

# GigaScience

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

<b>Manuscript Number:</b>	GIGA-D-21-00186R1	
<b>Full Title:</b>	GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes	
<b>Article Type:</b>	Technical Note	
<b>Funding Information:</b>	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
<b>Abstract:</b>	<p><b>Background:</b> 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><b>Findings:</b> 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><b>Conclusions:</b> 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>	
<b>Corresponding Author:</b>	Adam R Rivers, Ph.D. USDA Agricultural Research Service Gainesville, FL UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	USDA Agricultural Research Service	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Ravin Poudel, Ph.D.	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Ravin Poudel, Ph.D.	
	Lidimarie Trujillo Rodriguez, B.S.	
	Christopher R. Reisch, Ph.D.	
	Adam R Rivers, Ph.D.	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	GIGA-D-21-00186 GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model	

genomes

Ravin Poudel; Lidimarie Trujillo Rodriguez; Christopher R. Reisch; Adam R Rivers  
GigaScience

October 22, 2021

Dear Dr. Edmunds:

Thank you for synthesizing the comments from the reviewers, I found them fair and constructive and have implemented or addressed all suggestions, including feature requests and requests for additional data. For ease of review, I have placed all the reviewer comments and my responses in tabular format and attached that document as well. We have also registered GuideMaker with bio.tools (<https://bio.tools/guidemaker>) and SciCrunch.org (SCR\_021778) and included those identifiers in the paper.

Sincerely,

Adam Rivers

Point-by-point comments:

Reviewer 1

1. I tested the website and the tool, not finding any bugs and errors. Website is well made, congratulations!

Thanks.

2. Name of the tool: GuideMaker is not self-explanatory for what it is specialized for, which is pooled design. In the future consider naming your tools more distinctly as I am afraid that currently the tool will be buried under hundreds of other GuideSomething tools.

This is a good point. At this point we have a domain, website, preprint and users of the software so it would be pretty disruptive to change, but we will be more specific with future names.

Authors also claim to support Cas13 (page 3 line 65), but don't mention anything more specific about it. I mention that because design for RNA is vastly different from design for DNA and it should be explained how the tool designs for RNA.

We have removed mention of Cas13 since it was not evaluated for this application.

From my understanding the tool offers highly discriminatory settings towards off-target search for a quick resolution of the all vs all comparison problem, however authors ignore that CRISPR off-targets are not defined by the hamming distance, but levenshtein distance. This was proven already by many studies e.g. Tsai et al. 2015. I recommend that authors embrace this issue in the paper and explain why their design may be suitable, and for what kind of studies it would be alright to use hamming distance vs levenshtein distance instead of ignoring the problem.

We added an option to use Levenshtein distance to the command-line version of GuideMaker. We also evaluated the effect of using each distance metric on the guides selected there was virtually no difference in the guides selected for the bacterial genome tested. Results are reported in the manuscript on lines 127-130 and in Supplementary Table 4. We suspect this is because there were not many of indels in guides from bacterial genomes. Because hamming is much faster and gave equivalent results, we have kept it as the default distance metric.

5. Study could gain prominence by showing a couple figures and describing how the grid-optimization parameters were selected. This would be especially important for everyone that wants to use this tool for nonbacterial genomes (page 6, lines 128-131). Although script for optimization is included, it would be good to see what are the tradeoffs.

We have added graphical outputs to the grid optimization Jupyter notebook, along with instructions on how to run it for other genomes. We explain how to use this in the paper.

I believe that Figure 4 and all other AVX2 vs nonAVX2 comparisons are not interesting enough to include multiple times. AVX2 improvements are nice, but the tool is already plenty fast, and running time of 250 vs 220 seconds does not matter for normal users.

We have removed the redundant non-AVX figures since most processors now support AVX2. The AVX2 performance gains were larger previously but the NMSlib library improved its non-AVX performance so there is not much difference anymore. Supp fig 3 summarizes the effects of AVX2.

Similarly the number of cores does not seem to influence tool speed above 8 cores and one figure should be enough to explain that. We removed additional cores above 16, but retained 16 to show the flattening out in performance between 8 and 16.

Tool claims very fast running times, but does not compare to the running times of other similar tools for the design of the pooled screens, this could highlight its superiority. We now compare the tools to the command-line version of CHOPCHOP using *E. coli* in Supplementary Fig 5. Despite using precomputed mappings for CHOPCHOP, Guidemaker is about 100x faster and uses 60% less memory

CHOPCHOP is a general tool for the design of pooled screens while here it is used as a pooled screen tool due to its configurability. Additionally, CHOPCHOP also supports all PAM and all species, but on its python version available. That is a good point we have used the CLI version of CHOPCHOP for the added comparisons.

Comparisons to CHOPCHOP focus on the guides found, but I don't understand why consensus ratio between the tools should matter. What is more important is whether GuideMaker does indeed not filter any guides that are preferable for each gene (e.g. by CHOPCHOP ranking) and whether its hamming based filter is good enough to not cause significant unknown off-target effects (levenshtein distance off-targets not found by hamming distance filter). All it takes is one bulge and the hamming distance will become large, while levenshtein distance can even be as low as 1.

We used CHOPCHOP consensus because it is widely used and there is not a gold-standard ground truth data for this. Guidemaker reports about the same number of targets when these same filtering metrics are applied. We have also added Doench et al. 2016 scoring (Azimuth) and CFD scoring to evaluate on target and off target guides for Cas9, so user can sort output by these scores.

We added Supp fig 6 to show the effect of selection on On-Target and Off target scores filtered with GuideMaker parameters or unfiltered like CHOPCHOP. Our Isr filtering does not affect on-target scoring but does reduce off target scoring slightly. Our testing has revealed that using Levenshtein distance does not affect guide selection (see explanation above)

It is not clear to me why the tool can't be used with large genomes, filtering on the 11bp seed and hamming distance should be plenty fast for also very large genomes. It can be used for larger genomes just that HNSW loses its speed advantage around at around 1E9 guides. HNSW starts out much faster than indexing and searching the whole genome conventionally but the time per query grows more slowly for the conventional methods. Eventually it becomes faster to use conventional search rather than HNSW search.

Could it be that the tool should support other input, not only genbank file format? We have added support for importing sequence and annotation from GFF/GTF and Fasta files.

Reviewer 2

The author developed a software, GuideMaker, for designing CRISPR-Cas guide RNA pools in non-model genomes. Three bacterial genomes, a fungal genome, and a plant genome were used in performance benchmarking, which proves that the software supports the design of gRNAs in non-standard Cas enzymes for non-model organisms at the genome-scale. However, the advantages of this software are not well estimated nor presented compared to other tools like CHOPCHOP. We have improved our explanation of the advantages of GuideMaker relative to CHOPCHOP for its intended applications, including a performance evaluation in Supplementary figure 5 and a better explanation. We have also added both on-target and off-target scoring for NGG PAMs (the only PAMs for which training data is available), from Doench et al. 2016. Also, the software was mainly evaluated in three bacteria genomes, one fungus and Arabidopsis genome. There are no tests for non-model plant or animal genomes. Therefore, the "non-model genomes" in the title are exaggerated. I list more problems as follows.

We have added the genome of the 537 MG plant *Phaseolus vulgaris*. Our assertion that Guidemaker can be used for non-model organisms comes from the fact that it does not require precomputed reference genomes but rather computes guide pools quickly on the fly. This feature of the software can be shown without necessarily the

genomes of obscure organism. We have clarified this confusing part in the since Pseudomonas and Arabidopsis certainly are model organisms.  
The authors did not compare the computation resources and performance (running time, memory) with existing softwares like CHOPCHOP. Also, the authors need to compare the score rankings with CHOPCHOP to present the relative power of GuideMaker. Is there any score rankings concerning efficiency or off-target possibilities for the designed Guide RNAs This is a good suggestion; we have added Supp. Fig 5 that looks at the time and memory requirements for Guidemake and CHOPCHOP CLI.

We have added the same on target and off target ranking algorithms used by CHOPCHOP V3. Those algorithms are Azimuth and CDF from Doench et al. (2016).  
2. It is better to add support for gff formatted annotation input files since many non-model species do not have GenBank annotations. We have added support for importing sequence and annotation from GFF/GTF and Fasta files.

