

# CROPSR: An Automated Platform for Complex Genome-Wide CRISPR gRNA Design and Validation

**Hans Müller Paul**

University of Illinois Urbana-Champaign

**Dave D Istanto**

University of Illinois Urbana-Champaign

**Jacob Heldenbrand**

University of Illinois Urbana-Champaign

**Matthew Hudson** (✉ [mhudson@illinois.edu](mailto:mhudson@illinois.edu))

University of Illinois Urbana-Champaign

---

## Research Article

**Keywords:** CRISPR, gRNA design, Bioinformatics pipelines, Soybean, Miscanthus, Crops

**Posted Date:** October 8th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-927816/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **METHODOLOGY ARTICLE** 12  
3  
4 **CROPSR: An automated platform for complex** 4  
5  
6 **genome-wide CRISPR gRNA design and** 6  
7  
8 **validation** 8  
910 Hans Müller Paul<sup>1,2</sup>, Dave D Istanto<sup>2,3</sup>, Jacob Heldenbrand<sup>4</sup> and Matthew Hudson<sup>1,2,3,4\*</sup> 10  
1112  
13 **Abstract** 1314 **Background:** CRISPR/Cas9 technology has become an important tool to generate targeted, highly specific 14  
15 genome mutations. The technology has great potential for crop improvement, as crop genomes are tailored to 15  
16 optimize specific traits over generations of breeding. Many crops have highly complex and polyploid genomes, 16  
17 particularly those used for bioenergy or bioproducts. The majority of tools currently available for designing and 17  
18 evaluating gRNAs for CRISPR experiments were developed based on mammalian genomes that do not share 18  
19 the characteristics or design criteria for crop genomes. 1920 **Results:** We have developed the first open source tool for genome-wide design and evaluation of gRNA 20  
21 sequences for CRISPR experiments, CROPSR. The genome-wide approach provides a significant decrease in 21  
22 the time required to design a CRISPR experiment, including validation through PCR, at the expense of an 22  
23 overhead compute time required once per genome, at the first run. To better cater to the needs of crop 23  
24 geneticists, restrictions imposed by other packages on design and evaluation of gRNA sequences were lifted. A 24  
25 new machine learning model was developed to provide scores while avoiding situations in which the currently 25  
26 available tools sometimes failed to provide guides for repetitive, A/T-rich genomic regions. We show that our 26  
27 gRNA scoring model provides a significant increase in prediction accuracy over existing tools, even in non-crop 27  
28 genomes. 2829 **Conclusions:** CROPSR provides the scientific community with new methods and a new workflow for 29  
30 performing CRISPR/Cas9 knockout experiments. CROPSR reduces the challenges of working in crops, and 30  
31 helps speed gRNA sequence design, evaluation and validation. We hope that the new software will accelerate 31  
32 discovery and reduce the number of failed experiments. 3233 **Keywords:** CRISPR; gRNA design; Bioinformatics pipelines; Soybean; Miscanthus; Crops 33  
3435  
36 **Background** 3637  
38 Over the past decade, the CRISPR bacterial system 38  
39 harnessed from *Streptococcus pyogenes* [1, 2, 3] has 3937 \*Correspondence: [mhudson@illinois.edu](mailto:mhudson@illinois.edu)38 <sup>3</sup> Department of Crop Sciences, University of Illinois at Urbana-Champaign,  
39 Urbana, IL 61820, US

39 Full list of author information is available at the end of the article

<sup>1</sup>been optimized and become a revolutionary technique  
<sup>2</sup>for genome editing, leading to the addition of alter-  
<sup>3</sup>native CRISPR systems to the toolbox [1, 4, 5]. One  
<sup>4</sup>of the main challenges of this technology, since its in-  
<sup>5</sup>ception, has been optimizing its cutting efficiency and  
<sup>6</sup>specificity. Many bioinformatics tools have been devel-  
<sup>7</sup>oped to aid in this endeavour [6, 7, 8, 9, 10], with one  
<sup>8</sup>of the most popular tools being developed by Doench  
<sup>9</sup>*et al.* [11, 12]. One of the limitations of the current  
<sup>10</sup>CRISPR design tools is that, for the most part, they  
<sup>11</sup>revolve around the algorithm proposed by Doench *et*  
<sup>12</sup>*al.*. The Doench method is currently considered the  
<sup>13</sup>gold standard for CRISPR guide evaluation. This algo-  
<sup>14</sup>rithm was designed based on mammalian genomes, and  
<sup>15</sup>thus that task is where its performance is best. Many  
<sup>16</sup>crops, however, exhibit paleopolyploidy as a result of  
<sup>17</sup>genome duplication events. These events are evolution-  
<sup>18</sup>arily favorable as they help increase genome diversity  
<sup>19</sup>and the redundant alleles can undergo neofunction-  
<sup>20</sup>alization with low evolutionary pressure, improving  
<sup>21</sup>adaptability to new environment or stress conditions  
<sup>22</sup>[13]. Additionally, crop genomes are particularly af-  
<sup>23</sup>ected by these events, as specific, desirable traits, of-  
<sup>24</sup>ten associated with increased yield or stress resistance,  
<sup>25</sup>have been positively selected over generations of breed-  
<sup>26</sup>ing. RNA guide sequence evaluation models designed  
<sup>27</sup>around mammalian genomes often do not address mul-  
<sup>28</sup>tiple gene copies properly, or at all, as these patterns  
<sup>29</sup>are not frequently observed in these genomes. As a con-  
<sup>30</sup>sequence, the scoring efficiency of such models for use  
<sup>31</sup>on genomes of crop producing plants can be limited.

