

GigaScience

GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes --Manuscript Draft--

Manuscript Number:	GIGA-D-21-00186	
Full Title:	GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes	
Article Type:	Technical Note	
Funding Information:	agricultural research service (0500-00093-001-00-D)	Dr. Adam R Rivers
	agricultural research service (6066-21310-005-D)	Dr. Adam R Rivers
	agricultural research service (6066-21310-005-28-S)	Dr. Christopher R. Reisch
Abstract:	<p>Background: CRISPR-Cas systems have expanded the possibilities for gene editing in bacteria and eukaryotes. There are many excellent tools for designing the CRISPR-Cas guide RNAs for model organisms with standard Cas enzymes. GuideMaker is intended as a fast and easy-to-use design tool for atypical projects with 1) non-standard Cas enzymes, 2) non-model organisms, or 3) projects that need to design a panel of guide RNAs (gRNA) for genome-wide screens.</p> <p>Findings: GuideMaker can rapidly design gRNAs for gene targets across the genome from a degenerate protospacer adjacent motif (PAM) and a GenBank file. The tool applies Hierarchical Navigable Small World (HNSW) graphs to speed up the comparison of guide RNAs. This allows the user to design gRNAs targeting all genes in a typical bacterial genome in about 1-2 minutes.</p> <p>Conclusions: Guidemaker enables the rapid design of genome-wide gRNA for any CRISPR-Cas enzyme in non-model organisms. While GuideMaker is designed with prokaryotic genomes in mind, it can efficiently process smaller eukaryotic genomes as well. GuideMaker is available as command-line software, a stand-alone web application, and a tool in the CyCverse Discovery Environment. All versions are available under a Creative Commons CC0 1.0 Universal Public Domain Dedication.</p>	
Corresponding Author:	Adam R Rivers, Ph.D. USDA Agricultural Research Service Gainesville, FL UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	USDA Agricultural Research Service	
Corresponding Author's Secondary Institution:		
First Author:	Ravin Poudel, Ph.D.	
First Author Secondary Information:		
Order of Authors:	Ravin Poudel, Ph.D.	
	Lidimarie Trujillo Rodriguez, B.S.	
	Christopher R. Reisch, Ph.D.	
	Adam R Rivers, Ph.D.	
Order of Authors Secondary Information:		
Additional Information:		
Question	Response	

<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	<p>Yes</p>
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum</p>	<p>Yes</p>

[Standards Reporting Checklist?](#)



1 **GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes**

2 Ravin Poudel^{1,2}, Lidimarie Trujillo Rodriguez², Christopher R. Reisch², and Adam R. Rivers^{1*}

3 ¹Genomics and Bioinformatics Research Unit, USDA Agricultural Research Service, Gainesville, FL, 32608,
4 USA

5 ²Department of Microbiology and Cell Science, Institute of Food and Agricultural Sciences, University of
6 Florida, Gainesville, FL 326011, USA

7 * Correspondence: adam.rivers@usda.gov

8

9 **Abstract**

10 **Background:**

11 CRISPR-Cas systems have expanded the possibilities for gene editing in bacteria and eukaryotes. There are
12 many excellent tools for designing the CRISPR-Cas guide RNAs for model organisms with standard Cas
13 enzymes. GuideMaker is intended as a fast and easy-to-use design tool for atypical projects with 1) non-
14 standard Cas enzymes, 2) non-model organisms, or 3) projects that need to design a panel of guide RNAs
15 (gRNA) for genome-wide screens.

16 **Findings:**

17 GuideMaker can rapidly design gRNAs for gene targets across the genome from a degenerate protospacer
18 adjacent motif (PAM) and a GenBank file. The tool applies Hierarchical Navigable Small World (HNSW)
19 graphs to speed up the comparison of guide RNAs. This allows the user to design gRNAs targeting all genes
20 in a typical bacterial genome in about 1-2 minutes.

21 **Conclusions:**

22 Guidemaker enables the rapid design of genome-wide gRNA for any CRISPR-Cas enzyme in non-model
23 organisms. While GuideMaker is designed with prokaryotic genomes in mind, it can efficiently process
24 smaller eukaryotic genomes as well. GuideMaker is available as command-line software, a stand-alone web
25 application, and a tool in the CyCverse Discovery Environment. All versions are available under a Creative
26 Commons CC0 1.0 Universal Public Domain Dedication.

27

28 **Keywords** PAM, CRISPR-Cas, gRNA, HNSW

29

30 **Introduction**

31 CRISPR-Cas technology enables rapid and efficient genome editing in both prokaryotic and eukaryotic cells
32 [1,2]. CRISPR-based systems are set apart from other genome editing tools by the ease with which they can
33 be programmed to target specific sequences. Almost any DNA sequence in the cell can be targeted as long as
34 it possesses a compatible protospacer adjacent motif (PAM). The PAM is a conserved sequence that flanks
35 the DNA target site, known as the protospacer, and must be present for target recognition [3]. The target
36 specifying guide-RNA (gRNA) can be supplied as RNA, or encoded in DNA, depending on the organism
37 under investigation. Although CRISPR-Cas is often used to edit single genes in eukaryotes, it is increasingly
38 used for other purposes in prokaryotic and eukaryotic organisms, including non-model organisms [4].

39 The *Streptococcus pyogenes* Cas9 (SpCas9) was the first Cas described [5] and it is still the most widely
40 used enzyme in CRISPR in gene editing. Other Cas enzymes described early in the CRISPR revolution, such
41 as the *Staphylococcus aureus* Cas9 and the *Acidaminococcus* Cas12a, are also commonly used [6,7]. Accordingly, the
42 parameters for these enzymes are often included in computational tools to identify CRISPR target sites [8–
43 11]. Cas9 enzymes from other organisms and other Cas-associated proteins that can cleave dsDNA, ssDNA,
44 ssRNA, and insert transposon elements have also been described and have their place in molecular toolkits
45 [12–18]. Each of these enzymes generally possesses its own requirements, such as PAM sequence constraints,

46 PAM orientation, and protospacer length. Many of these CRISPR-Cas systems have been repurposed to
47 enable molecular genetics techniques like gene deletions, gene insertions, transcriptional depletion and
48 activation, and translational repression [12,19–22]. Some of these techniques can be scaled to the genome
49 level with chip-synthesized oligonucleotides and pooled approaches to screening [23]. In pooled screens,
50 high-throughput DNA sequencing is used to identify how the pool has changed over time to elucidate genes
51 that affect cells' fitness in specific conditions. Given the diversity of the CRISPR systems and their uses,
52 identifying appropriate target sites is not trivial, especially for the number of targets needed for genome-scale
53 experiments.