3. The authors mentioned GuideMaker can design gRNAs for any small to medium size genome (up to about 500 megabases). The maximum genome used in the article was Arabidopsis thaliana (114.1MB), which is obviously smaller than the described (up to about 500 megabases). We couldn't find the description whether the authors had investigated the larger genomes. Therefore, the detailed analysis or discussion of this problem is needed.

We have added the 537 MB Phaseolus vulgaris genome in Supp. Fig 4 to demonstrate this claim.

4. The authors stated GuideMaker to design CRISPR-Cas guide RNA pools in non-model genomes. Arabidopsis thaliana is a model organism and test in a non-model plant genome will be highly valuable.

We have added the genome of the 537 MG plant Phaseolus vulgaris. Our assertion that Guidemake can be used for non-model organisms comes from the fact that it does not require precomputed reference genomes but rather computes guide pools quickly on the fly. We have clarified this confusing part in the since Pseudomonas and Arabidopsis were model organisms.

5. It is also stated that GuideMaker can design gRNAs for any PAM sequence from any Cas system but the results of SaCas and StCad was described in only one sentence. This is now also shown in detail in Supplementary Figures 1-4. Guidemake allows any PAM to be chosen and more complex PAMs run faster, Supplementary Figure 1-2.

6. The source of the genomes was missing in the manuscript. In particular, some species have multiple genome versions in the same database. Therefore, to make the results more repeatable, the specific website and version number for each species are needed. This is a good point we have added the exact Accessions to the main text of manuscript (lines 176-179) and Supplementary Table 1.

Minor comments

1. Line 11, "bacteria" should be "bacterias".

It appears that "bacteria" is an acceptable plural form of the singular noun "bacterium", based on this explanation: <https://www.merriam-webster.com/dictionary/bacteria>

2. Line 38, delete the", including non-model organisms", prokaryotic and eukaryotic organisms include the non-model organisms.

Deleted.

3. Line 111, "candidates guides" should be "candidate guides".

Corrected.

4. Line 154, "gRNA identify with GuideMaker" should be "gRNA identified with GuideMaker".

Corrected.

5. Line 195, "The second way GuideMaker reduces..." should be "The second way that GuideMaker reduces...".

This section was rewritten so the text no longer exists.

6. Line 204, "and", no need for italics.

This was italicized for emphasis. I have removed the italics.

7. Line 207, "gRNA's" should be "gRNAs".

Corrected

8. Lines 209-210, "we anticipate performance will..." should be "we anticipate that performance will...". Added optional that.

9. Figure. 1. It seems that the font size of the description of Control gRNAs is inconsistent with others, please check.

The entire document has been reformatted to 12-point font.

10. Line 22,55,98,159,175,187,219 and 247, "Guidemaker" should be "GuideMaker". Thanks, the format is now consistent.  
11. Line 262, "CAS" should be "Cas".  
Corrected  
12. Supplementary Figure 4. Grammar mistake in sentence "the different number of logical cores with or without AVX2 settings are available". It should be "the different number of logical cores with or without AVX2 settings is available".  
This has been rewritten for clarity.

### Reviewer 3

Overall, the tool is very well documented and easy to use. In the current version of the manuscript, GuideMaker does not show a clear improvement over the state-of-the-art design tool, CHOPCHOP. The authors do not implement any existing on-target scoring methods to determine the targeting efficacy of the picked sgRNAs. This can lead to picking guides that are highly specific but not effective enough. We have improved our explanation of the advantages of GuideMaker relative to CHOPCHOP for its intended applications, including a performance evaluation in Supp. Fig. 5 and a better explanation in the text. We have also added both on-target and off-target scoring for the "NGG" PAM (the only PAM for which training data is available). Based on the model from Doench et al. 2016.

1. Implementing on-target scoring methods, at least for the Cas enzymes that have on-target efficacy information, can help improve the process of picking sgRNAs. This tool will probably be used more often with standard Cas enzymes and it will be useful to have on-target efficacy scores attached to the guide RNAs.

Good suggestion, we have implemented the Azimuth model for on-target scoring from Doench et al. 2016, specifically their "V3\_nopos" model. We have also refactored the original feature calling to improve speed, updated code to Python 3.9 and transferred their original model in pickle format to a safer, reproducible, cross platform compatible model in the Onnx runtime. We have also added the off target CFD scoring from the same paper.

2. The authors do a thorough analysis of the computational performance of GuideMaker with various genomes and Cas enzymes but including a comparison of the computational performance of GuideMaker vs. CHOPCHOP will strengthen the manuscript.

We have added this comparison, in Supp. Fig 5.

3. The authors define the PAM sequence of SaCas9 to be NGRRT whereas the canonical PAM sequence of SaCas9 is NNGRRT. This should be modified throughout the manuscript and analyses involving SaCas9 should be redone. We have fixed this issue.

A good addition to the tool would be to output a file with all the sequences that were designed targeting the region of interest with the specific PAM sequence. This gives the user a sense of the universe from which the final guides were picked.

The user can get this by filtering the current output file by the locus name.

5. Another useful input parameter would be to specify a target region that the user wants to focus on such as letting the user input genomic coordinates or a gene name or locus tag. For example, CRISPy by Blin et al., 2016 takes a GenBank file as input and allows the user to input features specific to the uploaded genome.

Minor Points We have added the "--filter\_by\_locus" option to filter results for this application.

1. "CyVerse" is misspelled as "CyCVerse" in multiple places in the manuscript. We have fixed this.

2. Reference Figure 2 in Line 92.

Added.

3. Line 154: "Ratios between tools were calculated by dividing the number of gRNA identified.." The sentence was rewritten for clarity.

4. In Supplementary Figure 3 "wit haVX2" should be "with aVX2".

Corrected.

5. GitHub link in Line 336 does not work.

Those links are fixed.

6. Line 225-226: "GuideMaker also creates off-target gRNAs for use as negative controls in high-throughput experiments." "Off-target gRNAs" is misleading in this context.

	We now refer to them as “off-target control RNA sequences” since they are not guides.
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. 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><b>Resources</b></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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?



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

2 Ravin Poudel<sup>1,2</sup>, Lidimarie Trujillo Rodriguez<sup>2</sup>, Christopher R. Reisch<sup>2</sup>, and Adam R. Rivers<sup>1\*</sup>

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

5 <sup>2</sup>Department 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](mailto: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 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 challenging 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 using a degenerate protospacer  
18 adjacent motif (PAM) and a genome. The tool applies Hierarchical Navigable Small World (HNSW) graphs  
19 to speed up the comparison of guide RNAs and optionally provides on-target and off-target scoring. This  
20 allows the user to design effective gRNAs targeting all genes in a typical bacterial genome in about 1-2  
21 minutes.

### 22 **Conclusions:**



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

28

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

30

### 31 **Introduction**

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

40 The *Streptococcus pyogenes* Cas9 (SpCas9) was the first Cas described [5] and it is still the most widely  
41 used enzyme in CRISPR gene editing. Other Cas enzymes described early in the CRISPR revolution, such as  
42 the *Staphylococcus aureus* Cas9 and the *Acidaminococcus* Cas12a, are also commonly used [6,7]. Accordingly, the  
43 parameters for these enzymes are often included in computational tools to identify CRISPR target sites [8–  
44 11]. Cas9 enzymes from other organisms and other Cas-associated proteins that can cleave dsDNA, ssDNA,  
45 ssRNA, and insert transposon elements have also been described and have their place in molecular toolkits  
46 [12–18]. Each of these enzymes generally has specific requirements, such as PAM sequence constraints, PAM

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

55         Here we introduce GuideMaker, a computational tool to identify target sites and design gRNA  
56 sequences that is not limited to any specific CRISPR system or organism. GuideMaker is most useful for a  
57 few kinds of CRISPR experiments. The first use case is designing pools of gRNAs for genome-wide  
58 screening experiments like Perturb-seq and CRISPR pool [23,24]. GuideMaker is optimized for making the  
59 all-versus-all comparisons necessary to design a genome-wide screen and return candidate gRNAs for every  
60 gene locus. The tool allows the user to filter targets based on their proximity to features of interest, like the  
61 start codon for any coding sequence. The second major use case is for researchers working with non-model  
62 organisms. Online gRNA design tools often have a limited number of preselected genomes available for  
63 analysis because most methods require PAM site positions to be precomputed. GuideMaker rapidly computes  
64 all guide positions on demand from user-provided GenBank files or a set of GFF/GTF (general feature  
65 format/general transfer format) files and fasta files from any organism. The third use case is for researchers  
66 working with Cas enzymes other than the canonical versions of Cas9, Cas12a (Cpf1), or Cas13 with different  
67 PAM and target site requirements. GuideMaker allows the user to specify a custom PAM with variable length,  
68 including degenerate nucleotides and allows the PAM to be on either the 3' or 5' side of the protospacer.  
69 These features allow GuideMaker to support any current or future CRISPR-Cas system. Since the  
70 determination of which CRISPR-Cas system functions best in any given organism is not predictable, this tool  
71 is highly relevant to researchers developing CRISPR tools in new species. For SgCas9 GuideMaker also  
72 implements on-target and off-target scoring from Doench et al. (2016). Because there is limited experimental

73 data on most Cas/organism combinations, cannot calculate target scoring for other Cas enzymes but instead  
74 uses design heuristics that prioritize uniqueness in the seed region of the guide.