<sup>32</sup>  
<sup>33</sup>The CRISPR system, as a genome editing tool, is  
<sup>34</sup>comprised of a proteic component in the form of a  
<sup>35</sup>CRISPR-associated protein (Cas) and a single RNA  
<sup>36</sup>molecule. The RNA has a structural domain that in-  
<sup>37</sup>terfaces with the Cas protein, and a domain called  
<sup>38</sup>guide RNA (gRNA) that can be designed to compli-  
<sup>39</sup>ment the DNA region that will be targeted. The Cas

protein identifies the DNA based on a Protospacer Ad-  
jacent Motif (PAM), and once this location is identified  
the DNA double helix is unwound and the gRNA is  
bound to it. The Cas protein then promotes a double-  
stranded break in the double helix of the DNA be-  
fore detaching from it. The cell's native DNA repair  
systems tries to repair the damage, but often end up  
introducing mismatches or causing deletions (fig. 1a).  
One of the main advantages to this system is that to  
target a different portion of the DNA, there's only a  
small part of the system that needs to be redesigned:  
the 20 base long fragment of the gRNA that com-  
plexes with the target DNA. Designing the optimal  
sequence for the guide RNA is important in ensuring  
that the experiment will be more likely to grant the  
expected outcome: a successful mutation. This is in-  
creasingly more important as the complexity of the  
target genome increases. In crops in particular, a com-  
plete CRISPR/Cas9 mutation experiment may take  
upwards of two years (fig. 1b). Inadequately designed  
guides may take weeks, or potentially months, to re-  
veal a failed mutation. This interval is more valuable  
for season-sensitive crops, in which case the window  
for growing the plants may be limited and, for a failed  
experiment to be repeated, there could be a delay until  
the next growing season.

<sup>26</sup>  
<sup>27</sup>  
<sup>28</sup>Here we present CROPSR, an open-source, auto-  
<sup>29</sup>ated platform that efficiently incorporates genome-  
<sup>30</sup>wide gRNA design and evaluation, as well as unique  
<sup>31</sup>primers to facilitate experimental validation of the  
<sup>32</sup>CRISPR/Cas9 knockouts through PCR, to increase  
<sup>33</sup>the accuracy, efficiency and accessibility of CRISPR/  
<sup>34</sup>Cas9 mutation experiment design. CROPSR has been  
<sup>35</sup>specifically developed to be straightforward install and  
<sup>36</sup>use on typical UNIX / HPC hardware, to be easily up-  
<sup>37</sup>datable to include new scoring models or Cas protein  
<sup>38</sup>PAM sites, and to match or exceed the performance  
<sup>39</sup>of current CRISPR/Cas9 gRNA design tools, particu-

larly when used on soybean, *Miscanthus* or other complex, repetitive and polyploid crop genomes. The newly developed model for evaluating gRNA sequences provided a substantial increase in quantity and quality of guide sequences provided compared to currently available tools, particularly for bioenergy crops such as the recently published *Miscanthus sinensis* genome [14].

## Methods

### Overview and required dependencies

CROPSR is a self-contained package written in Python 3, and it will not run on legacy versions (e.g. Python 2.7). Almost all necessary libraries are included with Python 3 as part of the standard library, with the two exceptions being Numpy (<https://numpy.org>) and Pandas (<https://pandas.pydata.org>). Both are widely used, well documented, and can be downloaded from the Python Package Index and installed using the "pip" package installer. Note that for the analysis of larger genomes (2Gb or larger), Python 3.8 or newer is recommended to prevent an issue with a garbage collector bug present on 3.7 and under (this was not an issue with smaller genomes, tested at 750Mb).