54 Here we introduce GuideMaker, a computational tool to identify target sites and design gRNA
55 sequences that is not limited to any specific CRISPR system or organism. Guidemaker is most useful for a
56 few kinds of CRISPR experiments. The first use case is designing pools of gRNAs for genome-wide
57 screening experiments like Perturb-seq and CRISPR pool [23,24]. GuideMaker is optimized for making the
58 all-versus-all comparisons necessary to design a genome-wide screen and return candidate gRNAs for every
59 gene locus. The tool allows the user to filter targets based on their proximity to features of interest, like the
60 start codon for any coding sequence. The second major use case is for researchers working with non-model
61 organisms. Online gRNA design tools often have a limited number of preselected genomes available for
62 analysis because most methods require PAM site positions to be precomputed. GuideMaker rapidly computes
63 all guide positions on demand so the user can provide a set of GenBank files from any organism for analysis.
64 The third use case is for researchers working with Cas enzymes other than the canonical versions of Cas9,
65 Cas12a (Cpf1), or Cas13 with different PAM and target site requirements. GuideMaker allows the user to
66 specify a custom PAM with variable length, including degenerate nucleotides and allows the PAM to be on
67 either the 3' or 5' side of the protospacer. These features allow GuideMaker to support any current or future
68 CRISPR-Cas system. Since the determination of which CRISPR-Cas system functions best in any given
69 organism is not predictable, this tool is highly relevant to researchers developing CRISPR tools in new
70 species. In some cases, GuideMaker may not be the best choice. There are mature tools for designing gRNAs
71 in model organisms with common CRISPR-Cas systems and targeting a small number of loci [25,26]. Some

72 of these employ sophisticated statistical models to select the best Cas9 gRNA candidates and may be a better
73 choice for well-studied systems [8]. Because there is limited experimental data on most Cas/organism
74 combinations, GuideMaker relies on design heuristics rather than machine learning-based identification
75 methods.

76

77 **Methods**

78 **Main features, input parameters, and workflow**

79 GuideMaker is designed to be easy to use as either a web application or a command-line utility. The key
80 features of GuideMaker are:

- 81 1. All the potential guides in a genome can be quickly designed in one run.
- 82 2. It can design gRNAs for any small to medium size genome (up to about 500 megabases).
- 83 3. It can design gRNAs for any PAM sequence from any Cas system.
- 84 4. Search is customizable through user-defined guide parameters (as highlighted in Figure 1). These
85 features are specific to organisms, CRISPR-Cas systems, and experiments. Tuning these parameters
86 can improve the sensitivity and specificity of gRNA.
- 87 5. Users can exclude specific restriction sites from guides to preserve those sites for downstream
88 experiments.
- 89 6. It creates control gRNAs based on the input genome. In CRISPR experiments it is often desirable to
90 create negative control gRNAs to evaluate off-target binding. GuideMaker provides the user with
91 realistic control gRNAs that are highly divergent from sequences adjacent to PAM sites.
- 92 7. Provides an interactive visualization and exploratory tool to evaluate the guides.
- 93 8. Provides tabular result files which can be used for the design and ordering of gRNAs.
- 94 9. The software can be run as a web application [27], a CyVerse application, or a command-line
95 application [28]. Server code is included for running local instances of the web application as well.

96 A typical workflow of GuideMaker involves three major steps (Figure 2). In the first step, the user
97 uploads the input genome in one or more .gbk or gzipped .gbk files and defines the PAM and gRNA
98 parameters (as highlighted in Figure 1). Guidemaker identifies and filters target sites, then returns summary
99 data to the graphical environment (Figure 2). Users can use the interactive plots to learn more about the
100 identified gRNAs and sort them by genome coordinates or locus tag. In the final step, GuideMaker provides
101 the results as downloadable files under the results section. These files are used for synthesizing guides. The
102 command-line version of GuideMaker has similar input parameters as the web application, with the flexibility
103 to generate plots and configure the underlying hyper-parameters for the Hierarchical Navigable Small World
104 (HNSW) graph, or to run the web application locally. To make the application easier to install we distribute
105 the application as a Bioconda environment [29], Docker container [30], Python package on Github [28],
106 through the Cyverse discovery environment [31] or as an online web application [27]. Detailed information
107 on accessing the software through various methods is available on the project homepage [32].

108 **Search method**

109 GuideMaker initially scans the genome, recording all candidate guide sequences adjacent to the
110 specified PAM sequence on both DNA strands (Figure 3). Candidate guides are then optionally checked for
111 the restriction sites. Next, the candidates guides are searched for a unique "seed region" closest to the PAM
112 site and candidate gRNAs that are not unique in their "seed region" are removed. Then, approximate nearest
113 neighbor search is used to remove candidate guides too similar to PAM adjacent sequences in the genome,
114 based on Hamming distance (the number of substitutions required to turn one DNA sequence into another
115 equal-length sequence). The approximate nearest neighbor search is performed using the Hierarchical
116 Navigable Small World (HNSW) graph method in the Non-Metric Space Library (NMSLIB) [33,34]. An
117 index of all the initial candidate guides is created using the bitwise Hamming distance metric. Each guide with
118 a unique "seed region" is compared to all candidate guides and any guides with Hamming distances below the
119 user-set threshold are removed. This differs from the standard procedure of indexing the genome and
120 mapping each candidate guide against the whole genome then parsing each result. HNSW has a search

121 complexity of $\mathcal{O}(\log N)$ and index complexity of $\mathcal{O}(N \cdot \log N)$ [33]. Finally, user-defined criteria are applied
122 specifying the proximity and orientation of guides relative to genomic features like genes. A list of guides is
123 then returned to the user with relevant information about the guide and its target genomic features.

