

GigaScience

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

Manuscript Number:	GIGA-D-21-00186R2	
Full Title:	GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes	
Article Type:	Technical Note	
Funding Information:	agricultural research service (6066-21310-005-D)	Dr. Adam R Rivers
	agricultural research service (6066-21310-005-28-S)	Dr. Christopher R. Reisch
	agricultural research service (0500-00093-001-00-D)	Dr. Adam R Rivers
Abstract:	<p>Background: CRISPR-Cas systems have expanded the possibilities for gene editing in bacteria and eukaryotes. There are many excellent tools for designing CRISPR-Cas guide RNAs for model organisms with standard Cas enzymes. GuideMaker is intended as a fast and easy-to-use design tool for challenging 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 using a degenerate protospacer adjacent motif (PAM) and a genome. The tool applies Hierarchical Navigable Small World (HNSW) graphs to speed up the comparison of guide RNAs and optionally provides on-target and off-target scoring. This allows the user to design effective 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 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:		
Response to Reviewers:	Dear Dr. Edmunds,	

We have completed minor revisions to the manuscript requested by the reviewers.

Specifically, we have:

Included the bio.tools identifier in the manuscript

Included the Scicrunch.org identifier in the manuscript

We have addressed these remaining comments:

“Maybe you intended to remove mentions of Cas13, but in the current version it still stands out. Page 3/line 66.”

Sorry, that last reference to Cas13 has now been removed.

“It is hard for me to believe that edit distance search for off-targets is equal to the hamming distance. This might be true for very small bacterial genomes, but for larger genomes (eg. human/mouse) this probably can't hold. It could also be that your implementation of the edit distance calculation for the guides could be flawed and therefore not reflecting the actuality. Consider adding tests for that "leven" option.”

We have addressed this in two ways:

We added a unit test (`test_levin_dist`) to the test code verifying that both Levenshtein and Hamming distance are being calculated as expected. This test code can be found here https://github.com/USDA-ARS-GBRU/GuideMaker/blob/main/tests/test_core.py#L319-L347

In that unit test we created a test sequence:

```
CGTAGCTAGTCACTAGCTGACAGCAAGGTTTTTCGTAGCTAGACACTAGCTGACA  
GCAAGGTTTTTCGTAGCTAGTCACTAGCTGACTAGCAAGG
```

That test sequence had three guide areas embedded in it (changes are shown with brackets and underscores):

1. CGTAGCTAG[T]CACTAGCTGACA_GCA|AGG
2. CGTAGCTAG[A]CACTAGCTGACA_GCA|AGG
3. CGTAGCTAG[T]CACTAGCTGACTAGCA|AGG

Guide 2 has 1 substitution (in brackets) and guide 3 has 1 insertion (underscore) relative to guide 1.

The Levenshtein distances for sequence 1 vs. [2, 3] are [1, 2], while the Hamming distances for sequence 1 vs. [2, 3] are [1,16].

The test code verifies that these edit distances are calculated correctly by the functions in Guidemaker. These edit distance calculations come directly from the highly-used NMSLIB library.

To address the concern that the guides designed with Leven and Hamming distance would diverge more for longer genomes, we tested the effect of using Levin and Hamming on the 537 MB genome of *Phaseolus vulgaris* (NC_023759). That data has been added to Supplementary Table 4.

Indeed, fewer guides were identical when Levin distance was used for the longer genomes, but the guides designed with Levin and Hamming were still 98% similar (versus 99.9% similar for *E. coli* MG 1655). For the larger *Phaseolus vulgaris* genome using Levin Distance with the “NGG” PAM took about twice as long, while for *E. coli* it took about 15x as long. This is likely because indexing, not distance computation, makes up a larger part of the compute time for larger genomes.

We agree that Levin distance the more biologically relevant measure of efficiency but think that for most users designing multiple guides per gene and working on smaller genomes the data supports the conclusion that Hamming is an appropriate distance approximation.

In the last revision we added Levin distance as an option for users who need it. We discuss the results in lines 233-242.

We have also added Supplementary Table 2 which summarizes the runtime to compute all guides for the PAMs “NGG”, “NNGRRT”, and “NNAGAAW” in the Homo sapiens (GRCh38.p13) genome. We added this benchmark for the large community of human researchers.

We have made additional improvements to the bibliography and abbreviation sections.

Sincerely,
Adam Rivers

Additional Information:

Question	Response
----------	----------

Are you submitting this manuscript to a special series or article collection?	No
---	----

<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>	Yes
--	-----

<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</p>	Yes
--	-----

Standards Reporting Checklist?	
<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 Standards Reporting Checklist?</p>	Yes

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 **ORCID IDs:**

9 Ravin Poudel: 0000-0003-2622-3889; Lidimarie T Rodriguez: 0000-0002-7611-6689; Christopher R Reisch:
10 0000-0002-6911-2452; Adam R Rivers: 0000-0002-3703-834X

11 **Abstract**

12 **Background:**

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

18 **Findings:**

19 GuideMaker can rapidly design gRNAs for gene targets across the genome using a degenerate protospacer
20 adjacent motif (PAM) and a genome. The tool applies Hierarchical Navigable Small World (HNSW) graphs
21 to speed up the comparison of guide RNAs and optionally provides on-target and off-target scoring. This
22 allows the user to design effective gRNAs targeting all genes in a typical bacterial genome in about 1-2
23 minutes.

24 **Conclusions:**

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

30

31 **Keywords** PAM, CRISPR-Cas, gRNA, Perturb-seq , Hierarchical Navigable Small World graph

32

33 **Introduction**

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

42 The *Streptococcus pyogenes* Cas9 (SpCas9) was the first Cas described [5] and it is still the most widely
43 used enzyme in CRISPR gene editing. Other Cas enzymes described early in the CRISPR revolution, such as
44 the *Streptococcus pyogenes* Cas9 and the *Acidaminococcus* Cas12a, are also commonly used [6,7]. Accordingly, the
45 parameters for these enzymes are often included in computational tools to identify CRISPR target sites [8–
46 11]. Cas9 enzymes from other organisms and other Cas-associated proteins that can cleave dsDNA, ssDNA,
47 ssRNA, and insert transposon elements have also been described and have their place in molecular toolkits

48 [12–18]. Each of these enzymes generally has specific requirements, such as PAM sequence constraints, PAM
49 orientation, and protospacer length. Many of these CRISPR-Cas systems have been repurposed to enable
50 molecular genetics techniques like gene deletions, gene insertions, transcriptional depletion and activation,
51 and translational repression [12,19–22]. Some of these techniques can be scaled to the genome level with
52 chip-synthesized oligonucleotides and pooled approaches to screening [23]. In pooled screens, high-
53 throughput DNA sequencing is used to identify how the pool has changed over time to elucidate genes that
54 affect cells' fitness in specific conditions. Given the diversity of the CRISPR systems and their uses,
55 identifying appropriate target sites is not trivial, especially for the number of targets needed for genome-scale
56 experiments.

