilad, J. Whole genome sequencing (WGS) for food-borne pathogen surveillance
282 and control - taking the pulse. Euro surveillance : bulletin européen sur les maladies
283 transmissibles = European communicable disease bulletin, 2017; 22(23), p.30547.

284
285

286

287 DATA BIBLIOGRAPHY

288

289 1. Assembled genomes used for the tutorial were downloaded from NCBI in August 2016 by
290 selecting those submitted as *Streptococcus agalactiae* taxon or sub-taxa. All the assemblies
291 have been deposited as a zip file in FigShare
292 (<https://figshare.com/s/9cbe1d422805db54cd52>), where a file with the original ftp link for
293 each NCBI directory is also available.

294

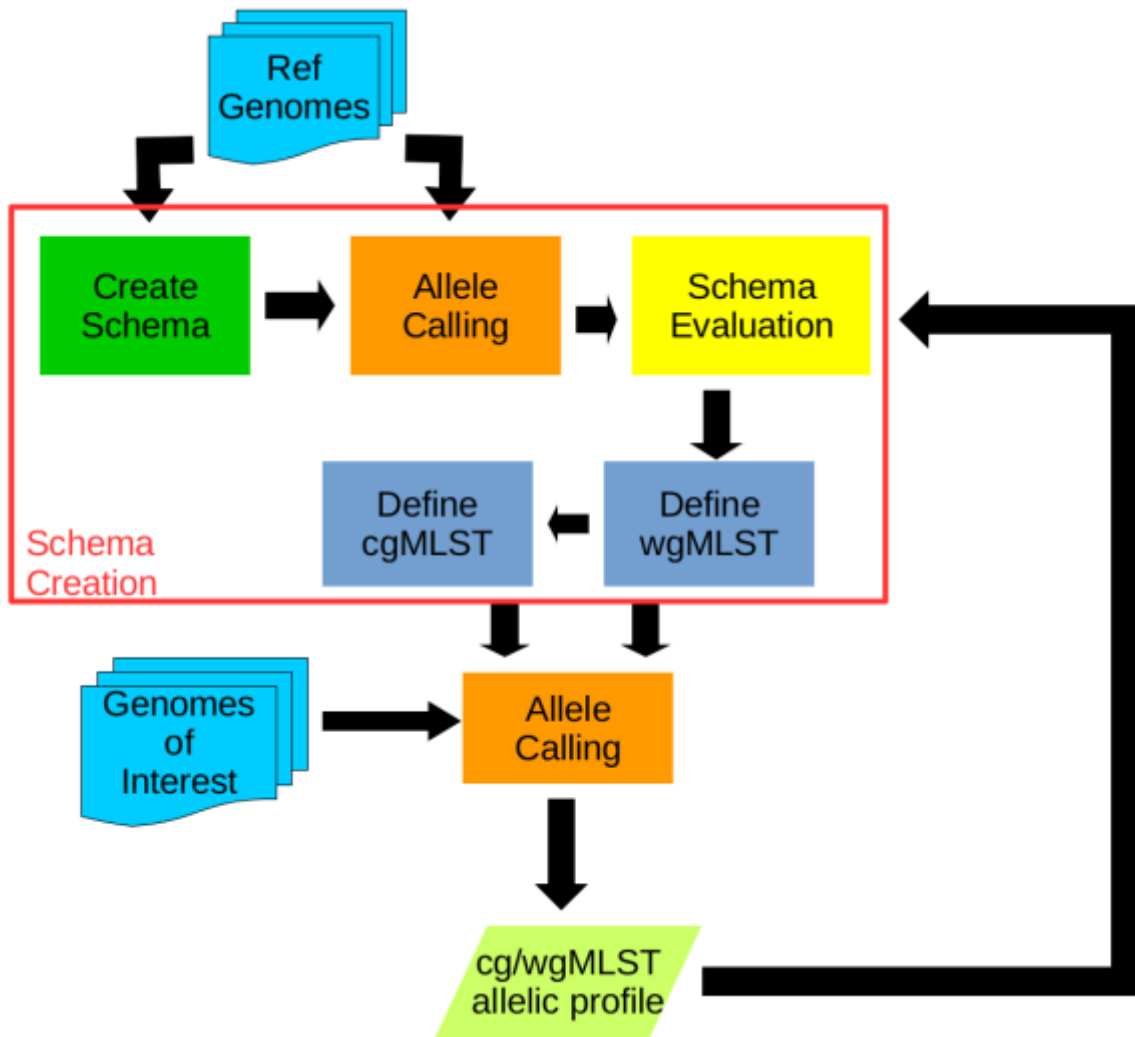
295 2. Code for the chewBBACA suite is available at <https://github.com/B-UMMI/chewBBACA>
296 while the tutorial example is found at https://github.com/B-UMMI/chewBBACA_tutorial.