124 The core of GuideMaker's search method is the HNSW method in NMSLIB [34]. The method
125 builds a multilayer graph index of the input data and has several parameters that can be optimized for index
126 building and search to trade-off speed and accuracy. Graph construction is the most time-consuming step in
127 our tests, and thus grid optimization was run to minimize run time while keeping recall above 99% relative to
128 the ground truth exact nearest-neighbor search. The grid-optimization parameters: [M, efc, ef, and post] used
129 in the HNSW graph for approximate nearest neighbor search have been optimized for bacterial genomes. A
130 script for re-optimization (flag `--config`) of these hyper-parameters is included in the command-line version of
131 the software.

132 **Computational performance**

133 Genomes of different sizes, GC content, and chromosome numbers were used to test the speed and
134 scalability of GuideMaker (Supplementary Table 1). For benchmarking the performance, the same parameters
135 were used unless a specific parameter was being tested: a PAM motif of 'NGG', 3' pam orientation, target
136 length of 20, lsr (length of seed region) of 11, before and after parameters of 500, knum of 10, controls of 10,
137 dist of 3 and threads of 32. We profiled the performance of GuideMaker with different threads [1, 2, 4, 8, 16,
138 and 32] in processors with and without the AVX2 processor instruction set. All tests were run on a single
139 compute node with 2 x 24 core Intel(R) Xeon(R) Platinum 8260 CPU @ 2.40 GHz with Cascade Lake
140 microarchitecture. Three bacterial genomes, a fungal genome, and a plant genome were used in performance
141 benchmarking: *Escherichia coli* (K12), *Pseudomonas aeruginosa* (PAO1), *Burkholderia thailandensis* (E264), *Arabidopsis*
142 *thaliana*, and *Aspergillus fumigatus*. For the gene or locus-specific comparisons, only the guides within the locus
143 coordinates (i.e. zero feature distance) were considered.

144 **Comparison to existing design method**

145 We compared the results of GuideMaker with the results of the online version of CHOPCHOP[35].
146 GuideMaker and CHOPCHOP parameters were set to approximate the same search. The length of the target
147 sequence was set to 20 and zero mismatches were allowed in the seed region (11bp) of the target. The
148 *Escherichia coli* (str. K-12/MG1655) genome was used with the online version of CHOPCHOP since it has a
149 limited number of genomes. Targets were searched in 40 Kbp increments to account for CHOPCHOP's size
150 limitations. Target sequences were searched across multiple 40 Kbp segments of *E.coli* genome
151 (NC_000913.3:2001-42000, NC_000913.3:80001-120000, NC_000913.3:160001-200000,
152 NC_000913.3:240001-280000, and NC_000913.3:320001-360000). We also searched for target sequences
153 and genes/locus_tags within 40Kbp of (NC_000913.3:2001-42000) to compare identifications at the locus
154 level. Ratios between tools were calculated by dividing the number of gRNA identify with GuideMaker by the
155 number of CHOPCHOP identified gRNA to represent the proportion of guides identified by both
156 GuideMaker and CHOPCHOP.

157

158 **Results**

159 The time for Guidemake to complete a typical run identifying all SpCas9 gRNAs (PAM 'NGG') in a bacterial
160 genome using 8 compute cores was 75 seconds for *E. coli* and 130 seconds for *P. aeruginosa* (Figure 4). For
161 SaCas9 and StCas9, which have a longer PAM sequence ('NGRRT' and 'NNAGAAW' respectively, with 3'
162 PAM orientation) and thereby fewer potential targets, the same genomes ran in 19 or 5 seconds
163 (Supplementary Figures 1). The fungus *Aspergillus fumigatus* (28MB) and plant *Arabidopsis thaliana* (114 MB)
164 have larger genomes but are still processed quickly. *A. fumigatus* processed between 23 – 304 seconds, while
165 *A. thaliana* processed in 250-921 seconds depending on the number of cores, AVX2 instructions, and PAM
166 sequence (Supplementary Figures 2). GuideMaker can take advantage of Advanced Vector Extensions
167 (AVX2) on newer x86 processors, which improves the search speed because HNSW search is accelerated
168 with AVX2 (Supplementary Figure 3). The acceleration was larger when fewer processors were available
169 (Supplementary Figure 3). With more processors, the run time was similar regardless of AVX2 use. The
170 HNSW algorithms are parallelized, and indexing-and-search takes most of the compute time in GuideMaker

171 so the software scales well when additional cores are added up to 8 cores (Supplementary Figure 3). In
172 practice it scaled up sub-linearly with genome size, globally estimating Cas9 guides for *E. coli* MG1655
173 (4.6MB) in 75 seconds and *A. thaliana* (114.1MB) in 921 seconds, both on 8 cores (Memory usage: 1.9GB for
174 *E. coli* and 15.4 GB for *A. thaliana*, Supplementary Figure 4).

175 The results of Guidemaker were compared with the popular guide design software CHOPCHOP
176 version 3 [35]. When GuideMaker's filtering settings are set to match CHOPCHOP, the results are very
177 similar and 99.9% of the targets identified by GuideMaker fall within 2bp of target coordinates returned by
178 CHOPCHOP. When GuideMaker's unique seed region criterion was not applied at the loci level, the average
179 number of guides identified by the two approaches was similar per locus (Mean GuideMaker = 116.8, Mean
180 CHOPCHOP = 113.6, p-value = 0.86, Supplementary Table 2). Although the number of guides identified
181 per gene locus differed, none of the genes were missed by either tool. GuideMaker's default requirement of a
182 seed region is more stringent than CHOPCHOP, and with it enabled, GuideMaker returns (count=1787)
183 38.4% (for 2Kbp-42Kbp regions) of the targets compared to CHOPCHOP (count=4651) *E. coli* K12. At the
184 sequence level, 96.7% of the identified gRNA (1729/1787) from both of the tools had identical sequences.
185 The more stringent filtering could potentially reduce off-targeting but that would need to be experimentally
186 validated in a range of organisms. The ratio of gRNA found by both the tools across the multiple 40Kbp
187 regions was 39.2% (sd= 1.9%, Supplementary Table 3) when using Guidemaker's more stringent default
188 settings. This ratio was calculated by dividing the number of gRNA from GuideMaker by the number from
189 CHOPCHOP for each 40Kb region.