75

## 76 **Methods**

### 77 **Main features, input parameters, and workflow**

78 GuideMaker is designed to be easy to use as either a web application (Figure 2) or a command-line utility. The  
79 key features of GuideMaker are:

- 80 1. All the potential guides in a genome can be quickly designed in one run.
- 81 2. It can design gRNAs for any PAM sequence from any Cas system.
- 82 3. Search is customizable through user-defined guide parameters (as highlighted in Figure 1). These  
83 features are specific to organisms, CRISPR-Cas systems, and experiments. Tuning these parameters  
84 can improve the sensitivity and specificity of gRNA.
- 85 4. Users can exclude specific restriction sites from guides to preserve those sites for downstream  
86 experiments.
- 87 5. It creates control sequences based on the input genome. In CRISPR experiments it is often desirable  
88 to create negative control sequences to evaluate off-target binding. GuideMaker provides the user  
89 with realistic control gRNAs that are highly divergent from sequences adjacent to PAM sites.
- 90 6. It provides an option to select the subset of results by locus tags of interest.
- 91 7. It provides off-target Cutting Frequency Determination (CFD) scores for gRNAs [8].
- 92 8. Provides on-target efficacy score for canonical “NGG” PAM. These efficiency scores are based on  
93 Azimuth algorithm[8].
- 94 9. Provides tabular result files which can be used for the design and ordering of gRNA pools.
- 95 10. Provides an interactive visualization and exploratory tool to evaluate the guides.
- 96 11. The software can be run as a web application [25], a CyVerse application, or a command-line  
97 application [26]. Server code is included for running local instances of the web application as well.

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

## 112 **Search method**

113           GuideMaker initially scans the genome, recording all candidate guide sequences adjacent to the  
114 specified PAM sequence on both DNA strands (Figure 3). Candidate guides are then optionally checked for  
115 the restriction sites. Next, the candidate guides are searched for a unique "seed region" closest to the PAM  
116 site and candidate gRNAs that are not unique in their "seed region" are removed. Then, approximate nearest  
117 neighbor search is used to remove candidate guides too similar to PAM adjacent sequences in the genome,  
118 based on Hamming distance by default (the number of substitutions required to turn one DNA sequence into  
119 another equal-length sequence). Levenshtein distance is optionally available on the command line and  
120 CyVerse versions of GuideMaker, but it is substantially slower and returns guides that are ~99.9% similar for  
121 tests in *E. coli*. (Supplementary Table 4). The approximate nearest neighbor search is performed using the  
122 Hierarchical Navigable Small World (HNSW) graph method in the Non-Metric Space Library (NMSLIB)  
123 [31,32]. An index of all the initial candidate guides is created using the bitwise Hamming distance metric.

124 Each guide with a unique "seed region" is compared to all candidate guides and any guides with edit distances  
125 (user can select either the Hamming or Levenshtein distance, with Hamming being the default) below the  
126 user-set threshold are removed. This differs from the standard procedure of indexing the genome and  
127 mapping each candidate guide against the whole genome then parsing each result. HNSW has a search  
128 complexity of  $\mathcal{O}(\log N)$  and index complexity of  $\mathcal{O}(N \cdot \log N)$  [31]. Finally, user-defined criteria are applied  
129 to specify the proximity and orientation of guides relative to genomic features like genes. A list of guides is  
130 then returned to the user with relevant information about the guide and its target genomic features.

131 GuideMaker provides an option to predict off-target CFD scores (flag *-ofd\_score*) on the predicted  
132 guides, and on-target scores (flag *-doench\_efficiency\_score*) for the canonical NGG PAM using the the "version3  
133 no position" gradient boosted regression treed model from Doench et al (2016). These models have been  
134 converted into json format and Open Neural Network exchange format respectively for reproducibility, and  
135 featurization and scoring has been refactored to run in Python 3 using the Onnx model runtime [33].

136 The core of GuideMaker's search method is the HNSW method in NMSLIB [32]. The method  
137 builds a multilayer graph index of the input data and has several parameters that can be optimized for index  
138 building and search to trade-off speed and accuracy. Graph construction is the most time-consuming step in  
139 our tests, and thus grid optimization was run to minimize run time while keeping recall above 99% relative to  
140 the ground truth exact nearest-neighbor search. The grid-optimization parameters: [M, efc, ef, and post] used  
141 in the HNSW graph for approximate nearest neighbor search have been optimized for bacterial genomes. A  
142 Jupyter notebook [34] script for re-optimization and visualization of these hyper-parameters is included in the  
143 test directory of the command-line version of the software and optimized parameters can be passed to  
144 GuideMaker with the *--config* flag.

### 145 **Computational performance**

146 Genomes of different sizes, GC content, and chromosome numbers were used to test the speed and  
147 scalability of GuideMaker (Supplementary Table 1). For benchmarking the performance, the same parameters  
148 were used unless a specific parameter was being tested: a PAM motif of 'NGG', 3' pam orientation, target

149 length of 20, lsr (length of seed region) of 11, before and after parameters of 500, knum of 10, controls of 10,  
150 dist of 3 and threads of 16. We profiled the performance of GuideMaker with different threads [1, 2, 4, 8, 16]  
151 in processors with and without the AVX2 processor instruction set. All tests were run on a single compute  
152 node with 2 x 24 core Intel® X®(R) Platinum 8260 CPU @ 2.40 GHz with Cascade Lake microarchitecture.  
153 Three bacterial genomes, a fungal genome, and a plant genome were used in performance benchmarking:  
154 *Escherichia coli* K12 (NC\_000913), *Pseudomonas aeruginosa* PAO1 (NC\_002516), *Burkholderia thailandensis* E264  
155 (NC\_007651), *Arabidopsis thaliana* (NC\_003070), *Aspergillus fumigatus* (NC\_007194), and *Phaseolus vulgaris*  
156 (NC\_023759). For the gene or locus-specific comparisons, only the guides within the locus coordinates (i.e.,  
157 zero feature distance) were considered.

### 158 **Comparison to existing design method**

159 We compared the results of GuideMaker with the results of the online and command-line versions of  
160 CHOPCHOP[35]. GuideMaker and CHOPCHOP parameters were set to approximate the same search. The  
161 length of the target sequence was set to 20 and zero mismatches were allowed in the seed region (11bp) of the  
162 target. The *Escherichia coli* (str. K-12/MG1655) genome was used with the online version of CHOPCHOP  
163 since it has a limited number of genomes. Targets were searched in 40 Kbp increments to account for  
164 CHOPCHOP's size limitations. Target sequences were searched across multiple 40 Kbp segments of *E.coli*  
165 genome (NC\_000913.3:2001-42000, NC\_000913.3:80001-120000, NC\_000913.3:160001-200000,  
166 NC\_000913.3:240001-280000, and NC\_000913.3:320001-360000 ). We also searched for target sequences  
167 and genes/locus\_tags within 40Kbp of (NC\_000913.3:2001-42000) to compare identifications at the locus  
168 level. The ratio between the tools was calculated by dividing the number of gRNA identified with  
169 GuideMaker by the number of guides identified by CHOPCHOP to represent the proportion of guides  
170 identified by both GuideMaker and CHOPCHOP.