57 Here we introduce GuideMaker, a computational tool to identify target sites and design gRNA
58 sequences that is not limited to any specific CRISPR system or organism. GuideMaker is most useful for a
59 few kinds of CRISPR experiments. The first use case is designing pools of gRNAs for genome-wide
60 screening experiments like Perturb-seq and CRISPR pool [23,24]. GuideMaker is optimized for making the
61 all-versus-all comparisons necessary to design a genome-wide screen and return candidate gRNAs for every
62 gene locus. The tool allows the user to filter targets based on their proximity to features of interest, like the
63 start codon for any coding sequence. The second major use case is for researchers working with non-model
64 organisms. Online gRNA design tools often have a limited number of preselected genomes available for
65 analysis because most methods require PAM site positions to be precomputed. GuideMaker rapidly computes
66 all guide positions on demand from user-provided GenBank files or a set of GFF/GTF (general feature
67 format/general transfer format) files and fasta files from any organism. The third use case is experiments with
68 Cas enzymes other than the canonical versions of Cas9 and Cas12a (Cpf1), that have atypical PAM and target
69 site requirements. GuideMaker allows the user to specify a custom PAM with variable length, including
70 degenerate nucleotides and allows the PAM to be on either the 3' or 5' side of the protospacer. These features
71 allow GuideMaker to support any current or future CRISPR-Cas system. Since the determination of which
72 CRISPR-Cas system functions best in any given organism is not predictable, this tool is highly relevant to
73 researchers developing CRISPR tools in new species. For SgCas9 GuideMaker also implements on-target and

74 off-target scoring from Doench et al. (2016). Because there is limited experimental data on most
75 Cas/organism combinations, cannot calculate target scoring for other Cas enzymes but instead uses design
76 heuristics that prioritize uniqueness in the seed region of the guide.

77

78 **Methods**

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

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

- 82 1. All the potential guides in a genome can be quickly designed in one run.
- 83 2. It can design gRNAs for any PAM sequence from any Cas system.
- 84 3. 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 4. Users can exclude specific restriction sites from guides to preserve those sites for downstream
88 experiments.
- 89 5. It creates control sequences based on the input genome. In CRISPR experiments it is often desirable
90 to create negative control sequences to evaluate off-target binding. GuideMaker provides the user
91 with realistic control gRNAs that are highly divergent from sequences adjacent to PAM sites.
- 92 6. It provides an option to select a subset of results by locus tags of interest.
- 93 7. It provides off-target Cutting Frequency Determination (CFD) scores for gRNAs [8].
- 94 8. Provides on-target efficacy score for canonical “NGG” PAM. These efficiency scores are based on
95 Azimuth algorithm[8].
- 96 9. Provides tabular result files which can be used for the design and ordering of gRNA pools.
- 97 10. Provides an interactive visualization and exploratory tool to evaluate the guides.

98 11. The software can be run as a web application [25], a CyVerse application, or a command-line
99 application [26]. Server code is included for running local instances of the web application as well.

100 A typical workflow of GuideMaker involves three major steps (Figure 2). In the first step, the user
101 uploads the input genome in one or more GenBank or GFF/GTF and fasta files (gzipped or uncompressed)
102 and defines the PAM and gRNA parameters (as highlighted in Figure 1). GuideMaker identifies and filters
103 target sites, then returns summary data to the graphical environment (Figure 2). Users can inspect the
104 interactive plots to learn more about the identified gRNAs and sort them by genome coordinates or locus tag.
105 In the final step, GuideMaker provides the results as downloadable files under the results section. These files
106 are used for synthesizing the guides. The command-line version of GuideMaker has similar input parameters
107 as the web application, with the flexibility to generate plots, configure the underlying hyper-parameters for
108 the Hierarchical Navigable Small World (HNSW) graph, filter the results by specific locus tag, select
109 Hamming or Levenshtein as the edit distance, predict on-target scores for “NGG” PAM, off-target CFD
110 scores, or to run the web application locally. To make the application easier to install we distribute the
111 application as a Bioconda environment[27], Docker container [28], Python package on Github [26], through
112 the CyVerse discovery environment [29] or as an online web application [25]. Detailed information on
113 accessing the software through various methods is available on the project homepage [30].

114 **Search method**

115 GuideMaker initially scans the genome, recording all candidate guide sequences adjacent to the
116 specified PAM sequence on both DNA strands (Figure 3). Candidate guides are then optionally checked for
117 the restriction sites. Next, the candidate guides are searched for a unique "seed region" closest to the PAM
118 site and candidate gRNAs that are not unique in their "seed region" are removed. Then, approximate nearest
119 neighbor search is used to remove candidate guides too similar to PAM adjacent sequences in the genome,
120 based on Hamming distance by default (the number of substitutions required to turn one DNA sequence into
121 another equal-length sequence). Users can also select Levenshtein distance in the command line version. The
122 approximate nearest neighbor search is performed using the Hierarchical Navigable Small World (HNSW)

123 graph method in the Non-Metric Space Library (NMSLIB) [31,32]. An index of all the initial candidate guides
124 is created using the selected edit distance. Each guide with a unique "seed region" is compared to all candidate
125 guides and any guides with edit distances below the user-set threshold are removed. This differs from the
126 standard procedure of indexing the genome and mapping each candidate guide against the whole genome
127 then parsing each result. HNSW has a search complexity of $\mathcal{O}(\log N)$ and index complexity of $\mathcal{O}(N \cdot$
128 $\log N)$ [31]. Finally, user-defined criteria are applied to specify the proximity and orientation of guides relative
129 to genomic features like genes. A list of guides is then returned to the user with relevant information about
130 the guide and its target genomic features.

131 The core of GuideMaker's search method is the HNSW method in NMSLIB [32]. The method
132 builds a multilayer graph index of the input data and has several parameters that can be optimized for index
133 building and search to trade-off speed and accuracy. Graph construction is the most time-consuming step in
134 our tests, and thus grid optimization was run to minimize run time while keeping recall above 99% relative to
135 the ground truth exact nearest-neighbor search. The grid-optimization parameters (M, efc, ef, and post) used
136 in the HNSW graph for approximate nearest neighbor search have been optimized for bacterial genomes. A
137 Jupyter notebook [33] script for re-optimization and visualization of these hyper-parameters is included in the
138 test directory of the command-line version of the software and optimized parameters can be passed to
139 GuideMaker with the *--config* flag.