297

298

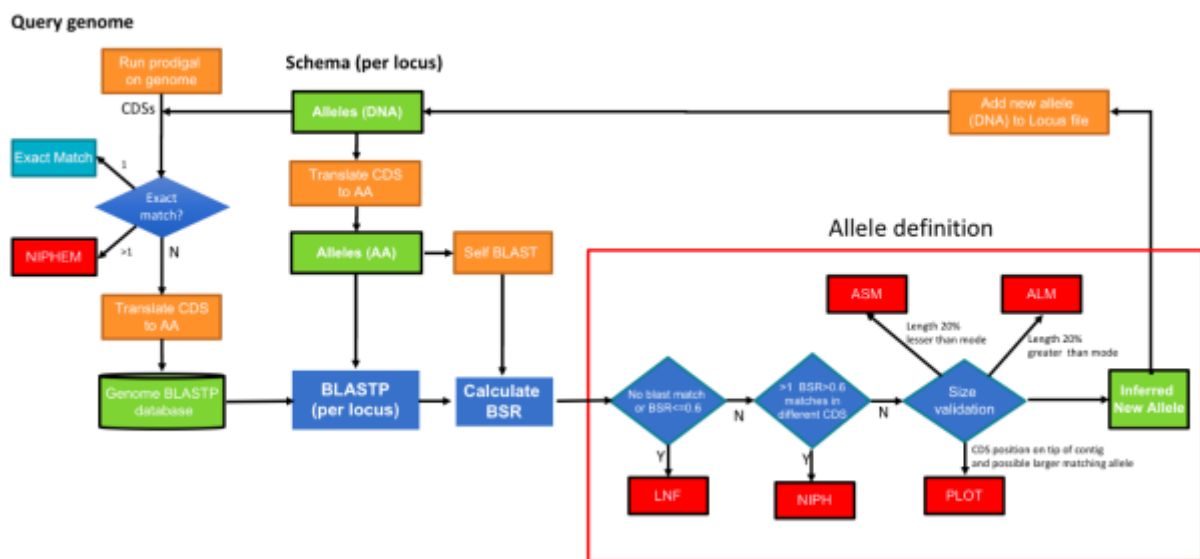
299 FIGURES AND TABLES

300



301
302

Figure 1 - chewBBACA workflow from schema definition to schema evaluation



303
304

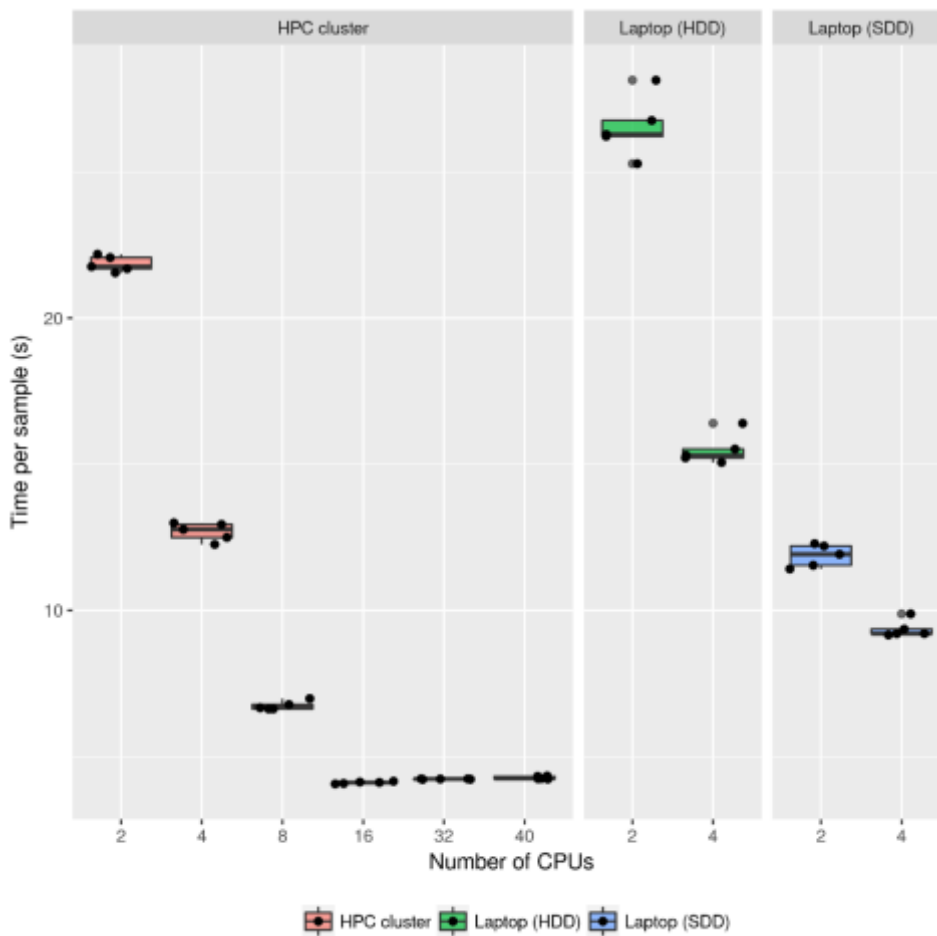
Figure 2 - chewBBACA allele calling algorithm.

305



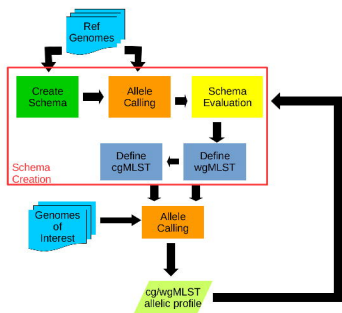
306
307
308
309

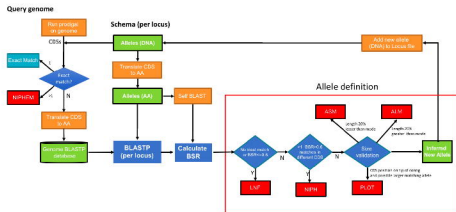
Figure 3 - chewBBACA Allele definition outputs. A) size exclusion of alleles 20 % smaller or larger than the allele length mode for the loci B) Detection of loci duplication on the draft genome C) detection of locus identified on the 5' or 3' ends of the contig D) Detection of paralogous loci



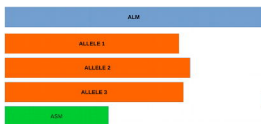
310
311
312
313

Figure 4 – Benchmarking of chewBBACA's allele calling algorithm for bacterial genome assemblies (~2Mb) using a cgMLST schema of 1264 loci on a HPC cluster and two laptops with different storage devices. The allele calling was executed 5 times for each CPU data point.

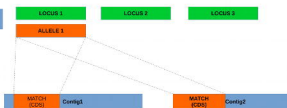




(a) ALM/ ASM



(b) NIPH/NIPHEM



(c) PLOT



(d) Paralog detection