190

191 **Discussion**

192 Designing gRNAs is a two-step process where GuideMaker first identifies potential guides adjacent to PAM
193 sequences and then filters the potential guides based on multiple criteria. The most important criterion is that
194 each guide has a minimum edit distance from any other sequence adjacent to a PAM site in the genome; this
195 decreases the likelihood of off-target binding. The second way GuideMaker reduces off-target binding is by

196 requiring that a set number of bases near the PAM site are unique from any other candidate guide. The 8
197 bases nearest the PAM are the most important for target specificity, and any mismatch is sufficient to prevent
198 binding [36,37]. The length of the unique region should be set with consideration for the size of the genome
199 since requiring short unique regions will limit the number of total guides that can be found. For example,
200 requiring that every gRNA be unique in the first 3 bp would only allow for $4^3 = 64$ possible guides to be
201 designed. For normal *--lsr* values of 9-12 this is only limiting for human-sized genomes and can be disabled by
202 setting *--lsr* to 0. All guides designed by GuideMaker are perfect matches to a single site in the genome.
203 Specificity is obtained by requiring all similar PAM-adjacent sequences to be unique in the critical "seed
204 region" *and* have a total number of mismatches that exceed the user-defined threshold. This double criterion
205 is expected to increase specificity.

206 The primary goal of the current version of our software is to support the design of gRNAs in non-
207 standard Cas enzymes for non-model organisms at the genome-scale. It is known that gRNA's do not
208 perform equally, thus empirical experiments will be needed to fully validate the functionality and efficacy of
209 gRNA predictions. Given the similarity in targets identified by GuideMaker and CHOPCHOP, we anticipate
210 performance will be similar to the current state of the art but applicable to more design use cases. When a
211 unique seed region and Hamming distance-based filters were applied, GuideMaker created guides more
212 conservatively, generating only about 40% of the guides created by CHOPCHOP. While CHOPCHOP has
213 an option to specify the maximum number of mismatches in the first 9 bp or the whole guide, it does not
214 allow the application of both criteria. While there are small differences in the number and position of guides
215 generated by GuideMaker, with GuideMaker being more conservative by default, both programs create
216 enough guides to target nearly all gene loci in the genome of *E. coli*. If experimentally validated data become
217 available from genome-wide screens with different Cas enzymes, the future versions of GuideMaker could
218 potentially incorporate scoring matrices to help rank candidate guides.

219 Guidemaker is a fast and flexible tool for designing guide RNA across the entire genome in non-
220 model organisms or with non-canonical Cas enzymes. It takes advantage of fast HNSW search to quickly

221 index and search new genomes. Several parameters can be tuned to ensure compatibility with the specific
222 application of the user. For example, GuideMaker checks the designed gRNA for a given restriction enzyme
223 site to prevent incompatibility with the cloning strategy. Second, the maximum distance from a target
224 sequence from the start of an annotated feature can be chosen to disrupt promoters or the beginning of the
225 coding sequence, since these sites are preferred for CRISPRi experiments. GuideMaker also creates off-target
226 gRNAs for use as negative controls in high-throughput experiments. Lastly, the program plots the results for
227 visual exploration of the targets and exports the data as .csv files. The software is available as a command-line
228 application, a web application, and is integrated into the CyVerse Discovery Environment to provide users
229 with a range of usage options.

230

231 **Availability and Requirements**

232 Project name: GuideMaker

233 Project home page: <https://guidemaker.org>

234 Operating system(s): Linux or MacOS

235 Programming language: Python >=3.6

236 Other requirements: 'pybedtools==0.8.2', 'nmslib>=2.0.6','altair', 'streamlit>0.80.0

237 License: CC0 1.0 Public Domain Dedication

238

239

240 **Competing Interests**

241 Authors declare no competing interests

242

243 **Data Availability**

244 The source code and command-line executables for GuideMaker are available at the Zenodo [38] and can be
245 installed directly from Github [28], Bioconda [29], or as a Docker container [30]. Data and code to reproduce
246 the analysis in the paper are available at Zenodo [39]. As a work of the United States Department of
247 Agriculture, Guidemaker is released to the public domain under a Creative Commons (CC0) public domain
248 attribution. The program is also available as a web application through the Cyverse discovery environment
249 [31], and as a stand-alone web application [27].

250

251 **Additional Files**

252 **Supplementary Figure 1.** Performance of GuideMaker for SaCas9 and StCas9.

253 **Supplementary Figure 2.** Performance of GuideMaker for SpCas9, SaCas9, and StCas9.

254 **Supplementary Figure 3.** Performance of GuideMaker with AVX2 settings.

255 **Supplementary Figure 4.** Memory usage of GuideMaker for SpCas9, SaCas9, and StCas9.

256 **Supplementary Table 1:** Organism features

257 **Supplementary Table 2:** Comparison of the average number of gRNA identified by GuideMaker and
258 CHOPCHOP.

259 **Supplementary Table 3:** Comparison of consensus ratio between GuideMaker and CHOPCHOP.

260

261 **List of abbreviations**

262 CAS: CRISPR-associated protein; CRISPR: Clustered Regularly Interspaced Short Palindromic Repeats;
263 gRNA: Guide RNA; HMSW: Hierarchical Navigable Small World; NMSLIB: Non-Metric Space Library;
264 PAM: Protospacer Adjacent Motif

265

266 **Funding**

267 The research was supported by the United States Department of Agriculture (USDA), Agricultural Research
268 Service (ARS) project number 6066-21310-005-D, and ARS cooperative agreement 6066-21310-005-28-S to
269 the University of Florida. This research used resources provided by the SCINet scientific computing initiative
270 of the USDA-ARS, ARS project number 0500-00093-001-00-D.

271

272 **Author Contributions**

273 R.P., L.T.R., C.R.R., and A.R.R. conceived and designed the study. R.P. and A.R.R developed and optimized
274 the software and performed the experiments. R.P., L.T.R., C.R.R., and A.R.R, tested the software, wrote, and
275 revised the manuscripts. All authors read and approved the final manuscript.