### Input files and data pre-processing

To run genome-wide guide design, CROPSR requires the user to provide reference files for both the genomic DNA sequence from a FASTA file, and the genomic annotation from a GFF file. Since these files are user-provided they can be obtained from a variety of sources, including publicly available databases such as NCBI (<https://www.ncbi.nlm.nih.gov>) or Phytozome (<https://phytozome.jgi.doe.gov>), as well as user-generated data for novel unpublished genomes.

### Genome sequence

Once the FASTA file containing the genomic sequence is imported, page breaks are removed and the text is

formatted to account for differences in file formatting<sup>1</sup> for FASTA files from different sources. The file is then<sup>2</sup> converted into a Python dictionary, so sequences can<sup>3</sup> be fetched by chromosome number.<sup>4</sup>

### Annotation files

For correlation of genomic features with the DNA sequence, CROPSR utilizes data from a GFF (preferably GFF3) file. Imported features, including genomic location, strand and phase, source of the GFF file, quality scores, and functional annotation are stored in a Pandas dataframe for ease of access later on. Although this is sufficient information to utilize genomes obtained from the NCBI, data originated from the Phytozome databases utilize a different structure and, thus, require an additional file.<sup>16</sup>

*Handling annotation on Phytozome genomes* Due to Phytozome's internal database organization, the GFF files provided are not associated with functional annotation. Rather, each entry is linked to a transcript name, which is then used to fetch functional annotation from a separate text file on the Phytozome website GUI. Hence, to link functional annotation to genomic locations, the "annotation.info.txt" for a specific genome is also required. This file can be found at the same folder as the GFF file when downloading from the database. When CROPSR identifies the source of the GFF as being Phytozome, it will look for the "annotation.info.txt" in the same folder, and with the same file name structure (e.g. a GFF file named "file\_name.gff3" would cause CROPSR to attempt to open a file named "file\_name.annotation.info.txt" from the same folder). Data from this secondary file is then appended to the function annotation field on the dataframe.<sup>35</sup>

### Genome-wide gRNA design

One of the main factors that set CROPSR apart from other gRNA design being widely utilized is its capa-<sup>38</sup>

1bility to generate guides genome-wide, besides being  
 2able to do so on a per-gene basis. A *gene* flag will  
 3prompt the input to require only a FASTA file, whereas  
 4a *genome* flag will require both FASTA and GFF. Most  
 5of the procedure described below is the same in both  
 6scenarios, but differences will be highlighted as needed.  
 7

### 8 *Identification of PAM sites and guide design*

9 PAM motifs for the specified CRISPR system will be  
 10 identified by regular expression across the sequence,  
 11 both in the sense and anti-sense strands, excluding the  
 12 final bases on both the 5' and 3' ends. For each motif  
 13 identified this way, a 20 base pair guide RNA will be  
 14 designed following the instruction set for that CRISPR  
 15 system (e.g. for Cas9, the guide will comprise the 20 bp  
 16 upstream of the PAM motif). Then, a longer version of  
 17 the guide sequence including 5 bp flanking on both the  
 18 5' and 3' ends is generated for on-site score calculation.  
 19

### 20 *Guide evaluation*

21 *On-site score* CROPSR has two different scoring al-  
 22 gorithms implemented to evaluate guide sequences.  
 23 The first is a modified version of the Doench algo-  
 24 rithm [11]. For each guide, the following procedure is  
 25 applied: first, the sequence is converted into two one-  
 26 hot matrices, one first-order and one second-order, as  
 27 per Doench's description. Then, the first order matrix  
 28 and second order matrix get multiplied by weight ma-  
 29 trices obtained from [11] to generate the first order  
 30 score matrix  $I$  and second order score matrix  $J$ , re-  
 31 spectively. Guides are then assigned a score  $k$  based  
 32 on their GC content ( $k = -0.2026259$  if GC content  
 33  $< 50\%$ , and  $k = -0.1665878$  if GC content is  $\geq 50\%$ ).  
 34 A final on-site score  $f(s)$  is then calculated as a logis-  
 35 tic regression, using 0.59763615 for the intercept [11],  
 36 as shown on equation 1.  
 37

$$38 \quad f(s) = \left[ 1 + \exp^{-(int+k+\sum I_{r,c}+\sum J_{r,c})} \right]^{-1} \quad (1)$$