171 The command-line version of CHOPCHOP was used to compare the memory usage and  
172 computation time of CHOPCHOP and GuideMaker over an entire genome. The *E. coli* K-12 genome was  
173 chosen for comparison because the precomputed 2bit genome files and Bowtie indexes were provided with

174 CHOPCHOP v 3. The matching GenBank file was downloaded for Guidemaker and both programs were  
175 run 5 times on the same machine using different numbers of processor cores [1, 2, 4, 8, 16].

176

## 177 **Results**

178 The time for GuideMaker to complete a typical run identifying all SpCas9 gRNAs (PAM 'NGG') in a  
179 bacterial genome using 8 compute cores was 75 seconds for *E. coli* and 130 seconds for *P. aeruginosa* (Figure  
180 4). For SaCas9 and StCas9, which have a longer PAM sequence ('NGRRT' and 'NNAGAAW' respectively,  
181 with 3' PAM orientation) and thereby fewer potential targets, the same genomes ran in 19 or 5 seconds  
182 (Supplementary Figures 1). The fungus *Aspergillus fumigatus* (28MB) and the plants *Arabidopsis thaliana* (114  
183 MB) and *Phaseolus vulgaris* (537MB) have larger genomes but are still processed quickly. *A. fumigatus* processed  
184 between 23-304 seconds, while *A. thaliana* processed in 250-921 and *P. vulgaris* processed in 333-4162 seconds  
185 depending on the number of cores, AVX2 instructions, and PAM sequence (Supplementary Figure 2).  
186 GuideMaker can take advantage of Advanced Vector Extensions (AVX2) on newer x86 processors, which  
187 improves the search speed because HNSW search is accelerated with AVX2 (Supplementary Figure 3). The  
188 acceleration was larger when fewer processors were available (Supplementary Figure 3). The HNSW  
189 algorithms are parallelized, and indexing-and-search takes most of the compute time in GuideMaker so the  
190 software scales well when additional cores are added up to 8 cores (Supplementary Figure 3). In practice it  
191 scaled up sub-linearly with genome size, globally estimating Cas9 guides for *E. coli* MG1655 (4.6MB) in 75  
192 seconds and *Phaseolus vulgaris* (537MB) in 1549 seconds, both on 8 cores (Memory usage: 1.9GB for *E. coli* and  
193 46.9GB for *P. vulgaris*, Supplementary Figure 4).

194 The results of GuideMaker were compared with the popular guide design software CHOPCHOP  
195 version 3 [35]. When GuideMaker's filtering settings are set to match CHOPCHOP, the results are very  
196 similar and 99.9% of the targets identified by GuideMaker fall within 2bp of target coordinates returned by  
197 CHOPCHOP. When GuideMaker's unique seed region criterion was not applied at the loci level, the average  
198 number of guides identified by the two approaches was similar per locus (Mean GuideMaker = 116.8, Mean

199 CHOPCHOP = 113.6, p-value = 0.86, Supplementary Table 2). Although the number of guides identified  
200 per gene locus differed, none of the genes were missed by either tool. GuideMaker's default requirement of a  
201 unique seed region is more stringent than CHOPCHOP, and with it enabled, GuideMaker returns  
202 (count=1787) 38.4% (for 2Kbp-42Kbp regions) of the targets compared to CHOPCHOP (count=4651) *E.*  
203 *coli* K12. At the sequence level, 96.7% of the identified gRNA (1729/1787) from both tools had identical  
204 sequences. The more stringent seed region filtering used by default in GuideMaker reduced the CFD scores  
205 of guides suggesting that it would reduce off target binding (Supplementary Figure 6), but that would need to  
206 be experimentally validated in a range of organisms. The ratio of gRNA found by both the tools across the  
207 multiple 40Kbp regions was 39.2% (sd= 1.9%, Supplementary Table 3) when using GuideMaker's more  
208 stringent default settings. This ratio was calculated by dividing the number of gRNA from GuideMaker by the  
209 number from CHOPCHOP for each 40Kb region. GuideMaker processed an entire *E. coli* genome about 60  
210 times faster than the command line version of CHOPCHOP v3. It also used almost 50% less memory across  
211 all the compared processor cores (Supplementary Figure 5).

212

## 213 Discussion

214 Designing gRNAs is a two-step process where GuideMaker first identifies potential guides adjacent to PAM  
215 sequences and then filters the potential guides based on multiple criteria. The most important criterion is that  
216 each guide has a minimum edit distance from any other sequence adjacent to a PAM site in the genome; this  
217 decreases the likelihood of off-target binding. The second way GuideMaker reduces off-target binding is by  
218 requiring that a set number of bases near the PAM site are unique from any other candidate guide. The 8  
219 bases nearest the PAM are the most important for target specificity, and any mismatch is sufficient to prevent  
220 binding [36,37]. The length of the unique region should be set with consideration for the size of the genome  
221 since requiring short unique regions will limit the number of total guides that can be found. For example,  
222 requiring that every gRNA be unique in the first 3 bp would only allow for  $4^3 = 64$  possible guides to be  
223 designed. For normal *-lvr* values of 9-12 this is only limiting for human-sized genomes and can be disabled by



224 setting `--lr` to 0. All guides designed by GuideMaker are perfect matches to a single site in the genome.  
225 Additional specificity is obtained by requiring all similar PAM-adjacent sequences to be unique in the critical  
226 "seed region" and have a total number of mismatches that exceed the user-defined threshold. This double  
227 criterion is expected to increase specificity.

228         The primary goal of the current version of our software is to support the design of gRNAs for non-  
229 standard Cas enzymes or non-model organisms at the genome scale. Guide RNAs do not perform equally,  
230 thus empirical experiments will be needed to fully validate the functionality and efficacy of gRNA predictions.  
231 Given the similarity in targets identified by GuideMaker and CHOPCHOP, we anticipate that performance is  
232 similar to the current state of the art but applicable to more design use cases. When a unique seed region and  
233 edit distance-based filters were applied, GuideMaker created guides more conservatively, generating only  
234 about 40% of the guides created by CHOPCHOP. While CHOPCHOP has an option to specify the  
235 maximum number of mismatches in the first 9 bp or the whole guide, it does not allow the application of  
236 both criteria. While there are small differences in the number and position of guides generated by  
237 GuideMaker, with GuideMaker being more conservative by default, both programs create enough guides to  
238 target nearly all gene loci in the genome of *E. coli*. The current version of the GuideMaker provides options to  
239 predict off-target CFD scores and on-target scores for the canonical NGG PAM. Both scoring approaches  
240 are based on the publicly available models trained on empirical data with SpCas9. If experimentally validated  
241 data become available from genome-wide screens with different Cas enzymes, future versions of GuideMaker  
242 could potentially incorporate new scoring models to help rank candidate guides.

243         GuideMaker is a fast and flexible tool for designing guide RNA across the entire genome in non-  
244 model organisms or with non-canonical Cas enzymes. It takes advantage of fast HNSW search to quickly  
245 index and search new genomes. Several parameters can be tuned to ensure compatibility with the specific  
246 application of the user. For example, GuideMaker checks the designed gRNA for a given restriction enzyme  
247 site to prevent incompatibility with the cloning strategy. Second, the maximum distance from a target  
248 sequence from the start of an annotated feature can be chosen to disrupt promoters or the beginning of the  
249 coding sequence, since these sites are preferred for CRISPRi experiments. GuideMaker also creates off-target

250 control RNA sequences for use as negative controls in high-throughput experiments. Lastly, the program  
251 plots the results for visual exploration of the targets and exports the data as .csv files. The software is  
252 available as a command-line application, a web application, and is integrated into the CyVerse Discovery  
253 Environment to provide users with a range of usage options. Guidemaker is a fast, flexible design tool for the  
254 creation of challenging guide RNA pools.

255

## 256 **Availability and Requirements**

257 Project name: GuideMaker

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

259 Operating system(s): Linux or macOS

260 Programming language: Python >=3.6

261 Other requirements: pybedtools>=0.8.2, nmslib>=2.0.6, streamlit>1.2.6, biopython >=1.79, pandas >= 1.0,

262 onnx>= 1.8.1

263 License: CC0 1.0 Public Domain Dedication

264

265

## 266 **Competing Interests**

267 Authors declare no competing interests

268

## 269 **Data Availability**

270 The source code and command-line executables for GuideMaker are available at the Zenodo [38] and can be  
271 installed directly from Github [26], Bioconda [27], or as a Docker container [28]. Data and code to reproduce  
272 the analysis in the paper are available at Zenodo [38]. As a work of the United States Department of  
273 Agriculture, GuideMaker is released to the public domain under a Creative Commons (CC0) public domain  
274 attribution. The program is also available as a web application through the CyVerse discovery environment  
275 [29], and as a stand-alone web application [25].