140 **Target specificity**

141 Estimating the on-target and off-target performance of a guide requires experimental data, while this
142 is not available for most Cas systems it is available for SpCas9. Guidemaker re-implements two gRNA
143 scoring methods from [8] to provide on-target and off-target scoring for the common SpCas9 enzyme with
144 25 nt guides. The on-target scoring method is the Doench Rule Set 2 method, specifically the "Azimuth
145 Version 3 no position" model. The model applies boosted regression trees to nucleotide features. The
146 featurization script was rewritten and parallelized for increased speed and updated to Python 3. The original
147 Python Pickle model data object was converted to Open Neural Network Exchange (ONNX) format [34],

148 and parameters were moved to a JSON file for better reproducibility and security. GuideMaker uses the
149 ONNX Runtime [35] rather than Scikit-Learn [36] to make predictions from the model. For off-target scoring
150 GuideMaker calculates Cutting Frequency Distribution (CFD) scores using the scoring matrix from [8],
151 converted to JSON format for better reproducibility and security.

152

153 **Computational performance**

154 Genomes of different sizes, GC content, and chromosome numbers were used to test the speed and
155 scalability of GuideMaker (Supplementary Table 1). For benchmarking the performance, the same parameters
156 were used unless a specific parameter was being tested: a PAM motif of 'NGG', 3' pam orientation, target
157 length of 20, lsr (length of seed region) of 11, before and after parameters of 500, knum of 10, controls of 10,
158 dist of 3 and threads of 16. We profiled the performance of GuideMaker with different threads (1, 2, 4, 8, and
159 16) in processors with and without the AVX2 processor instruction set. The human genome was run with
160 separate parameters described in Supplementary Table 2. All tests were run on a single compute node with 2
161 x 24 core Intel Xeon® Platinum 8260 CPU @ 2.40 GHz with Cascade Lake microarchitecture. Three
162 bacterial genomes, a fungal genome, two plant genomes and a human genome were used in performance
163 benchmarking: *Escherichia coli* K12 (NC_000913), *Pseudomonas aeruginosa* PAO1 (NC_002516), *Burkholderia*
164 *thailandensis* E264 (NC_007651), *Aspergillus fumigatus* (NC_007194), *Arabidopsis thaliana* (NC_003070), *Phaseolus*
165 *vulgaris* (NC_023759), and *Homo sapiens* (GRCh38.p13). For the gene or locus-specific comparisons, only the
166 guides within the locus coordinates (i.e., zero feature distance) were considered.

167 **Comparison to existing design method**

168 We compared the results of GuideMaker with the results of the online and command-line versions of
169 CHOPCHOP (RRID:SCR_015723)[37]. GuideMaker and CHOPCHOP parameters were set to approximate
170 the same search. The length of the target sequence was set to 20 and zero mismatches were allowed in the
171 seed region (11nt) of the target. The *Escherichia coli* (str. K-12/MG1655) genome was used with the online
172 version of CHOPCHOP. Targets were searched in 40 Kbp increments to account for CHOPCHOP's online

173 size limitations. Target sequences were searched across multiple 40 Kbp segments of *E. coli* genome
174 (NC_000913.3:2001-42000, NC_000913.3:80001-120000, NC_000913.3:160001-200000,
175 NC_000913.3:240001-280000, and NC_000913.3:320001-360000). We also searched for target sequences
176 and genes/locus_tags within 40Kbp of (NC_000913.3:2001-42000) to compare identifications at the locus
177 level. The ratio between the tools was calculated by dividing the number of gRNA identified with
178 GuideMaker by the number of guides identified by CHOPCHOP to represent the proportion of guides
179 identified by both GuideMaker and CHOPCHOP.

180 The command-line version of CHOPCHOP was used to compare the memory usage and
181 computation time of CHOPCHOP and GuideMaker over an entire genome. The *E. coli* K-12 genome was
182 chosen for comparison because the precomputed 2bit genome files and Bowtie indexes were provided with
183 CHOPCHOP v 3. The matching GenBank file was downloaded for Guidemake and both programs were
184 run 5 times on the same machine using different numbers of processor cores [1, 2, 4, 8, 16].

185

186 **Results**

187 The time for GuideMaker to complete a typical run identifying all SpCas9 gRNAs (PAM 'NGG') in a
188 bacterial genome using 8 compute cores was 75 seconds for *E. coli* and 130 seconds for *P. aeruginosa* (Figure
189 4). For SaCas9 and StCas9, which have a longer PAM sequence (“NGRRT” and “NNAGAAW” respectively,
190 with 3' PAM orientation) and thereby fewer potential targets, the same genomes ran in 19 or 5 seconds
191 (Supplementary Figure 1). The fungus *Aspergillus fumigatus* (28MB) and the plants *Arabidopsis thaliana* (114 MB)
192 and *Phaseolus vulgaris* (537MB) have larger genomes but are still processed quickly. *A. fumigatus* processed
193 between 23-304 seconds, while *A. thaliana* processed in 250-921 and *P. vulgaris* processed in 333-4162 seconds
194 depending on the number of cores, AVX2 instructions, and PAM sequence (Supplementary Figure 2).
195 Guidemake designed guides for the entire human genome in 2-22 hours depending on the PAM used,
196 Supplementary Table 2.

197 GuideMaker can take advantage of Advanced Vector Extensions (AVX2) on newer x86 processors,
198 which improves the search speed because HNSW search is accelerated with AVX2 (Supplementary Figure 3).
199 The acceleration was larger when fewer processors were available (Supplementary Figure 3). The HNSW
200 algorithms are parallelized, and indexing-and-search takes most of the compute time in GuideMaker so the
201 software scales well when additional cores are added up to 8 cores (Supplementary Figure 3). In practice it
202 scaled up sub-linearly with genome size, globally estimating Cas9 guides for *E. coli* MG1655 (4.6MB) in 75
203 seconds and *Phaseolus vulgaris* (537MB) in 1549 seconds, both on 8 cores (Memory usage: 1.9GB for *E. coli* and
204 46.9GB for *P. vulgaris*, Supplementary Figure 4).