The second is an algorithm based on a linear support  
 vector regressor (SVR), designed for CROPSR to rem-  
 edy problems identified in the algorithm described by  
 Doench *et al.* [11]. The initial procedure is retained:  
 the sequence is converted into two one-hot matrices,  
 one first-order and one second-order, as previously de-  
 scribed. The two one-hot matrices are then converted  
 into a single vector, which is fed into a linear sup-  
 port vector regressor as the feature set. The model  
 was trained utilizing the same dataset used by Doench  
*et al.* [11], and the feature set was modeled to predict  
 the gene % rank, a continuous variable, and generate  
 a score, rather than defining a threshold for classifica-  
 tion (thus allowing lower-scoring guides to be accessed  
 by users if required).  
 15  
 16

*Note on off-site score* CROPSR does not compute  
 off-site scores when designing guides. Instead, guides  
 that align to multiple locations on the genome are  
 tagged, so they can be identified in the database, and  
 either filtered out or used to edit paralogous regions in  
 parallel.  
 22  
 23

### 24 *Data and associated metadata storage*

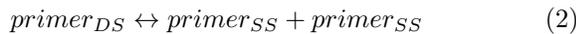
25 Metadata including start and end positions, chromo-  
 26 some number, CRISPR system information, guide se-  
 27 quence, on-site score and a randomly generated unique  
 28 ID for each guide are stored on a Pandas dataframe for  
 29 ease of access.  
 30

### 31 *Functional annotation*

32 The cutsite position, determined based on the specific  
 CRISPR system, is utilized to fetch functional anno-  
 33 tation by cross-referencing with the dataframe con-  
 34 taining data from the GFF file input. Each entry on  
 the GFF dataframe has both a start and end position  
 35 fields, as well as chromosome number. Each cutsite has  
 36 its functional annotation field appended by every GFF  
 37 entry containing the cutsite in its start-end interval.  
 38  
 39

<sup>1</sup>Primer design

<sup>2</sup>For each guide, a pair of PCR primers is also de-  
<sup>3</sup>signed for quick experimental verification. Primer de-  
<sup>4</sup>sign is conducted utilizing an in-house algorithm on  
<sup>5</sup>the *prmrdsn* module, rather than depending on avail-  
<sup>6</sup>able pre-existing tools such as Primer3 [15, 16]. The  
<sup>7</sup>module can accept either FASTA files or a MongoDB  
<sup>8</sup>entry as inputs, with predominantly the same pro-  
<sup>9</sup>cedures being conducted. This grants the user flexi-  
<sup>10</sup>bility to utilize this module as a standalone primer  
<sup>11</sup>design tool, or as part of the CROPSR suite. For  
<sup>12</sup>each input, pools of possible forward and reverse  
<sup>13</sup>primer sequences are generated using regular expres-  
<sup>14</sup>sion. Melting temperatures ( $T_m$ ) for each candidate in  
<sup>15</sup>the pools is determined through thermodynamics pa-  
<sup>16</sup>rameters, assuming hybridization as a two-state pro-  
<sup>17</sup>cess [17, 18]. Under this assumption, we can describe  
<sup>18</sup>the thermodynamic parameters for forming single-  
<sup>19</sup>stranded (*primer<sub>SS</sub>*) primers from double-stranded  
<sup>20</sup>primer (*primer<sub>DS</sub>*) (equation 2).



<sup>24</sup>The equilibrium constant for this reaction is given  
<sup>25</sup>by  $K = \frac{[\text{primer}_{SS}][\text{primer}_{SS}]}{[\text{primer}_{DS}]}$ . Van't Hoff's equation  
<sup>26</sup>defines the relation between free energy,  $\Delta G$ , and  $K$   
<sup>27</sup>as  $\Delta G^\circ = -RT \ln K$ , where  $R$  is the ideal gas con-  
<sup>28</sup>stant,  $T$  is the reaction temperature, in kelvin. Thus,  
<sup>29</sup>we have the derived equation 3. The melting temper-  
<sup>30</sup>ature,  $T_m$ , is the point in which when half of the  
<sup>31</sup>double-stranded has been dissociated, which means  
<sup>32</sup> $[\text{primer}_{DS}]$  is then equal to half of its initial con-  
<sup>33</sup>centration, and  $[\text{primer}_{DS}] = [\text{primer}_{SS}]$ . Plugging  
<sup>34</sup>these values and isolating for  $T_m$ , we have a way to  
<sup>35</sup>predict the melting temperature based on the Gibbs'  
<sup>36</sup>free energy (equation 4) [19, 20].