276

## 277 **Additional Files**

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

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

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

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

282 **Supplementary Figure 5.** Performance and memory usage comparisons between CHOPCHOP (CLI  
283 version) and GuideMaker.

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

286 **Supplementary Table 1:** Organisms features

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

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

290 **Supplementary Table 4:** Comparison of processing times and results for Levenshtein and Hamming  
291 distances

292

293 **List of abbreviations**

294 BWA: Burrows-Wheeler Aligner; Cas: CRISPR-associated protein; CDS: CoDing Sequences; CRISPR:  
295 Clustered Regularly Interspaced Short Palindromic Repeats; gRNA: Guide RNA; HMSW: Hierarchical  
296 Navigable Small World; NMSLIB: Non-Metric Space Library; PAM: Protospacer Adjacent Motif

297

298 **Funding**

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

303

304 **Author Contributions**

305 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  
306 the software and performed the experiments. R.P., L.T.R., C.R.R., and A.R.R, tested the software, wrote, and  
307 revised the manuscripts. All authors read and approved the final manuscript.

308

309 **References**

- 310 1. Jiang W, Bikard D, Cox D, Zhang F, Marraffini LA. RNA-guided editing of bacterial genomes using  
311 CRISPR-Cas systems. Nat Biotechnol. 2013; doi: 10.1038/nbt.2508.
- 312 2. Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. Genome engineering using the CRISPR-Cas9  
313 system. Nat Protoc. 2013; doi: 10.1038/nprot.2013.143.

- 314 3. Mojica FJM, Díez-Villaseñor C, García-Martínez J, Almendros C. Short motif sequences determine the  
315 targets of the prokaryotic CRISPR defence system. *Microbiology*. 2009; doi: 10.1099/mic.0.023960-0.
- 316 4. Pickar-Oliver A, Gersbach CA. The next generation of CRISPR–Cas technologies and applications. *Nat*  
317 *Rev Mol Cell Biol*. 2019; doi: 10.1038/s41580-019-0131-5.
- 318 5. Deltcheva E, Chylinski K, Sharma CM, Gonzales K, Chao Y, Pírzada ZA, et al.. CRISPR RNA maturation  
319 by trans-encoded small RNA and host factor RNase III. *Nature*. 2011; doi: 10.1038/nature09886.
- 320 6. Ran FA, Cong L, Yan WX, Scott DA, Gootenberg JS, Kriz AJ, et al.. In vivo genome editing using  
321 *Staphylococcus aureus* Cas9. *Nature*. 2015; doi: 10.1038/nature14299.
- 322 7. Zetsche B, Gootenberg JS, Abudayyeh OO, Slaymaker IM, Makarova KS, Essletzbichler P, et al.. Cpf1 Is a  
323 single RNA-guided endonuclease of a class 2 CRISPR-Cas system. *Cell*. 2015; doi: 10.1016/j.cell.2015.09.038.
- 324 8. Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, et al.. Optimized sgRNA design  
325 to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat Biotechnol*. 2016; doi:  
326 10.1038/nbt.3437.
- 327 9. Hiranniramol K, Chen Y, Liu W, Wang X. Generalizable sgRNA design for improved CRISPR/Cas9  
328 editing efficiency. Luigi Martelli P, editor. *Bioinformatics*. 2020; doi: 10.1093/bioinformatics/btaa041.
- 329 10. Xu H, Xiao T, Chen CH, Li W, Meyer CA, Wu Q, et al.. Sequence determinants of improved CRISPR  
330 sgRNA design. *Genome Res*. 2015; doi: 10.1101/gr.191452.115.
- 331 11. Perez AR, Pritykin Y, Vidigal JA, Chhangawala S, Zamparo L, Leslie CS, et al.. GuideScan software for  
332 improved single and paired CRISPR guide RNA design. *Nat Biotechnol*. 2017; doi: 10.1038/nbt.3804.
- 333 12. Anzalone A V., Koblan LW, Liu DR. Genome editing with CRISPR–Cas nucleases, base editors,  
334 transposases and prime editors. *Nat Biotechnol*. 2020; doi: 10.1038/s41587-020-0561-9.
- 335 13. Hale CR, Zhao P, Olson S, Duff MO, Graveley BR, Wells L, et al.. RNA-guided RNA cleavage by a  
336 CRISPR RNA-Cas protein complex. *Cell*. 2009; doi: 10.1016/j.cell.2009.07.040.

- 337 14. Abudayyeh OO, Gootenberg JS, Konermann S, Joung J, Slaymaker IM, Cox DBT, et al.. C2c2 is a single-  
338 component programmable RNA-guided RNA-targeting CRISPR effector. *Science*. 2016; doi:  
339 10.1126/science.aaf5573.
- 340 15. Ma E, Harrington LB, O'Connell MR, Zhou K, Doudna JA. Single-stranded DNA cleavage by divergent  
341 CRISPR-Cas9 enzymes. *Mol Cell*. 2015; doi: 10.1016/j.molcel.2015.10.030.
- 342 16. Klompe SE, Vo PLH, Halpin-Healy TS, Sternberg SH. Transposon-encoded CRISPR–Cas systems direct  
343 RNA-guided DNA integration. *Nature*. 2019; doi: 10.1038/s41586-019-1323-z.
- 344 17. Strecker J, Ladha A, Gardner Z, Schmid-Burgk JL, Makarova KS, Koonin E V, et al.. RNA-guided DNA  
345 insertion with CRISPR-associated transposases. *Science*. 2019; doi: 10.1126/science.aax9181.
- 346 18. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A programmable dual-RNA-  
347 guided DNA endonuclease in adaptive bacterial immunity. *Science*. 2012; doi: 10.1126/science.1225829.
- 348 19. Hsu PD, Lander ES, Zhang F. Development and applications of CRISPR–Cas9 for genome engineering.  
349 *Cell*. 2014; doi: 10.1016/j.cell.2014.05.010.
- 350 20. Qi LS, Larson MH, Gilbert LA, Doudna JA, Weissman JS, Arkin AP, et al.. Repurposing CRISPR as an  
351 RNA-guided platform for sequence-specific control of gene expression. *Cell*. 2013; doi:  
352 10.1016/j.cell.2013.02.022.
- 353 21. Cox DBT, Gootenberg JS, Abudayyeh OO, Franklin B, Kellner MJ, Joung J, et al.. RNA editing with  
354 CRISPR-Cas13. *Science*. 2017; doi: 10.1126/science.aaq0180.
- 355 22. Yan WX, Chong S, Zhang H, Makarova KS, Koonin E V., Cheng DR, et al.. Cas13d Is a compact RNA-  
356 targeting type VI CRISPR effector positively modulated by a WYL-domain-containing accessory protein. *Mol*  
357 *Cell*. 2018; doi: 10.1016/j.molcel.2018.02.028.

358 23. Dixit A, Parnas O, Li B, Chen J, Fulco CP, Jerby-Arnon L, et al.. Perturb-Seq: Dissecting molecular  
359 circuits with scalable single-cell RNA profiling of pooled genetic screens. *Cell*. 2016; doi:  
360 10.1016/j.cell.2016.11.038.

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

363 25. Zhu LJ. Overview of guide RNA design tools for CRISPR-Cas9 genome editing technology. *Front Biol*.  
364 2015; doi: 10.1007/s11515-015-1366-y.

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

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

368 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)  
369 [GBRU/GuideMaker/releases/tag/v0.2.0](https://github.com/USDA-ARS-GBRU/GuideMaker/releases/tag/v0.2.0).

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

371 30. GuideMaker. The GuideMaker docker container. [https://github.com/orgs/USDA-ARS-](https://github.com/orgs/USDA-ARS-GBRU/packages?repo_name=GuideMaker)  
372 [GBRU/packages?repo\\_name=GuideMaker](https://github.com/orgs/USDA-ARS-GBRU/packages?repo_name=GuideMaker).