205 The results of GuideMaker were compared with the popular guide design software CHOPCHOP
206 version 3 [37]. When GuideMaker's filtering settings are set to match CHOPCHOP, the results are very
207 similar and 99.9% of the targets identified by GuideMaker fall within 2 nt of target coordinates returned by
208 CHOPCHOP. When GuideMaker's unique seed region criterion was not applied at the loci level, the average
209 number of guides identified by the two approaches was similar per locus (Mean GuideMaker = 116.8, Mean
210 CHOPCHOP = 113.6, p-value = 0.86, Supplementary Table 3). Although the number of guides identified
211 per gene locus differed, none of the genes were missed by either tool. GuideMaker's default requirement of a
212 unique seed region is more stringent than CHOPCHOP, and with it enabled, GuideMaker returned
213 (count=1787) 38.4% of the targets compared to CHOPCHOP (count=4651) over a 2Kbp-42Kbp test region
214 in *E. coli* str. K12 substr. MG1655. At the sequence level, 96.7% of the identified gRNA (1729/1787) from
215 both tools had identical sequences. The ratio of gRNA found by both the tools across the multiple 40Kbp
216 regions was 39.2% (sd= 1.9%, Supplementary Table 4) when using GuideMaker's more stringent default
217 settings. This ratio was calculated by dividing the number of gRNA from GuideMaker by the number from
218 CHOPCHOP for each 40Kb region. The effect of the stringent filtering heuristic used by GuideMaker was
219 investigated computationally by applying on target and off target scoring to the guides designed by
220 GuideMaker with and without the filtering heuristic (Supplementary Figure 5). As expected, the filtering
221 heuristic did not affect on-target scoring but did reduce the off-target CFD scores, suggesting that
222 GuideMaker heuristics could decrease off-target binding. This result remains to be validated experimentally.

223 The speed and memory usage of the command line versions of CHOPCHOP and Guidemaker were also
224 compared. When using 8 cores to process the *Escherichia coli* str. K-12 substr. MG165 genome, Guidemaker
225 was 65 times faster and used 2.7 times less memory than CHOPCHOP (Supplementary Figure 6).

226

227 **Discussion**

228 Designing gRNAs is a two-step process where GuideMaker first identifies potential guides adjacent to PAM
229 sequences and then filters the potential guides based on multiple criteria. The most important criterion is that
230 each guide has a minimum edit distance from any other sequence adjacent to a PAM site in the genome; this
231 decreases the likelihood of off-target binding. The second way GuideMaker reduces off-target binding is by
232 requiring that a set number of bases near the PAM site are unique from any other candidate guide. The 8
233 bases nearest the PAM are the most important for target specificity, and any mismatch is sufficient to prevent
234 binding [38,39]. The length of the unique region should be set with consideration for the size of the genome
235 since requiring short unique regions will limit the number of total guides that can be found. For example,
236 requiring that every gRNA be unique in the first 3 nt would only allow for $4^3 = 64$ possible guides to be
237 designed. For normal *--lvr* values of 9-12 this is only limiting for human-sized genomes and can be disabled by
238 setting *--lvr* to 0. All guides designed by GuideMaker are perfect matches to a single site in the genome.
239 Additional specificity is obtained by requiring all similar PAM-adjacent sequences to be unique in the critical
240 "seed region" and have a total number of mismatches that exceed the user-defined threshold. This double
241 criterion is expected to increase specificity.

242 The primary goal of the current version of our software is to support the design of gRNAs for non-
243 standard Cas enzymes or non-model organisms at the genome scale. Guide RNAs do not perform equally,
244 thus empirical experiments will be needed to fully validate the functionality and efficacy of gRNA predictions.
245 Given the similarity in targets identified by GuideMaker and CHOPCHOP, we anticipate that performance is
246 similar to the current state of the art but applicable to more design use cases. When a unique seed region and
247 edit distance-based filters were applied, GuideMaker created guides more conservatively, generating only

248 about 40% of the guides created by CHOPCHOP. While CHOPCHOP has an option to specify the
249 maximum number of mismatches in the first 9 nt or the whole guide, it does not allow the application of
250 both criteria. While there are differences in the number and position of guides generated by GuideMaker,
251 with GuideMaker being more conservative by default, both programs create enough guides to target nearly all
252 gene loci in the genome of *E. coli*. The current version of the GuideMaker provides options to predict off-
253 target CFD scores and on-target scores for the canonical “NGG” PAM. Both scoring approaches are based
254 on the publicly available models trained on empirical data with SpCas9. If experimentally validated data
255 become available from genome-wide screens with different Cas enzymes, future versions of GuideMaker
256 could potentially incorporate new scoring models to help rank candidate guides.

257 GuideMaker is a fast and flexible tool for designing guide RNA across the entire genome in non-
258 model organisms or with non-canonical Cas enzymes. It takes advantage of fast HNSW search to quickly
259 index and search new genomes. Several parameters can be tuned to ensure compatibility with the specific
260 application of the user. For example, GuideMaker checks the designed gRNA for a given restriction enzyme
261 site to prevent incompatibility with the cloning strategy. Second, the maximum distance from a target
262 sequence from the start of an annotated feature can be chosen to disrupt promoters or the beginning of the
263 coding sequence, since these sites are preferred for CRISPRi experiments. GuideMaker also creates off-target
264 control RNA sequences for use as negative controls in high-throughput experiments. Lastly, the program
265 plots the results for visual exploration of the targets and exports the data as .csv files. The software is
266 available as a command-line application, a web application, and is integrated into the CyVerse Discovery
267 Environment to provide users with a range of usage options. Guidemake is a fast, flexible design tool for the
268 creation of challenging guide RNA pools.

269

270 **Availability and Requirements**

271 Project name: GuideMaker

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

273 Operating system(s): Linux or macOS
274 Programming language: Python >=3.6
275 Other requirements:
276 License: CC0 1.0 Public Domain Dedication
277 RRID: SCR_021778
278 biotoolsID: guidemaker

279

280

281

282 **Competing Interests**

283 Authors declare no competing interests

284

285 **Data Availability**

286 The source code and command-line executables for GuideMaker are available and can be installed directly
287 from Github [26], Bioconda [27], or as a Docker container [28]. Data and code to reproduce the analysis in
288 the paper are available at Zenodo [40]. As a work of the United States Department of Agriculture,
289 GuideMaker is released to the public domain under a Creative Commons (CC0) public domain attribution.
290 The program is also available as a web application through the CyVerse discovery environment [29], and as a
291 stand-alone web application [25].