$$\Delta G^\circ = RT \ln \frac{[\text{primer}_{SS}][\text{primer}_{SS}]}{[\text{primer}_{DS}]} \quad (3)$$

$$T_m = \frac{\Delta G^\circ}{RT \ln \frac{2 \times [\text{primer}_{SS}]}{[\text{primer}_{DS}]}} = \frac{\Delta G^\circ}{RT \ln 2} \quad (4)$$

Interactions between bases on different strands are  
influenced by neighboring bases, which is why the free  
energy at  $37^\circ C$ ,  $\Delta G_{37}^\circ$ , is calculated based on the  
nearest-neighbor model, as proposed by SantaLucia  
[21] and shown in equation 5.  $\Delta G_{37}^\circ(\text{init})$  is the sum  
of the free energies for bases on the extremities, and  
each  $\Delta G_{37}^\circ(i)$  is one of the ten possible adjacent-base  
combinations, as listed on table 1 [21].

$$\Delta G_{37}^\circ(\text{pred}) = \Delta G_{37}^\circ(\text{init}) + \sum_{i=1}^{10} n_i \Delta G_{37}^\circ(i) \quad (5)$$

**Table 1 Nearest-neighbor parameters for DNA/DNA duplexes**

| Nearest-neighbor sequence<br>(5' → 3' / 3' → 5') | $\Delta G_{37}^\circ$<br>(kJ × mol <sup>-1</sup> ) |
|--|--|
| AA/TT  | -4.26  |
| AT/TA  | -3.67  |
| AC/TG  | -6.09  |
| AG/TC  | -5.40  |
| TA/AT  | -2.50  |
| TC/AG  | -5.51  |
| TG/AC  | -6.12  |
| CG/GC  | -9.07  |
| GC/CG  | -9.36  |
| CC/GG  | -7.66  |
| Terminal A/T                                     | 4.31   |
| Terminal C/G                                     | 4.05   |

Primers from the forward and reverse pools then get  
matched based on their melting temperatures to form  
primer pairs that work under the same PCR condi-  
tions.

*Selecting unique primer pairs*

For each primer pair, the amplified fragment length  
and sequence are recorded. The amplified sequence is  
aligned against the genome using bowtie [], to iden-  
tify whether this pair generates a unique amplicon of  
the desired length. If multiple fragments appear with  
the same length and similar sequences, the primers are  
tagged so they can be verified for homeologous copies

<sup>1</sup>of the cloned fragment. Whenever possible, a primer  
<sup>2</sup>pair without the homeologous tag will be selected and  
<sup>3</sup>added to the dataframe for a specific guide RNA.

#### <sup>4</sup> *Using as a standalone tool*

<sup>5</sup>When using the *prmrdsn* module of CROPSR as a  
<sup>7</sup>standalone tool, a FASTA file for the reference genome  
<sup>8</sup>is required. The procedure will be the same, except  
<sup>9</sup>the alignment will be performed against the provided  
<sup>10</sup>reference genome, and the primers will be provided as  
<sup>11</sup>an output .CSV file instead of added to a dataframe.

#### <sup>12</sup> *Output files*

<sup>13</sup>CROPSR outputs to a Mongo database by default, but  
<sup>14</sup>an option to, instead, output a .CSV file is provided.

## <sup>15</sup> **Results**

<sup>16</sup>We have developed CROPSR, a tool to design, eval-  
<sup>17</sup>uate, and validate CRISPR sgRNA in biofuel crop  
<sup>18</sup>genomes. The code for CROPSR was written in  
<sup>19</sup>Python 3.7, and is available at [https://github.com/  
<sup>20</sup>cabbi-bio/cropsr](https://github.com/cabbi-bio/cropsr). To the best of our knowledge,  
<sup>21</sup>CROPSR is the first tool developed for genome-wide  
<sup>22</sup>generation and validation of CRISPR sgRNA in crop  
<sup>23</sup>genomes. One of the main advantages of CROPSR  
<sup>24</sup>is that it was developed as a single, modular, self-  
<sup>25</sup>contained tool that performs all steps in the experi-  
<sup>26</sup>mental design - from identifying the PAM sites located  
<sup>27</sup>in the genome, to the design of PCR primers for ex-  
<sup>28</sup>perimental validation (fig. 2).