276

277 **References**

- 278 1. Jiang W, Bikard D, Cox D, Zhang F, Marraffini LA. RNA-guided editing of bacterial genomes using
279 CRISPR-Cas systems. *Nat Biotechnol.* 2013; doi: 10.1038/nbt.2508.
- 280 2. Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. Genome engineering using the CRISPR-Cas9
281 system. *Nat Protoc.* 2013; doi: 10.1038/nprot.2013.143.
- 282 3. Mojica FJM, Díez-Villaseñor C, García-Martínez J, Almendros C. Short motif sequences determine the
283 targets of the prokaryotic CRISPR defence system. *Microbiology.* 2009; doi: 10.1099/mic.0.023960-0.
- 284 4. Pickar-Oliver A, Gersbach CA. The next generation of CRISPR–Cas technologies and applications. *Nat*
285 *Rev Mol Cell Biol.* 2019; doi: 10.1038/s41580-019-0131-5.

- 286 5. Deltcheva E, Chylinski K, Sharma CM, Gonzales K, Chao Y, Pirzada ZA, et al.. CRISPR RNA maturation
287 by trans-encoded small RNA and host factor RNase III. *Nature*. 2011; doi: 10.1038/nature09886.
- 288 6. Ran FA, Cong L, Yan WX, Scott DA, Gootenberg JS, Kriz AJ, et al.. In vivo genome editing using
289 *Staphylococcus aureus* Cas9. *Nature*. 2015; doi: 10.1038/nature14299.
- 290 7. Zetsche B, Gootenberg JS, Abudayyeh OO, Slaymaker IM, Makarova KS, Essletzbichler P, et al.. Cpf1 Is a
291 single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell*. 2015; doi: 10.1016/j.cell.2015.09.038.
- 292 8. Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, et al.. Optimized sgRNA design
293 to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat Biotechnol*. 2016; doi:
294 10.1038/nbt.3437.
- 295 9. Hiranniramol K, Chen Y, Liu W, Wang X. Generalizable sgRNA design for improved CRISPR/Cas9
296 editing efficiency. Luigi Martelli P, editor. *Bioinformatics*. 2020; doi: 10.1093/bioinformatics/btaa041.
- 297 10. Xu H, Xiao T, Chen CH, Li W, Meyer CA, Wu Q, et al.. Sequence determinants of improved CRISPR
298 sgRNA design. *Genome Res*. 2015; doi: 10.1101/gr.191452.115.
- 299 11. Perez AR, Pritykin Y, Vidigal JA, Chhangawala S, Zamparo L, Leslie CS, et al.. GuideScan software for
300 improved single and paired CRISPR guide RNA design. *Nat Biotechnol*. 2017; doi: 10.1038/nbt.3804.
- 301 12. Anzalone A V., Koblan LW, Liu DR. Genome editing with CRISPR–Cas nucleases, base editors,
302 transposases and prime editors. *Nat Biotechnol*. 2020; doi: 10.1038/s41587-020-0561-9.
- 303 13. Hale CR, Zhao P, Olson S, Duff MO, Graveley BR, Wells L, et al.. RNA-guided RNA cleavage by a
304 CRISPR RNA-Cas protein complex. *Cell*. 2009; doi: 10.1016/j.cell.2009.07.040.
- 305 14. Abudayyeh OO, Gootenberg JS, Konermann S, Joung J, Slaymaker IM, Cox DBT, et al.. C2c2 is a single-
306 component programmable RNA-guided RNA-targeting CRISPR effector. *Science*. 2016; doi:
307 10.1126/science.aaf5573.

- 308 15. Ma E, Harrington LB, O'Connell MR, Zhou K, Doudna JA. Single-stranded DNA cleavage by divergent
309 CRISPR-Cas9 enzymes. *Mol Cell*. 2015; doi: 10.1016/j.molcel.2015.10.030.
- 310 16. Klompe SE, Vo PLH, Halpin-Healy TS, Sternberg SH. Transposon-encoded CRISPR–Cas systems direct
311 RNA-guided DNA integration. *Nature*. 2019; doi: 10.1038/s41586-019-1323-z.
- 312 17. Strecker J, Ladha A, Gardner Z, Schmid-Burgk JL, Makarova KS, Koonin E V, et al. RNA-guided DNA
313 insertion with CRISPR-associated transposases. *Science*. 2019; doi: 10.1126/science.aax9181.
- 314 18. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A programmable dual-RNA-
315 guided DNA endonuclease in adaptive bacterial immunity. *Science*. 2012; doi: 10.1126/science.1225829.
- 316 19. Hsu PD, Lander ES, Zhang F. Development and applications of CRISPR-Cas9 for genome engineering.
317 *Cell*. 2014; doi: 10.1016/j.cell.2014.05.010.
- 318 20. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, et al. Repurposing CRISPR as an
319 RNA-guided platform for sequence-specific control of gene expression. *Cell*. 2013; doi:
320 10.1016/j.cell.2013.02.022.
- 321 21. Cox DBT, Gootenberg JS, Abudayyeh OO, Franklin B, Kellner MJ, Joung J, et al. RNA editing with
322 CRISPR-Cas13. *Science*. 2017; doi: 10.1126/science.aag0180.
- 323 22. Yan WX, Chong S, Zhang H, Makarova KS, Koonin E V., Cheng DR, et al. Cas13d Is a compact RNA-
324 targeting type VI CRISPR effector positively modulated by a WYL-domain-containing accessory protein. *Mol*
325 *Cell*. 2018; doi: 10.1016/j.molcel.2018.02.028.
- 326 23. Dixit A, Parnas O, Li B, Chen J, Fulco CP, Jerby-Arnon L, et al. Perturb-Seq: Dissecting molecular
327 circuits with scalable single-cell RNA profiling of pooled genetic screens. *Cell*. 2016; doi:
328 10.1016/j.cell.2016.11.038.
- 329 24. Peters JM, Colavin A, Shi H, Czarny TL, Larson MH, Wong S, et al. A comprehensive, CRISPR-based
330 functional analysis of essential genes in bacteria. *Cell*. 2016; doi: 10.1016/j.cell.2016.05.003.

331 25. Zhu LJ. Overview of guide RNA design tools for CRISPR-Cas9 genome editing technology. *Front Biol.*
332 2015; doi: 10.1007/s11515-015-1366-y.