292

293 **Additional Data**

294 **Supplementary Figure 1.** Performance of GuideMaker for SaCas9 and StCas9 in selected bacteria

295 **Supplementary Figure 2.** Performance of GuideMaker for SpCas9, SaCas9, and StCas9 in selected
296 eukaryotes

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

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

299 **Supplementary Figure 5.** Comparison of efficiency and CFD scores with or without GuideMaker based
300 filters

301 **Supplementary Figure 6.** Performance and memory usage comparisons between CHOPCHOP (CLI
302 version) and GuideMaker

303 **Supplementary Table 1:** Organism features

304 **Supplementary Table 2:** Comparison of processing times and the number of gRNAs with different PAMs
305 in *Homo sapiens* (GRCh38.p13)

306 **Supplementary Table 3:** Comparison of the average number of gRNAs predicted by GuideMaker and
307 CHOPCHOP

308 **Supplementary Table 4:** Comparison of consensus ratio between GuideMaker and CHOPCHOP

309 **Supplementary Table 5:** Comparison of processing times and guide similarity for Levenshtein and
310 Hamming distances with different PAMs in *Escherichia coli* and *Phaseolus vulgaris*.

311

312 **List of abbreviations**

313 AVX2: Advanced Vector Extensions 2; bp: base pair; Cas: CRISPR-associated protein; Cas12a: CRISPR
314 associated protein 12a (previously known as Cpf1); CFD: Cutting Frequency Determination; Cpf1: See

315 Cas12a; CRISPR: Clustered Regularly Interspaced Short Palindromic Repeats; GFF: General Feature Format;
316 gRNA: Guide RNA; GTF: General Transfer Format; HNSW: Hierarchical Navigable Small World; kbp:
317 kilobase pairs; MB: megabases; NMSLIB: Non-Metric Space Library; nt : nucleotides; ONNX: Open Neural
318 Network Exchange; PAM: Protospacer Adjacent Motif; SaCas9: *Streptococcus aureus* CRISPR-associated protein
319 9; SpCas9: *Streptococcus pyogenes* CRISPR-associated protein 9

320 **Funding**

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

325

326 **Author Contributions**

327 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
328 the software and performed the experiments. R.P., L.T.R., C.R.R., and A.R.R, tested the software, wrote, and
329 revised the manuscripts. All authors read and approved the final manuscript.

330

331 **References**

- 332 1. Jiang W, Bikard D, Cox D, Zhang F, Marraffini LA. RNA-guided editing of bacterial genomes using
333 CRISPR-Cas systems. *Nat Biotechnol.* 2013; doi: 10.1038/nbt.2508.
- 334 2. Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. Genome engineering using the CRISPR-Cas9
335 system. *Nat Protoc.* 2013; doi: 10.1038/nprot.2013.143.
- 336 3. Mojica FJM, Díez-Villaseñor C, García-Martínez J, Almendros C. Short motif sequences determine the
337 targets of the prokaryotic CRISPR defence system. *Microbiology.* 2009; doi: 10.1099/mic.0.023960-0.

- 338 4. Pickar-Oliver A, Gersbach CA. The next generation of CRISPR–Cas technologies and applications. *Nat*
339 *Rev Mol Cell Biol.* 2019; doi: 10.1038/s41580-019-0131-5.
- 340 5. Deltcheva E, Chylinski K, Sharma CM, Gonzales K, Chao Y, Pirzada ZA, et al.. CRISPR RNA maturation
341 by trans-encoded small RNA and host factor RNase III. *Nature.* 2011; doi: 10.1038/nature09886.
- 342 6. Ran FA, Cong L, Yan WX, Scott DA, Gootenberg JS, Kriz AJ, et al.. In vivo genome editing using
343 *Staphylococcus aureus* Cas9. *Nature.* 2015; doi: 10.1038/nature14299.
- 344 7. Zetsche B, Gootenberg JS, Abudayyeh OO, Slaymaker IM, Makarova KS, Essletzbichler P, et al.. Cpf1 Is a
345 single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell.* 2015; doi: 10.1016/j.cell.2015.09.038.
- 346 8. Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, et al.. Optimized sgRNA design
347 to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat Biotechnol.* Nature Publishing
348 Group; 2016; doi: 10.1038/nbt.3437.
- 349 9. Hiranniramol K, Chen Y, Liu W, Wang X. Generalizable sgRNA design for improved CRISPR/Cas9
350 editing efficiency. Luigi Martelli P, editor. *Bioinformatics.* Oxford University Press; 2020; doi:
351 10.1093/bioinformatics/btaa041.
- 352 10. Xu H, Xiao T, Chen CH, Li W, Meyer CA, Wu Q, et al.. Sequence determinants of improved CRISPR
353 sgRNA design. *Genome Res.* 2015; doi: 10.1101/gr.191452.115.
- 354 11. Perez AR, Pritykin Y, Vidigal JA, Chhangawala S, Zamparo L, Leslie CS, et al.. GuideScan software for
355 improved single and paired CRISPR guide RNA design. *Nat Biotechnol.* 2017; doi: 10.1038/nbt.3804.
- 356 12. Anzalone A V., Koblan LW, Liu DR. Genome editing with CRISPR–Cas nucleases, base editors,
357 transposases and prime editors. *Nat Biotechnol.* Springer US; 2020; doi: 10.1038/s41587-020-0561-9.
- 358 13. Hale CR, Zhao P, Olson S, Duff MO, Graveley BR, Wells L, et al.. RNA-Guided RNA Cleavage by a
359 CRISPR RNA-Cas Protein Complex. *Cell.* Cell; 2009; doi: 10.1016/j.cell.2009.07.040.
- 360 14. Abudayyeh OO, Gootenberg JS, Konermann S, Joung J, Slaymaker IM, Cox DBT, et al.. C2c2 is a single-

361 component programmable RNA-guided RNA-targeting CRISPR effector. *Science*. American Association for
362 the Advancement of Science; 2016; doi: 10.1126/science.aaf5573.

363 15. Ma E, Harrington LB, O'Connell MR, Zhou K, Doudna JA. Single-stranded DNA cleavage by divergent
364 CRISPR-Cas9 enzymes. *Mol Cell*. Cell Press; 2015; doi: 10.1016/j.molcel.2015.10.030.

365 16. Klompe SE, Vo PLH, Halpin-Healy TS, Sternberg SH. Transposon-encoded CRISPR–Cas systems direct
366 RNA-guided DNA integration. *Nature*. Nature Publishing Group; 2019; doi: 10.1038/s41586-019-1323-z.

367 17. Strecker J, Ladha A, Gardner Z, Schmid-Burgk JL, Makarova KS, Koonin E V, et al.. RNA-guided DNA
368 insertion with CRISPR-associated transposases. *Science (80-)*. American Association for the Advancement of
369 Science; 2019; doi: 10.1126/science.aax9181.

370 18. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A programmable dual-RNA-
371 guided DNA endonuclease in adaptive bacterial immunity. *Science (80-)*. 2012; doi: 10.1126/science.1225829.