#### <sup>29</sup> *Feature comparison*

<sup>30</sup>Although there are tools available that provide a sim-  
<sup>31</sup>ilar set of functions for gene-based analysis and that  
<sup>32</sup>allow the task to be performed in batch [6], they are  
<sup>33</sup>comprised of a pipeline of preexisting methods and,  
<sup>34</sup>thus, have limitations similar to those of the pack-  
<sup>35</sup>ages they include. Additionally, batch-based analysis  
<sup>36</sup>is limited by a gene-based design, which doesn't neces-

<sup>37</sup>sarily consider applications for regulatory, non-coding<sup>1</sup>  
<sup>38</sup>regions.<sup>2</sup>

CROPSR was designed to be a self-contained tool,<sup>3</sup>  
 requiring minimal dependencies to function. Its mod-<sup>4</sup>  
 ular nature brings more control and flexibility to the<sup>5</sup>  
 user. Most parameters can be set via arguments in the<sup>6</sup>  
 command line, which is desirable for supercomputer<sup>7</sup>  
 applications.<sup>8</sup>

#### <sup>9</sup> *Primer design*

<sup>10</sup>Existing tools, such as Primer3 [15, 16], are widely<sup>11</sup>  
 adopted for designing PCR primers for bacterial, fun-<sup>12</sup>  
 gal, and mammalian genomes. For this reason they<sup>13</sup>  
 are often included in pipeline-based genomic tools for<sup>14</sup>  
 CRISPR sgRNA design, such as Chopchop [6, 7, 8].<sup>15</sup>  
 Due to the polyploid nature of many crop genomes,<sup>16</sup>  
 primers need to be designed in a way that allows in-<sup>17</sup>  
 dependent validation of mutations in different gene<sup>18</sup>  
 copies. Implementations of primer design on currently<sup>19</sup>  
 available tools either do not validate whether the am-<sup>20</sup>  
 plicon has matches elsewhere in the genome, or per-<sup>21</sup>  
 form a simple alignment and discard primer pairs<sup>22</sup>  
 that have multiple hits. The first scenario may cause<sup>23</sup>  
 primers that target more than one location in a poly-<sup>24</sup>  
 ploid genome, whereas the second may cause the tool<sup>25</sup>  
 to be unable to provide primers for a specific region.<sup>26</sup>  
 We have included a PCR primer design module within<sup>27</sup>  
 CROPSR's algorithm in an attempt to help mitigate<sup>28</sup>  
 these issues by providing information on additional<sup>29</sup>  
 gene copies that may be affected by the same guide,<sup>30</sup>  
 while also potentially identifying and validating guides<sup>31</sup>  
 that may affect only a single copy of a duplicated gene,<sup>32</sup>  
 when possible.<sup>33</sup>

Additionally, CROPSR's primer design module cal-<sup>34</sup>  
 culates the melting temperatures ( $T_m$ ) based on the<sup>35</sup>  
 nearest-neighbor model proposed by SantaLucia [21].<sup>36</sup>  
 This method provides more accurate temperatures<sup>37</sup>  
 for longer sequences when compared to Bolton and<sup>38</sup>  
 McCarthy [22], used in Primer3. Additionally, the<sup>39</sup>

<sup>1</sup>nearest-neighbor method adopted provides better con-  
<sup>2</sup>sistency in regions with varying GC contents, which is  
<sup>3</sup>often an issue in repetitive, A/T-rich crop genomes.

<sup>4</sup>

#### <sup>5</sup>*Revised model for sgRNA scoring*

<sup>6</sup>There are two well-established tools that are most  
<sup>7</sup>commonly used for designing sgRNA for CRISPR  
<sup>8</sup>knockout experiments in plants, each with its advan-  
<sup>9</sup>tages. Chopchop [6, 7, 8] provides a pipeline that fa-  
<sup>10</sup>cilitates guide sequence design and calculates a score  
<sup>11</sup>that is predictive of the guide’s on-target efficiency,  
<sup>12</sup>as well as its off-target potential. This score is based  
<sup>13</sup>on the methods proposed by Doench *et al.* [11, 12],  
<sup>14</sup>which utilized a support vector machine classifier to  
<sup>15</sup>segregate the designed guides in two groups: the top  
<sup>16</sup>20% performing, classified as optimal, from the bottom  
<sup>17</sup>80%. Additionally, the method penalizes guides that  
<sup>18</sup>can target the genome in more than a single location,  
<sup>19</sup>as this characteristic can promote undesirable, non-  
<sup>20</sup>specific edits. One of the shortcomings of this method,  
<sup>21</sup>in crops specifically, is that many crops exhibit pale-  
<sup>22</sup>opolyploid genomes. In such cases, most essential genes  
<sup>23</sup>will have multiple copies throughout the genome, ei-  
<sup>24</sup>ther due to the presence of multiple alleles, or dupli-  
<sup>25</sup>cation events in the evolutionary history of that crop.  
<sup>26</sup>Thus, in crop genetics applications, the elimination of  
<sup>27</sup>guides that have %AT outside the targeted range, or  
<sup>28</sup>hit in more than one location, may result in the soft-  
<sup>29</sup>ware being unable to output any guides. Additionally,  
<sup>30</sup>for some experiments it may be desirable to target all  
<sup>31</sup>of the paralogs in a particular group [23].