373 31. Merchant N, Lyons E, Goff S, Vaughn M, Ware D, Micklos D, et al.. The iPlant collaborative:  
374 cyberinfrastructure for enabling data to discovery for the life sciences. *PLOS Biol*. 2016; doi:  
375 10.1371/journal.pbio.1002342.

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

377 33. Malkov YA, Yashunin DA. Efficient and robust approximate nearest neighbor search using hierarchical  
378 navigable small world graphs. *IEEE Trans Pattern Anal Mach Intell*. 2020; doi:  
379 10.1109/TPAMI.2018.2889473.

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

382 35. Labun K, Montague TG, Krause M, Torres Cleuren YN, Tjeldnes H, Valen E. CHOPCHOP v3:  
383 Expanding the CRISPR web toolbox beyond genome editing. *Nucleic Acids Res.* 2019; doi:  
384 [10.1093/nar/gkz365](https://doi.org/10.1093/nar/gkz365).

385 36. Semenova E, Jore MM, Datsenko KA, Semenova A, Westra ER, Wanner B, et al.. Interference by  
386 clustered regularly interspaced short palindromic repeat (CRISPR) RNA is governed by a seed sequence. *Proc*  
387 *Natl Acad Sci U S A.* 2011; doi: [10.1073/pnas.1104144108](https://doi.org/10.1073/pnas.1104144108).

388 37. Hsu PD, Scott DA, Weinstein JA, Ran FA, Konermann S, Agarwala V, et al.. DNA targeting specificity of  
389 RNA-guided Cas9 nucleases. *Nat Biotechnol.* 2013; doi: [10.1038/nbt.2647](https://doi.org/10.1038/nbt.2647).

390 38. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. Source code and command-line executables for  
391 GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes. doi:  
392 [10.5281/zenodo.4849258](https://doi.org/10.5281/zenodo.4849258).

393 39. Poudel R, Rodriguez LT, Reisch CR, Rivers AR. Supporting material and analysis code for GuideMaker:  
394 Software to design CRISPR-Cas guide RNA pools in non-model genomes. doi: [10.5281/zenodo.4898253](https://doi.org/10.5281/zenodo.4898253).

395 40. Hirano S, Abudayyeh OO, Gootenberg JS, Horii T, Ishitani R, Hatada I, et al.. Structural basis for the  
396 promiscuous PAM recognition by *Corynebacterium diphtheriae* Cas9. *Nat Commun.* 2019; doi: [10.1038/s41467-](https://doi.org/10.1038/s41467-019-09741-6)  
397 [019-09741-6](https://doi.org/10.1038/s41467-019-09741-6).

398 41. Gleditsch D, Pausch P, Müller-Esparza H, Özcan A, Guo X, Bange G, et al.. PAM identification by  
399 CRISPR-Cas effector complexes: diversified mechanisms and structures. *RNA Biol.* 2019; doi:  
400 [10.1080/15476286.2018.1504546](https://doi.org/10.1080/15476286.2018.1504546).

401 42. Fu Y, Sander JD, Reyon D, Cascio VM, Joung JK. Improving CRISPR-Cas nuclease specificity using  
402 truncated guide RNAs. *Nat Biotechnol.* 2014; doi: [10.1038/nbt.2808](https://doi.org/10.1038/nbt.2808).



403 43. Wu X, Kriz AJ, Sharp PA. Target specificity of the CRISPR-Cas9 system. Sander JD, Joung JK, editors.  
404 Quant Biol. 2014; doi: 10.1007/s40484-014-0030-x.

405

406

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

408

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

414

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

416

417 **Figure 4. Performance of GuideMaker for SpCas9.** Evaluating the performance of GuideMaker across  
418 three bacterial genomes using the “**NGG**” PAM motif with a target length of 20, unique zone of 11, 3prime  
419 PAM orientation, before and into parameters of 500, knum of 10, controls of 10, and dist of 3. The mean of  
420 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.gbk.gz</i> , <i>Carsonella_ruddii.gbk</i>
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: --cfid_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

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

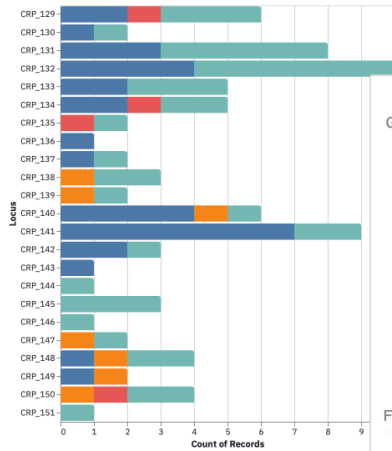
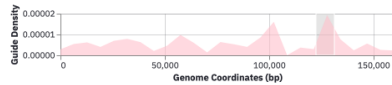
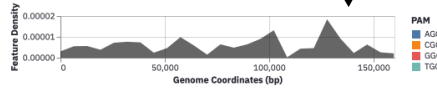
1000 - +

# 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 --lsl 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**

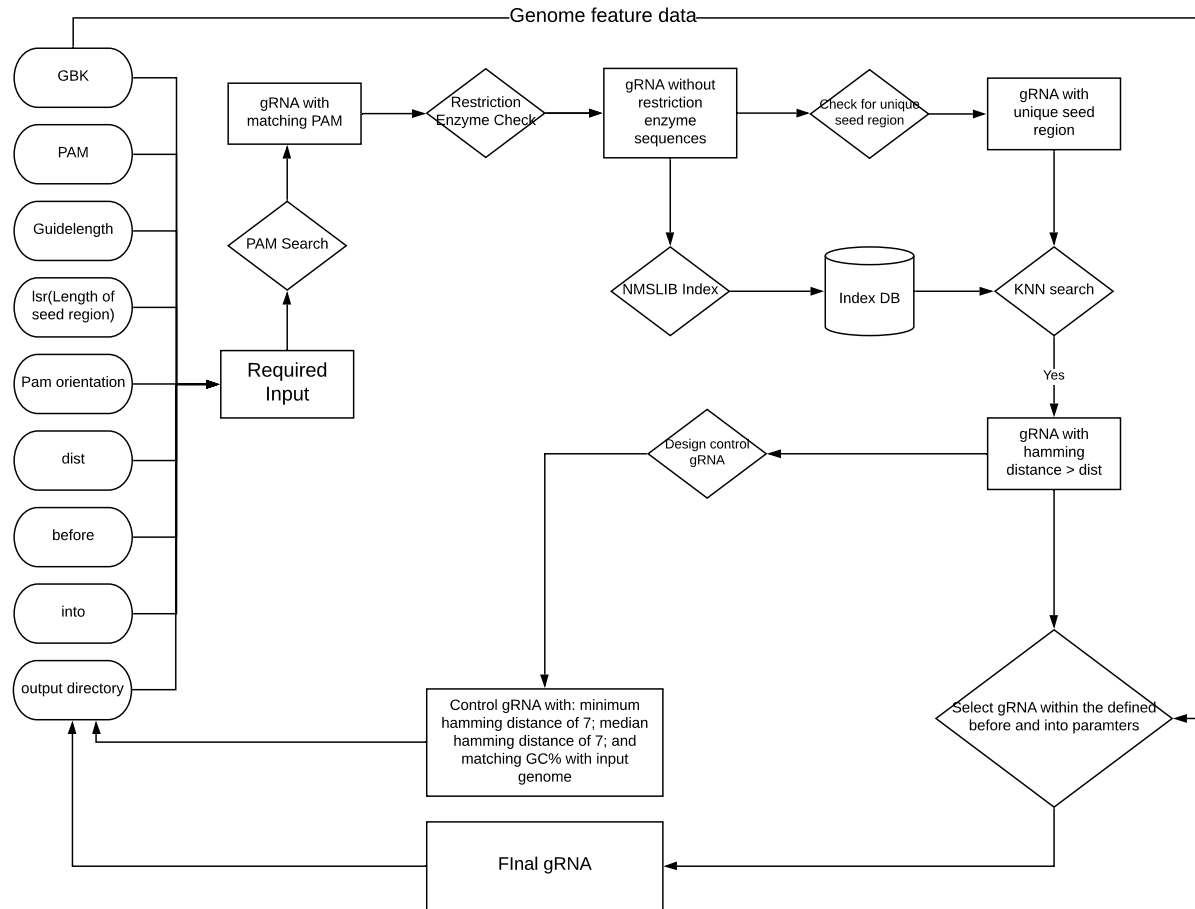
- Parameter Dictionary +
- Designing Experiments with GuideMaker Results +