372 19. Hsu PD, Lander ES, Zhang F. Development and applications of CRISPR–Cas9 for genome engineering.
373 *Cell*. 2014; doi: 10.1016/j.cell.2014.05.010.

374 20. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, et al.. Repurposing CRISPR as an
375 RNA-guided platform for sequence-specific control of gene expression. *Cell*. 2013; doi:
376 10.1016/j.cell.2013.02.022.

377 21. Cox DBT, Gootenberg JS, Abudayyeh OO, Franklin B, Kellner MJ, Joung J, et al.. RNA editing with
378 CRISPR-Cas13. *Science (80-)*. 2017; doi: 10.1126/science.aaq0180.

379 22. Yan WX, Chong S, Zhang H, Makarova KS, Koonin E V., Cheng DR, et al.. Cas13d Is a compact RNA-
380 targeting type VI CRISPR effector positively modulated by a WYL-domain-containing accessory protein. *Mol*
381 *Cell*. 2018; doi: 10.1016/j.molcel.2018.02.028.

382 23. Dixit A, Parnas O, Li B, Chen J, Fulco CP, Jerby-Arnon L, et al.. Perturb-Seq: dissecting molecular
383 circuits with scalable single-cell RNA profiling of pooled genetic screens. *Cell*. 2016; doi:

384 10.1016/j.cell.2016.11.038.

385 24. Peters JM, Colavin A, Shi H, Czarny TL, Larson MH, Wong S, et al.. A comprehensive, CRISPR-based
386 functional analysis of essential genes in bacteria. *Cell*. Elsevier; 2016; doi: 10.1016/j.cell.2016.05.003.

387 25. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. 2022. The GuideMaker web application.
388 <https://guidemaker.app.scinet.usda.gov>. Accessed 06 Jan 2022.

389 26. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. (2021). GuideMaker (Version 0.3.4). Zenodo.
390 <https://doi.org/10.5281/zenodo.5655842>.

391 27. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. (2022). The GuideMaker Bioconda installation (version
392 0.3.4) . <https://anaconda.org/bioconda/guidemaker>.

393 28. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. (2022). The GuideMaker Docker container.
394 <https://github.com/USDA-ARS-GBRU/GuideMaker/releases>.

395 29. Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, et al.. The iPlant collaborative:
396 cyberinfrastructure for enabling data to discovery for the life sciences. *PLOS Biol*. 2016; doi:
397 10.1371/journal.pbio.1002342.

398 30. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. 2022. The GuideMaker project homepage.
399 <https://guidemaker.org>. Accessed 06 Jan 2022.

400 31. Malkov YA, Yashunin DA. Efficient and robust approximate nearest neighbor search using Hierarchical
401 Navigable Small World graphs. 2016; arXiv <https://arxiv.org/abs/1603.09320>

402 32. Naidan B, Boytsov L, Malkov Y, Novak D. Non-metric space library manual. 2015; aXiv
403 <http://arxiv.org/abs/1508.05470>.

404 33. Kluyver T, Ragan-Kelley B, Pérez F, Granger B, Bussonnier M, Frederic J, et al.. Jupyter Notebooks - a
405 publishing format for reproducible computational workflows. In: Loizides F, Schmidt B, editors. *Position Power*
406 *Acad Publ Play Agents Agendas*. Netherlands: IOS Press; p. 87–90.

- 407 34. Microsoft. (2021). ONNX (Version 1.10.0) <https://github.com/onnx/onnx/releases/tag/v1.10.0>
- 408 35. Microsoft. (2021). ONNX Runtime (Version
409 1.8.1) <https://github.com/microsoft/onnxruntime/releases/tag/v1.8.1>
- 410 36. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al.. Scikit-learn: Machine
411 Learning in Python. *J Mach Learn Res.* 12:2825–302011;
- 412 37. Labun K, Montague TG, Krause M, Torres Cleuren YN, Tjeldnes H, Valen E. CHOPCHOP v3:
413 Expanding the CRISPR web toolbox beyond genome editing. *Nucleic Acids Res.* 2019; doi:
414 10.1093/nar/gkz365.
- 415 38. Semenova E, Jore MM, Datsenko KA, Semenova A, Westra ER, Wanner B, et al.. Interference by
416 clustered regularly interspaced short palindromic repeat (CRISPR) RNA is governed by a seed sequence. *Proc*
417 *Natl Acad Sci U S A.* 2011; doi: 10.1073/pnas.1104144108.
- 418 39. Hsu PD, Scott DA, Weinstein JA, Ran FA, Konermann S, Agarwala V, et al.. DNA targeting specificity of
419 RNA-guided Cas9 nucleases. *Nat Biotechnol.* 2013; doi: 10.1038/nbt.2647.
- 420 40. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. (2022). Code to reproduce analyses in the GuideMaker
421 paper (version 0.4.0). Zenodo. <https://doi.org/10.5281/zenodo.5825817>

422

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

424

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

430

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

432

433 **Figure 4. Performance of GuideMaker for SpCas9.** Evaluating the performance of GuideMaker across
434 three bacterial genomes using the “**NGG**” PAM motif with a target length of 20, unique zone of 11, 3prime
435 PAM orientation, before and into parameters of 500, knum of 10, controls of 10, and dist of 3. The mean of
436 10 runs was used for the evaluation, where dot and bar represent the mean and standard error, respectively.

Inputs	Descriptions	Notes/Examples
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	E.g. <i>Carsonella_ruddii</i> .gbk.gz, <i>Carsonella_ruddii</i> .gbk
Fasta File	One or more fasta or gzipped fasta files for a single genome. If using a fasta, a GFF/GTF file must also be provided but not a genbank file. Option: --fasta	E.g. <i>Carsonella_ruddii</i> .fasta
Gff File	One or more GFF or GTF files (optionally gzipped) for a single genome. If using a GFF/GTF a fasta file must also be provided but not a genbank file. Option: --gff	E.g. <i>Carsonella_ruddii</i> .gff
PAM	The Protospacer Adjacent Motif (PAM) is the 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	E.g. NGG (SpCas9) NGRRT (SaCas9)
Restriction Enzymes	It can be useful to avoid sequences with restriction endonuclease recognition sites for used cloning guide library. GuideMaker allows users to provide a list of defined or degenerate restriction site sequences to avoid targeting. Option: --restriction_enzyme_list.	E.g. NGRT; Default: None
PAM Orientation	The PAM orientation parameter defines PAM position relative to the protospacer. Depending on the CRISPR-Cas system, the orientation of PAM could be 5' or 3' to 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.	
Guidelength	Guidelength 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. GuideMaker divides each guide into the seed and distal regions (Figure A and B). 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.	
Edit Distance	Edit distance defines the number of substitutions required to turn one DNA sequence into another sequence. GuideMaker calculates the pairwise edit 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
Distance type	Defines the edit distance type. GuideMaker provides two edit distance type: hamming ; and leven. Option: -- dtype	Options:[hamming, leven]; Default hamming
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
Locus tag	List of locus tag for subsetting the final output so the gRNA specific to the listed locus tag are retrieved. Option: --filter_by_locus	Default: None
CFD score	Cutting Frequency Determination (CFD) score for accessing off-target activity of gRNAs. Option: --cfd_score	Default: None
Efficiency score	On-target efficiency score predicted based on Azimuth 2.0.– only for NGG PAM. Option: --doench_efficiency_score	Default: None
Similar guides	Retrieves the number of sequences similar to the gRNA. Option: --knum	Options: [2 – 20]; Default: 3
Control gRNAs	Provides the set number of random control gRNAs. Option: --controls	Default: 1000