<sup>32</sup>CRISPR-P [9, 10] is a web-based tool developed for  
<sup>33</sup>designing sgRNA for CRISPR systems in plants. Its al-  
<sup>34</sup>gorithm [24] calculates a score based on the sgRNA’s  
<sup>35</sup>chance to hit at off-target positions. A score is then  
<sup>36</sup>calculated for a guide based on the number of matches  
<sup>37</sup>it displays throughout the genome, including a limited  
<sup>38</sup>number of mismatches. This approach helps mitigate  
<sup>39</sup>the effect of discarding guides that may hit at many lo-

cations, improving the user capability to design guides<sup>1</sup>  
for crop genes that exhibit multiple copies. One of the<sup>2</sup>  
shortcomings of this method, however, is that it does<sup>3</sup>  
not calculate a score for the on-site binding, adopted<sup>4</sup>  
by other tools like Chopchop.<sup>5</sup>

<sup>5</sup>

<sup>6</sup>

Although the aforementioned models have important,<sup>7</sup>  
qualities and can correctly design functional guides for<sup>8</sup>  
many applications, both present limitations for a num-<sup>9</sup>  
ber of potential CRISPR/Cas9 editing applications in<sup>10</sup>  
crop systems, for example those employed by Dong,<sup>11</sup>  
*et al.* [23]. Hence, to circumvent these limitations, we<sup>12</sup>  
developed a model based on a slight variation of the<sup>13</sup>  
algorithm proposed by Doench *et al.* [11] and adopted<sup>14</sup>  
by Chopchop. For our model, we utilized the training<sup>15</sup>  
data set provided by Doench *et al.*, that was utilized to<sup>16</sup>  
generate the original model they describe, and followed<sup>17</sup>  
their methodology. Through this process we identified<sup>18</sup>  
parameters that may be interfering with the guide de-<sup>19</sup>  
sign for A/T-rich crop genomes, such as melting tem-<sup>20</sup>  
peratures at distinct portions of the 20 base sequence,<sup>21</sup>  
and a GC-content threshold, and removed these pa-<sup>22</sup>  
rameters from the input of our modified model. We<sup>23</sup>  
then compared a number of alternative scoring algo-<sup>24</sup>  
rithms to the one employed by Doench *et al.*. The data,<sup>25</sup>  
was fed into four different types of supervised learning<sup>26</sup>  
models, a support vector machine (SVM) classifier as<sup>27</sup>  
used by Doench *et al.*, a linear support vector regressor<sup>28</sup>  
(SVR), a random forest regressor (RFR), and a multi-<sup>29</sup>  
layer perceptron regressor. Using leave-one-group-out<sup>30</sup>  
as the cross-validation method, we evaluated and com-<sup>31</sup>  
pared each scoring algorithm. Based on these results,<sup>32</sup>  
we opted for the SVR scoring algorithm to incorporate<sup>33</sup>  
into CROPSR, as it performed best with the available<sup>34</sup>  
data while retaining many of the desirable character-<sup>35</sup>  
istics of the original method.<sup>36</sup>

<sup>36</sup>

Additionally, we do not use the off-site scoring model<sup>37</sup>  
used in other approaches (Chopchop uses Doench *et al.*<sup>38</sup>  
*et al.*, 2016 [12], CRISPR-P uses Hsu *et al.*, 2013 [24])<sup>39</sup>

where the off-site hit generates a large penalty in the score of the guide. As this penalizes the score directly and prevents the design of guides that target multiple genes, it prevents a key use of CRISPR/Cas9 in crop systems. We have addressed the presence of potential off-site hits in our output by adding a tag and instructing the user to check for potential additional copies of the targeted region. This is based on an assumption that designing a guide that is capable of editing all copies of a given gene, instead of a single copy, is a likely goal when generating gene knockouts in crop genomes. Separating the off-site hits from the scoring algorithm allows us to provide the user with the power to decide whether a guide that hits the genome at more than one locus is desirable or undesirable.