**API documentation** 📖

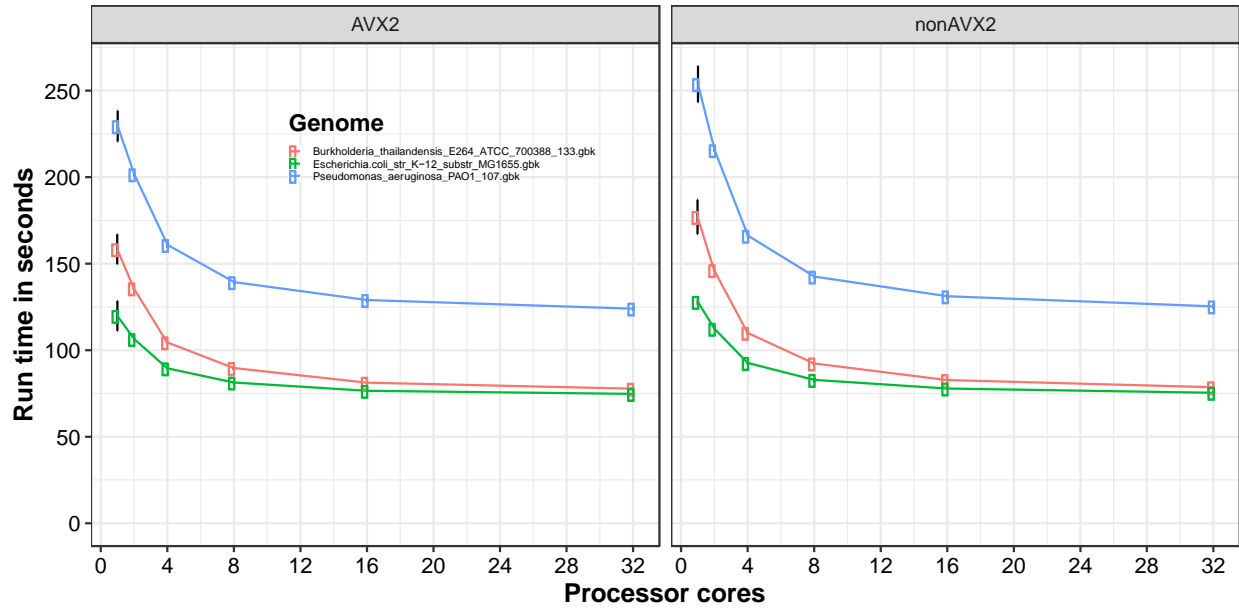
API documentation for the module can be found [here](#)

**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)



Update fig with additional option – may be not –as its very specific to NGG



Remove Non-avx2 to supplement



[Click here to access/download](#)

**Supplementary Material**

[Additional\\_Files\\_GigaScience\\_after\\_Review\\_v2.docx](#)



GIGA-D-21-00186

GuideMaker: Software to design CRISPR-Cas guide RNA pools in non-model genomes  
Ravin Poudel; Lidimarie Trujillo Rodriguez; Christopher R. Reisch; Adam R Rivers  
GigaScience

October 22, 2021

Dear Dr. Edmunds:

Thank you for synthesizing the comments from the reviewers, I found them fair and constructive and have implemented or addressed all suggestions, including feature requests and requests for additional data. For ease of review, I have placed all the reviewer comments and my responses in tabular format. We have also registered GuideMaker with bio.tools (<https://bio.tools/guidemaker>) and SciCrunch.org (SCR\_021778) and included those identifiers in the paper.

Sincerely,

Adam Rivers

Reviewer 1

1. I tested the website and the tool, not finding any bugs and errors. Website is well made, congratulations!	Thanks.
2. Name of the tool: GuideMaker is not self-explanatory for what it is specialized for, which is pooled design. In the future consider naming your tools more distinctly as I am afraid that currently the tool will be buried under hundreds of other GuideSomething tools.	This is a good point. At this point we have a domain, website, preprint and users of the software so it would be pretty disruptive to change, but we will be more specific with future names.
Authors also claim to support Cas13 (page 3 line 65), but don't mention anything more specific about it. I mention that because design for RNA is vastly different from design for DNA and it should be explained how the tool designs for RNA.	We have removed mention of Cas13 since it was not evaluated for this application.
From my understanding the tool offers highly discriminatory settings towards off-target search for a quick resolution of the all vs all comparison problem, however authors ignore that CRISPR off-targets are not defined by the hamming distance, but levenshtein distance. This was proven already by many studies e.g. Tsai et al. 2015. I recommend that authors embrace this issue in the paper and explain why their design may be suitable, and for what kind of studies it	We added an option to use Levenshtein distance to the command-line version of GuideMaker. We also evaluated the effect of using each distance metric on the guides selected there was virtually no difference in the guides selected for the bacterial genome tested. Results are reported in the manuscript on lines 127-130 and in Supplemental Table 4. We suspect this is because there were not many of indels in guides from bacterial genomes. Because hamming is much

would be alright to use hamming distance vs levenshtein distance instead of ignoring the problem.	faster and gave equivalent results, we have kept it as the default distance metric.
5. Study could gain prominence by showing a couple figures and describing how the grid-optimization parameters were selected. This would be especially important for everyone that wants to use this tool for nonbacterial genomes (page 6, lines 128-131). Although script for optimization is included, it would be good to see what are the tradeoffs.	We have added graphical outputs to the grid optimization Jupyter notebook, along with instructions on how to run it for other genomes. We explain how to use this in the paper.
I believe that Figure 4 and all other AVX2 vs nonAVX2 comparisons are not interesting enough to include multiple times. AVX2 improvements are nice, but the tool is already plenty fast, and running time of 250 vs 220 seconds does not matter for normal users.	We have removed the redundant non-AVX figures since most processors now support AVX2. The AVX2 performance gains were larger previously but the NMSlib library improved its non-AVX performance so there is not much difference anymore. Supp fig 3 summarizes the effects of AVX2.
Similarly the number of cores does not seem to influence tool speed above 8 cores and one figure should be enough to explain that.	We removed additional cores above 16, but retained 16 to show the flattening out in performance between 8 and 16.
Tool claims very fast running times, but does not compare to the running times of other similar tools for the design of the pooled screens, this could highlight its superiority.	We now compare the tools to the command-line version of CHOPCHOP using E. coli in Supplemental Fig 5. Despite using precomputed mappings for CHOPCHOP, Guidemaker is about 100x faster and uses 60% less memory
CHOPCHOP is a general tool for the design of pooled screens while here it is used as a pooled screen tool due to its configurability. Additionally, CHOPCHOP also supports all PAM and all species, but on its python version available	That is a good point we have used the CLI version of CHOPCHOP for the added comparisons.
Comparisons to CHOPCHOP focus on the guides found, but I don't understand why consensus ratio between the tools should matter. What is more important is whether GuideMaker does indeed not filter any guides that are preferable for each gene (e.g. by CHOPCHOP ranking) and whether its hamming based filter is good enough to not cause significant unknown off-target effects (levenshtein distance off-targets not found by hamming distance filter). All it takes is one bulge and the hamming distance will become large, while levenshtein distance can even be as low as 1.	We used CHOPCHOP consensus because it is widely used and there is not a gold-standard ground truth data for this. Guidemaker reports about the same number of targets when these same filtering metrics are applied. We have also added Doench et al. 2016 scoring (Azimuth) and CFD scoring to evaluate on target and off target guides for Cas9, so user can sort output by these scores. We added Supp fig 6 to show the effect of selection on On-Target and Off target scores filtered with GuideMaker parameters or unfiltered like CHOPCHOP. Our lsr filtering does not affect on-target scoring but does reduce off target scoring slightly. Our testing has revealed that using Levenshtein distance does not affect guide selection (see explanation above)



It is not clear to me why the tool can't be used with large genomes, filtering on the 11bp seed and hamming distance should be plenty fast for also very large genomes.	It can used for larger genomes just that HNSW loses it speed advantage around at around 1E10 guides. HNSW starts out much faster than indexing and searching the whole genome conventionally but the time per query grows more slowly for the conventional methods. Eventually it becomes faster to use conventional search rather than HNSW search.
Could it be that the tool should support other input, not only genbank file format?	We have added support for importing sequence and annotation from GFF/GTF and Fasta files.