Select Parameters to Design gRNAs

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

Drag and drop files here
Limit 200MB per file • GBK, GZ, GBFF

1

Upload one or more fasta file [.fasta, .fasta.gz]

Drag and drop files here
Limit 200MB per file • FASTA, GZ, FNA

Upload gff/gtf file if you are using fasta [.gff, .gtf]

Drag and drop files here
Limit 200MB per file • GFF, GTF

OR Use Demo GBK

Carsonella_ruddii.gbk.gz

Input PAM Motif [E.g. NGG]

NGG

Restriction Enzymes[e.g. NGRT]:

Enter to add more

PAM Orientation [Options: 3prime, 5prime]

3prime

Guidelength [Options: 10 - 27]

20

Length of seed region[Options: 0 - 27]

10

Edit Distance [Options: 0 - 5]

2

Before [Options: 1 - 500]

100

Into [Options: 1 - 500]

200

Similar Guides[Options: 2 - 20]

3

Control RNAs

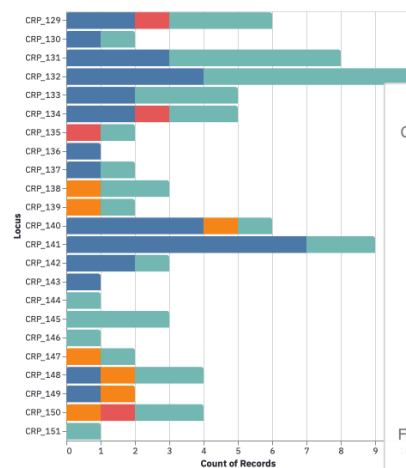
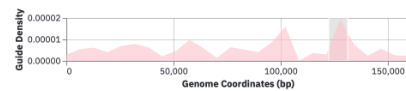
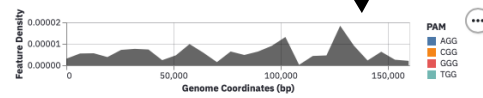
10