333 26. Cui Y, Xu J, Cheng M, Liao X, Peng S. Review of CRISPR/Cas9 sgRNA design tools. *Interdiscip Sci*
334 *Comput Life Sci.* 2018; doi: 10.1007/s12539-018-0298-z.

335 27. GuideMaker. The GuideMaker web app. <https://guidemaker.app.scinet.usda.gov>. Accessed 2021 May 27.

336 28. GuideMaker 2021. GuideMaker (Version 0.2.0). [https://github.com/USDA-ARS-](https://github.com/USDA-ARS-GBRU/GuideMaker/releases/tag/v0.2.0)
337 [GBRU/GuideMaker/releases/tag/v0.2.0](https://github.com/USDA-ARS-GBRU/GuideMaker/releases/tag/v0.2.0).

338 29. GuideMaker. The GuideMaker bioconda installation. <https://anaconda.org/bioconda/guidemaker>.

339 30. GuideMaker. The GuideMaker docker container. [https://github.com/orgs/USDA-ARS-](https://github.com/orgs/USDA-ARS-GBRU/packages?repo_name=GuideMaker)
340 [GBRU/packages?repo_name=GuideMaker](https://github.com/orgs/USDA-ARS-GBRU/packages?repo_name=GuideMaker).

341 31. Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, et al.. The iPlant collaborative:
342 cyberinfrastructure for enabling data to discovery for the life sciences. *PLOS Biol.* 2016; doi:
343 10.1371/journal.pbio.1002342.

344 32. GuideMaker. The GuideMaker project homepage. <https://guidemaker.org>. Accessed 2021 June 18.

345 33. Malkov YA, Yashunin DA. Efficient and robust approximate nearest neighbor search using hierarchical
346 navigable small world graphs. *IEEE Trans Pattern Anal Mach Intell.* 2020; doi:
347 10.1109/TPAMI.2018.2889473.

348 34. Naidan B, Boytsov L, Malkov Y, Novak D. Non-metric space library manual. 2015; doi:
349 <http://arxiv.org/abs/1508.05470>.

350 35. Labun K, Montague TG, Krause M, Torres Cleuren YN, Tjeldnes H, Valen E. CHOPCHOP v3:
351 Expanding the CRISPR web toolbox beyond genome editing. *Nucleic Acids Res.* 2019; doi:
352 10.1093/nar/gkz365.

- 353 36. Semenova E, Jore MM, Datsenko KA, Semenova A, Westra ER, Wanner B, et al.. Interference by
354 clustered regularly interspaced short palindromic repeat (CRISPR) RNA is governed by a seed sequence. Proc
355 Natl Acad Sci U S A. 2011; doi: 10.1073/pnas.1104144108.
- 356 37. Hsu PD, Scott DA, Weinstein JA, Ran FA, Konermann S, Agarwala V, et al.. DNA targeting specificity of
357 RNA-guided Cas9 nucleases. Nat Biotechnol. 2013; doi: 10.1038/nbt.2647.
- 358 38. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. Source code and command-line executables for
359 GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes. doi:
360 10.5281/zenodo.4849258.
- 361 39. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. Supporting material and analysis code for GuideMaker:
362 Software to design CRISPR-Cas guide RNA pools in non-model genomes. doi: 10.5281/zenodo.4898253.
- 363 40. Hirano S, Abudayyeh OO, Gootenberg JS, Horii T, Ishitani R, Hatada I, et al.. Structural basis for the
364 promiscuous PAM recognition by *Corynebacterium diphtheriae* Cas9. Nat Commun. 2019; doi: 10.1038/s41467-
365 019-09741-6.
- 366 41. Gleditsch D, Pausch P, Müller-Esparza H, Özcan A, Guo X, Bange G, et al.. PAM identification by
367 CRISPR-Cas effector complexes: diversified mechanisms and structures. RNA Biol. 2019; doi:
368 10.1080/15476286.2018.1504546.
- 369 42. Fu Y, Sander JD, Reyon D, Cascio VM, Joung JK. Improving CRISPR-Cas nuclease specificity using
370 truncated guide RNAs. Nat Biotechnol. 2014; doi: 10.1038/nbt.2808.
- 371 43. Wu X, Kriz AJ, Sharp PA. Target specificity of the CRISPR-Cas9 system. Sander JD, Joung JK, editors.
372 Quant Biol. 2014; doi: 10.1007/s40484-014-0030-x.

373

374

375

376 **Figure 1. Input parameters for GuideMaker**

377

378 **Figure 2. A typical workflow of GuideMaker:** 1) A user uploads the input genome (single or multiple) as
379 Genbank file, then defines the PAM sequence along with all the associated parameters and submits them to
380 run the program. 2) GuideMaker processes the input files and generates the interactive plots. Users can use
381 these interactive plots to explore the results and sort them by locus tag and genome coordinates. 3)
382 GuideMaker provides all the results and log files as downloads under the “Results” section.

383

384 **Figure 3. Entity Relationship Diagram showing the operation of the GuideMaker core program.**

385

386 **Figure 4. Performance of GuideMaker for SpCas9.** Evaluating the performance of GuideMaker across
387 three bacterial genomes using the “**NGG**” PAM motif with a target length of 20, unique zone of 11, 3prime
388 PAM orientation, before and into parameters of 500, knum of 10, controls of 10, and dist of 3. The mean of
389 10 runs was used for the evaluation, where dot and bar represent the mean and standard error, respectively.