Reviewer 2

The author developed a software, GuideMaker, for designing CRISPR-Cas guide RNA pools in non-model genomes. Three bacterial genomes, a fungal genome, and a plant genome were used in performance benchmarking, which proves that the software supports the design of gRNAs in non-standard Cas enzymes for non-model organisms at the genome-scale. However, the advantages of this software are not well estimated nor presented compared to other tools like CHOPCHOP.	We have improved our explanation of the advantages of GuideMaker relative to ChopChop for its intended applications, including a performance evaluation in Supplemental figure 5 and a better explanation. We have also added both on-target and off-target scoring for NGG PAMs (the only PAMs for which training data is available), from Doench et al. 2016.
Also, the software was mainly evaluated in three bacteria genomes, one fungus and Arabidopsis genome. There are no tests for non-model plant or animal genomes. Therefore, the "non-model genomes" in the title are exaggerated. I list more problems as follows.	We have added the genome of the 537 MG plant <i>Phaseolus vulgaris</i> . Our assertion that Guidemaker can be used for non-model organisms comes from the fact that it does not require precomputed reference genomes but rather computes guide pools quickly on the fly. This feature of the software can be shown without necessarily the genomes of obscure organism. We have clarified this confusing part in the since Pseudomonas and Arabidopsis certainly are model organisms.
The authors did not compare the computation resources and performance (running time, memory) with existing softwares like CHOPCHOP. Also, the authors need to compare the score rankings with CHOPCHOP to present the relative power of GuideMaker. Is there any score rankings concerning efficiency or off-target possibilities for the designed Guide RNAs	This is a good suggestion; we have added Supp. Fig 5 that looks at the time and memory requirements for Guidemaker and CHOPCHOP CLI.  We have added the same on target and off target ranking algorithms used by CHOPCHOP V3. Those algorithms are Azimuth and CDF from Doench et al. (2016).
2. It is better to add support for gff formatted annotation input files since many non-model species do not have GenBank annotations.	We have added support for importing sequence and annotation from GFF/GTF and Fasta files.
3. The authors mentioned GuideMaker can design gRNAs for any small to medium size genome (up to about 500 megabases). The maximum genome used in the article was Arabidopsis thaliana	We have added the 537 MB <i>Phaseolus vulgaris</i> genome in Supp. Fig 4 to demonstrate this claim.

(114.1MB), which is obviously smaller than the described (up to about 500 megabases). We couldn't find the description whether the authors had investigated the larger genomes. Therefore, the detailed analysis or discussion of this problem is needed.	
4. The authors stated GuideMaker to design CRISPR-Cas guide RNA pools in non-model genomes. Arabidopsis thaliana is a model organism and test in a non-model plant genome will be highly valuable.	We have added the genome of the 537 MG plant <i>Phaeoecolis vulgaris</i> . Our assertion that Guidemaker can be used for non-model organisms comes from the fact that it does not require precomputed reference genomes but rather computes guide pools quickly on the fly. We have clarified this confusing part in the since Pseudomonas and Arabidopsis were model organisms.
5. It is also stated that GuideMaker can design gRNAs for any PAM sequence from any Cas system but the results of SaCas and StCad was described in only one sentence.	This is now also shown in detail in Supp. Fig 4. Guidemaker allows any PAM to be chosen and more complex PAMs actually run faster.
6. The source of the genomes was missing in the manuscript. In particular, some species have multiple genome versions in the same database. Therefore, to make the results more repeatable, the specific website and version number for each species are needed.	This is a good point we have added the exact Accessions to the manuscript.
Minor comments	
1. Line 11, "bacteria" should be "bacterias".	It appears that "bacteria" is an acceptable plural form of the singular noun "bacterium", based on this explanation: <a href="https://www.merriam-webster.com/dictionary/bacteria">https://www.merriam-webster.com/dictionary/bacteria</a>
2. Line 38, delete the", including non-model organisms", prokaryotic and eukaryotic organisms include the non-model organisms.	Deleted.
3. Line 111, "candidates guides" should be "candidate guides".	Corrected.
4. Line154, "gRNA identify with GuideMaker" should be "gRNA identified with GuideMaker".	Corrected.
5. Line 195, "The second way GuideMaker reduces..." should be "The second way that GuideMaker reduces..."	This section was rewritten so the text no longer exists.
6. Line 204, "and", no need for italics.	This was italicized for emphasis. I have removed the italics.
7. Line 207, "gRNA's" should be "gRNAs".	Corrected

8. Lines 209-210, "we anticipate performance will..." should be "we anticipate that performance will...".	Added optional that.
9. Figure. 1. It seems that the font size of the description of Control gRNAs is inconsistent with others, please check.	The entire document has been reformatted to 12-point font.
10. Line 22,55,98,159,175,187,219 and 247, "Guidemaker" should be "GuideMaker".	Thanks, the format is now consistent.
11. Line 262, "CAS" should be "Cas".	Corrected
12. Supplementary Figure 4. Grammar mistake in sentence "the different number of logical cores with or without AVX2 settings are available". It should be "the different number of logical cores with or without AVX2 settings is available".	This has been rewritten for clarity.

### Reviewer 3

Overall, the tool is very well documented and easy to use. In the current version of the manuscript, GuideMaker does not show a clear improvement over the state-of-the-art design tool, CHOPCHOP. The authors do not implement any existing on-target scoring methods to determine the targeting efficacy of the picked sgRNAs. This can lead to picking guides that are highly specific but not effective enough.	We have improved our explanation of the advantages of GuideMaker relative to ChopChop for its intended applications, including a performance evaluation in Supp. Fig. 5 and a better explanation. We have also added both on-target and off-target scoring for NGG PAMs (the only PAMs for which training data is available). Based on the model from Doench et al. 2016.
1. Implementing on-target scoring methods, at least for the Cas enzymes that have on-target efficacy information, can help improve the process of picking sgRNAs. This tool will probably be used more often with standard Cas enzymes and it will be useful to have on-target efficacy scores attached to the guide RNAs.	Good suggestion, we have implemented the Azimuth model for on-target scoring from Doench et al. 2016, specifically their "V3_nopos" model. We have also refactored the original feature calling to improve speed, updated code to Python 3.9 and transferred their original model in pickle format to a safer, reproducible, cross platform compatible model in the Onnx runtime. We have also added the off target CFD scoring from the same paper.
2. The authors do a thorough analysis of the computational performance of GuideMaker with various genomes and Cas enzymes but including a comparison of the computational performance of GuideMaker vs. CHOPCHOP will strengthen the manuscript.	We have added this comparison, in Supp. Fig 5.

3. The authors define the PAM sequence of SaCas9 to be NGRRT whereas the canonical PAM sequence of SaCas9 is NNGRRT. This should be modified throughout the manuscript and analyses involving SaCas9 should be redone	We have fixed this issue.
A good addition to the tool would be to output a file with all the sequences that were designed targeting the region of interest with the specific PAM sequence. This gives the user a sense of the universe from which the final guides were picked.	The user can get this by filtering the current output file by the locus name.
5. Another useful input parameter would be to specify a target region that the user wants to focus on such as letting the user input genomic coordinates or a gene name or locus tag. For example, CRISPy by Blin et al., 2016 takes a GenBank file as input and allows the user to input features specific to the uploaded genome. Minor Points	We have added the "--filter_by_locus" option to filter results for this application.
1. "CyVerse" is misspelled as "CyCVerse" in multiple places in the manuscript.	We have fixed this.
2. Reference Figure 2 in Line 92.	Added.
3. Line 154: "Ratios between tools were calculated by dividing the number of gRNA identified.."	The sentence was rewritten for clarity
4. In Supplementary Figure 3 "wit haVX2" should be "with aVX2".	Corrected
5. GitHub link in Line 336 does not work.	Those links are fixed
6. Line 225-226: "GuideMaker also creates off-target gRNAs for use as negative controls in high-throughput experiments." "Off-target gRNAs" is misleading in this context.	We now refer to them as "off-target control RNA sequences" since they are not guides.