GuideMaker

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

```
Running:: 'guidemaker -i 69777e3d-da06-4414-90a3-42f5035febf8 -p NGG --
guidelength 20 --pam_orientation 3prime --lrs 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



Results

[Target Data](#)

[Control Data](#)

[Log File](#)

Parameter Dictionary

Designing Experiments with GuideMaker Results

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

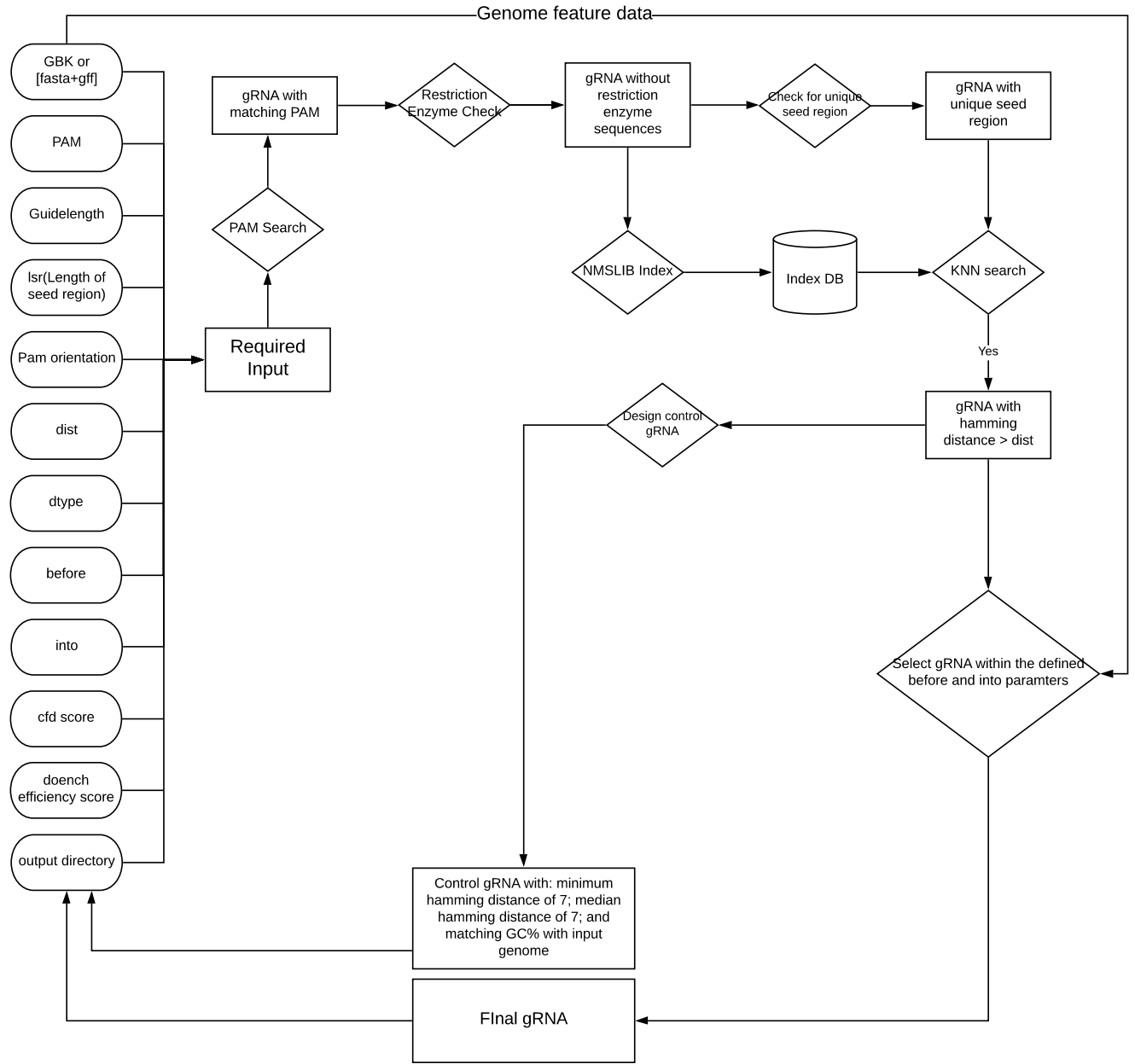
db_xref: GI:116235315

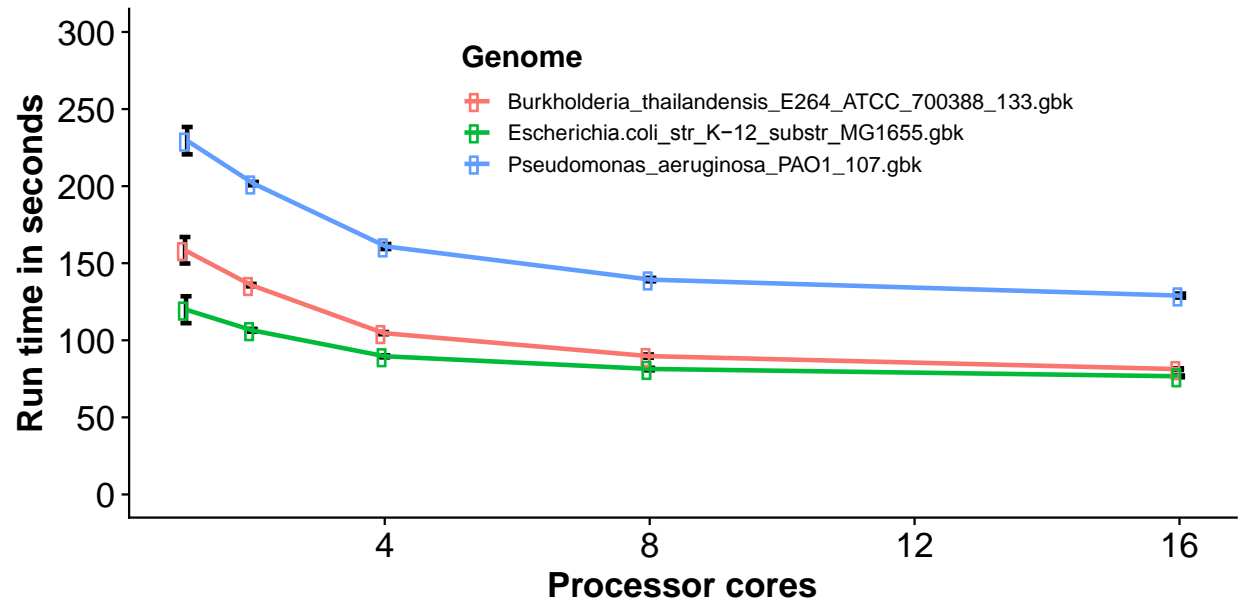
API documentation

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)



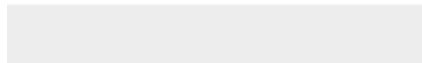




[Click here to access/download](#)

Supplementary Material

[Additional_Files_GigaScience_after_accepted.docx](#)



Dear Dr. Edmunds,

We have completed minor revisions to the manuscript requested by the reviewers.

Specifically, we have:

Included the bio.tools identifier in the manuscript

Included the Scicrunch.org identifier in the manuscript

We have addressed these remaining comments:

“Maybe you intended to remove mentions of Cas13, but in the current version it still stands out. Page 3/line 66.”

Sorry, that last reference to Cas13 has now been removed.

“It is hard for me to believe that edit distance search for off-targets is equal to the hamming distance. This might be true for very small bacterial genomes, but for larger genomes (eg. human/mouse) this probably can't hold. It could also be that your implementation of the edit distance calculation for the guides could be flawed and therefore not reflecting the actuality. Consider adding tests for that "leven" option.”

We have addressed this in two ways:

We added a unit test (test_levin_dist) to the test code verifying that both Levenshtein and Hamming distance are being calculated as expected. This test code can be found here https://github.com/USDA-ARS-GBRU/GuideMaker/blob/main/tests/test_core.py#L319-L347

In that unit test we created a test sequence:

```
CGTAGCTAGTCACTAGCTGACAGCAAGG TTTTTCGTAGCTAGACACTAGCTGACAGCAAGG TTTTTCGTAGCTA  
GTCAGCTAGCTGACTAGCAAGG
```

That test sequence had three guide areas embedded in it (changes are shown with brackets and underscores):

1. CGTAGCTAG[T]CACTAGCTGACA_GCA|AGG
2. CGTAGCTAG[A]CACTAGCTGACA_GCA|AGG
3. CGTAGCTAG[T]CACTAGCTGACTAGCA|AGG

Guide 2 has 1 substitution and guide 3 has 1 insertion relative to guide 1.

The Levenshtein distances for sequence 1 vs. [2, 3] are [1, 2], while the Hamming distances for sequence 1 vs. [2, 3] are [1,16].

The test code verifies that these edit distances are calculated correctly by the functions in Guidemaker.

To address the concern that the guides designed with Leven and Hamming distance would diverge more for longer genomes, we tested the effect of using Levin and Hamming on the 537 MB genome of *Phaseolus vulgaris* (NC_023759). That data has been added to Supplementary Table 4.

Indeed, fewer guides were identical when Levin distance was used for the longer genomes, but the guides designed with Levin and Hamming were still 98% similar (versus 99.9% similar for *E coli*. MG 1655). For the larger *Phaseolus vulgaris* genome using Levin Distance with the “NGG” PAM took about twice as long, while. for *E coli* it took about 15x as long. This is likely because indexing, not distance computation, makes up a larger part of the compute time for larger genomes.

We agree that Levin distance the more biologically relevant measure of efficiency but think that for most users designing multiple guides per gene and working on smaller genomes the data supports the conclusion that Hamming is an appropriate distance approximation.

In the last revision we added Levin distance an an option for users who need it. We discuss the results in lines 233-242.

We have also added Supplementary Table 2 which summarizes the runtime to compute all guides for the PAMs “NGG”, “NNGRRT “, and “NNAGAAW” in the *Homo sapiens* (GRCh38.p13) genome. We added this benchmark for the large community of human researchers.

We have made additional improvements to the bibliography and abbreviation sections.

Sincerely,

Adam Rivers