Input	Description	Note/Example
Genome File	GuideMaker accepts one or more Genbank (.gbk or gzipped .gbk.gz) files with sequence data from a single genome as an input. GuideMaker extracts all the required information from the Genbank file to identify gRNAs and genomic features, allowing users to globally create gRNAs without preprocessed mapping files. Option: --genbank	<i>Carsonella_ruddii.gbk.gz</i> , <i>Carsonella_ruddii.gbk</i>
PAM	The Protospacer Adjacent Motif (PAM) is a short, generally 2-8 bp, sequence essential for binding by the Cas protein[3,40,41]. GuideMaker provides users the flexibility to define the PAM sequence for any Cas protein, enabling usage of new CRISPR-Cas systems. Degenerate PAM sequences are allowed. Option: --pamseq	NGG (SpCas9) NGRRT (SaCas9)
Restriction Enzymes	GuideMaker allows users to provide a list of defined or degenerate restriction site sequences to exclude from guides, which may be needed for some vector systems. Option: --restriction_enzyme_list.	GAATTC; Default: None
PAM Orientation	The PAM orientation parameter defines PAM position relative to the protospacer. Depending on the CRISPR-Cas system, the PAM could be 5' or 3' side of the guide sequence. For instance, SpCas9 recognizes 'NGG' PAM on the 3' end of the guide (i.e. 5'-[guide][pam]-3'), whereas the Cpf1 PAM is on the 5' end of the guide sequence (i.e. 5'-[pam][guide]-3'). To accommodate such differences, GuideMaker offers flexibility to define the PAM orientation. Option: --pam_orientation.	
Guide length	Guide length defines the length of gRNA. Changing the guide length allows the user to adjust the gRNA efficacy and specificity [42]. GuideMaker allows users to select the length of gRNA within 10-27 bp. Option: --guidelength.	
Length of seed region	The seed region is the guide sequence closest to the PAM recognition site, and the distal region is the region furthest from the PAM. For instance, if the guide length is 22bp, and the length of the seed region is 10, then the size of the seed and the distal regions is 10 and 12, respectively. It has been shown that the region close to PAM is sensitive [36,43], and non-uniqueness in this region can lead to off-target matches; however, the importance of the seed region is specific to the CRISPR-Cas system and the organism. Thus, GuideMaker allows the user to define the seed region with the maximum length of 27 bp; although, the length of the seed region must be less than or equal to the Guidelength. Additionally, the length of the seed region should not be too small because the total number of possible guides is limited to 4 raised to the power of the seed length. Option: --lsr.	
Hamming Distance	Hamming distance defines the number of substitutions required to turn one DNA sequence into another equal-length sequence. GuideMaker calculates the Hamming distance between all the candidate gRNAs and all sequences adjacent to a PAM site. gRNAs with a distance less than or equal to the user-defined value are considered too similar and removed to minimize off-targeting. Option: --dist	Options: [0 – 5]; Default: 2
Before	Before parameter allows user to select gRNAs that are upstream of a feature's start site. For example, if "before" is set to 100, each gRNA within 100 bp upstream of a feature will be retrieved. Option: --before	Options: [1 – 500]; Default: 100
Into	The into parameter allows the user to select gRNAs that are downstream of a feature's start. For example, if "into" is set to 100, each gRNA within 100 bp downstream of a feature will be retrieved. Option: --into.	Options: [1 – 500]; Default: 200
Similar guides	The number of sequences similar to the gRNA to include in the design report. Option: --knum	Options: [2 – 20]; Default: 3
Control gRNAs	Provides the set number of random control gRNAs. Option: --controls	Default: 1000

Upload one or more Genome file [.gbk, .gbk.gz]

Drag and drop files here
Limit 500MB per file • GBK, GZ

[Browse files](#)

Input PAM Motif [E.g. NGG] **1**

NGG

Restriction Enzymes[e.g. NGRT]:

NGRT

PAM Orientation [Options: 3prime, 5prime]

3prime

Guidelength [Options: 10 - 27]

20 - +

Length of seed region[Options: 0 - 27]

10 - +

Hamming Distance [Options: 0 - 5]

2 - +

Before [Options: 1 - 500]

100 - +

Into [Options: 1 - 500]

200 - +

Similar Guides[Options: 2 - 20]

3 - +

Control RNAs

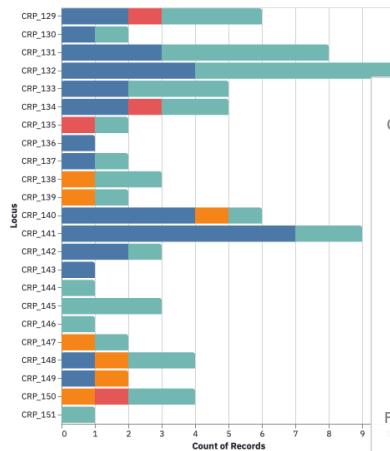
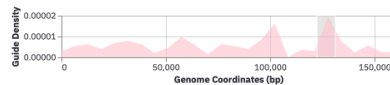
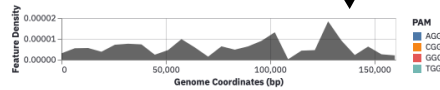
1000 - +

GuideMaker

Software to design CRISPR-Cas guide RNA pools in non-model genomes  

```
Running: 'guidemaker -i 69777e3d-da06-4414-90a3-42f5035feb8 -p NGG --
guidelength 20 --pam_orientation 3prime --lslr 10 --dist 2 --outdir 199f81c2-
bf42-11eb-ac6f-acde48001122 --log 199f81c2-bf42-11eb-ac6f-
acde48001122_log.txt --into 200 --before 100 --knum 3 --controls 1000 --
threads 2 --restriction_enzyme_list NGRT'
```

Accession: AP009180.1



2

Guide name: 8c758d7ab0babb1770874e4d064...

Guide sequence: TACAAAATATATATAATTA

GC: 0.05

Accession: AP009180.1

Guide start: 123916

Guide end: 123935

Guide strand: -

PAM: TGG

Feature id: fb10569bb9c3db0bdbcfefa55269f5...

Feature start: 123662

Feature end: 123916

Feature strand: -

Feature distance: 0

Similar guides: TTAACAGGAAATAACGGAAC;TC...

Similar guide 0;6;6

distances:

locus_tag: CRP_132

codon_start: 1

transl_table: 11

product: ribosomal protein L27

protein_id: BAF35163.1

db_xref: GI:116235315

Results

[Target Data](#)

[Control Data](#)

[Log File](#)

3

Parameter Dictionary +

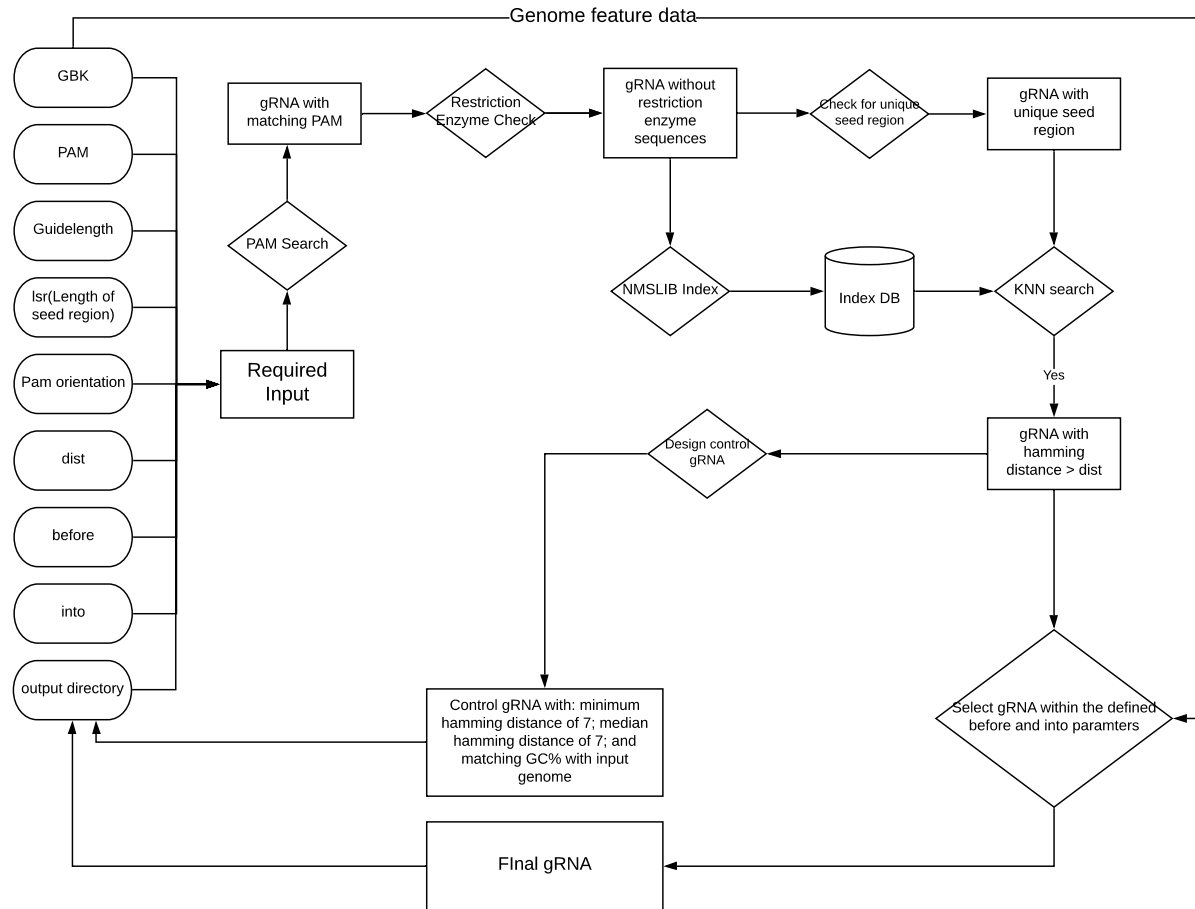
Designing Experiments with GuideMaker Results +

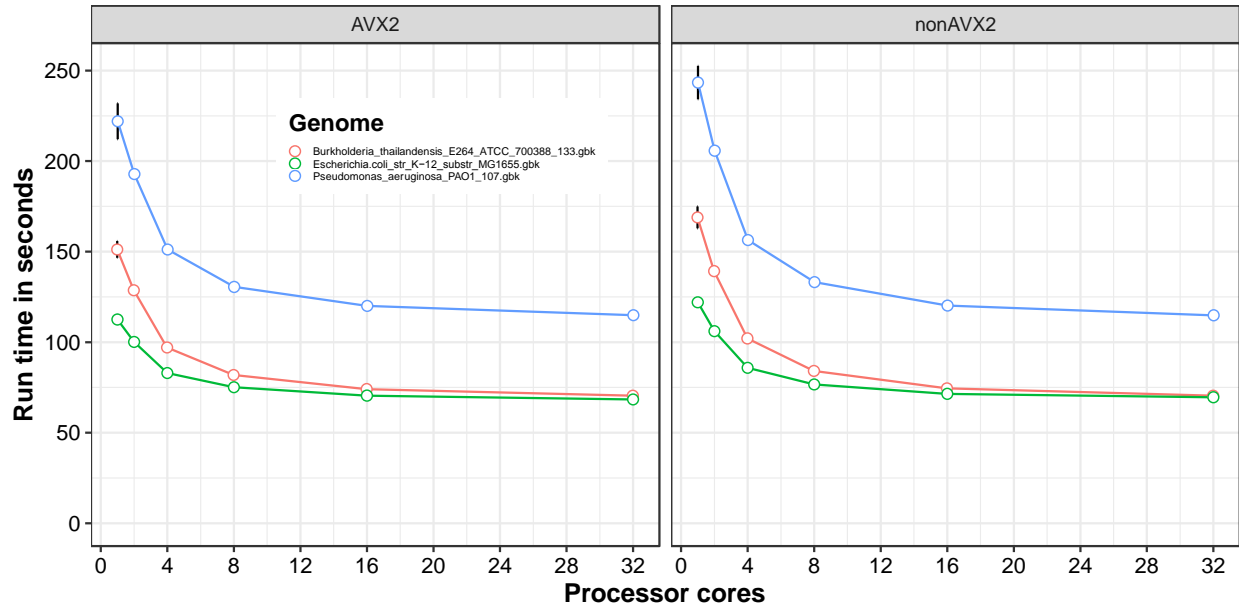
[API documentation](#)

API documentation for the module can be found [here](#)

License information

Guidemaker was created by the United States Department of Agriculture - Agricultural Research Service (USDA-ARS). As a work of the United States Government this software is available under the CC0 1.0 Universal Public Domain Dedication (CC0 1.0)







[Click here to access/download](#)

Supplementary Material

[Additional_Files_GigaScience_06182021.pdf](#)