### Benchmarking

We compared the guide sequences designed by CROPSR, Chopchop, and CRISPR-P for four syntaxin genes from soybean, previously characterized by our group using manually designed gRNAs after existing software failed to design guides [23]. These genes were problematic because they require guides in A/T-rich regions of a crop genome, and have highly similar sequences to one another. In addition, for this project it was necessary to target both each gene individually, and all of the genes simultaneously. CROPSR outperformed the other software tools for every gene in either total number of guides designed, number of guides designed with a score  $\geq 0.8$ , or in both (table 2). The choice of score cutoff was based on the threshold defined in the algorithm by [11], utilized by Chopchop [6, 7, 8], to discard under-performing guides. The only gene for which CROPSR did not suggest the largest total number of guides was Syn16, where CRISPR-P suggested one guide that CROPSR did not. However, CROPSR outperformed CRISPR-P in genes with scores above 0.8 by the same amount. For the other three genes, CROPSR and CRISPR-P suggested the

total number of guides, however CROPSR's scoring algorithm awarded a larger number of guides with scores above 0.8. Chopchop designed fewer guides for all four genes compared to the other two tools, and was outperformed by CROPSR in number of guides with scores above 0.8 for all genes.

**Table 2** Number of guide sequences generated by CROPSR and other currently available software for four syntaxin genes from soybean

| Soybean syntaxin gene | Number of guides with score $\geq 0.8$ (total number of guides) |          |          |
|-----------------------|---|----------|----------|
|                       | CROPSR  | Chopchop | CRISPR-P |
| Syn02                 | 2(141)  | 0(137)   | 0(141)   |
| Syn12                 | 9(217)  | 3(196)   | 2(207)   |
| Syn13                 | 13(396)   | 2(386)   | 2(396)   |
| Syn16                 | 6(386)  | 2(378)   | 5(387)   |

### Improvements to the sgRNA scoring model

The adoption of a regression-based model for the prediction of the gene % rank for sgRNAs provided an improvement of approximately 19% to the variance explained by the model, as represented by the  $r^2$  increase from 0.223 for the model adopted by Chopchop (fig. 3A) to 0.278 obtained for CROPSR (fig. 3B). Likewise, the Pearson's correlation between the independent variable, gene % rank, and the dependable variables, the scores, was improved by 11.5% for CROPSR compared to Chopchop. Another benefit of the model implemented on CROPSR is that the adoption of a quantitative approach improved the correlation with rank across the entire range, but the effect is noticeably larger in the high-scoring ( $\geq 0.8$ ) gRNAs compared to the model implemented by Chopchop (fig. 3C for Chopchop, fig. 3D for CROPSR). Additionally, a comparison of the variances within the target bin (80% – 100%, the highest-performing quintile) utilizing the Kruskal-Wallis method confirmed that the scores generated by Chopchop and CROPSR for that particular group are not part of the same distribution (p-value = 0.02788).

<sup>1</sup> *Whole-genome evaluation*

<sup>2</sup> One of the goals when developing CROPSR was to pro-  
<sup>3</sup> vide users with the ability to design guide sequences  
<sup>4</sup> targeting regulatory regions of crop genomes. Regula-  
<sup>5</sup> tory or non-coding regions often exhibit repetitive se-  
<sup>6</sup> quences and their genetic structure can widely vary  
<sup>7</sup> compared to coding, gene regions. These traits can  
<sup>8</sup> cause currently available software tools to have dif-  
<sup>9</sup> ficulty designing guide sequences in non-coding re-  
<sup>10</sup> gions. Although differences in parameters such as melt-  
<sup>11</sup> ing temperatures and G/C content have an effect  
<sup>12</sup> when evaluating on-target activity of a guide RNA se-  
<sup>13</sup> quence, the presence of repetitive sequences in non-  
<sup>14</sup> coding regions has a bigger impact on the prediction  
<sup>15</sup> of off-target activity. Due to the repetitive nature of  
<sup>16</sup> the DNA in regulatory regions, it is more likely that  
<sup>17</sup> guide RNA sequences originated in such regions will be  
<sup>18</sup> matched to multiple different locations in the genome.  
<sup>19</sup> This can cause the scoring algorithms used in these  
<sup>20</sup> software tools to heavily penalize guide sequences de-  
<sup>21</sup> signed to target non-coding regions, which may lead  
<sup>22</sup> to no guide sequences being considered viable by such  
<sup>23</sup> programs.  
<sup>24</sup>

<sup>25</sup> The approach adopted in CROPSR to circumvent  
<sup>26</sup> these issues is to take the whole genome as an input  
<sup>27</sup> rather than the sequence of a single gene or region. By  
<sup>28</sup> splitting the genome by chromosome, it is possible to  
<sup>29</sup> identify all potential target locations by scanning for  
<sup>30</sup> the PAM sites, and then designing the guide sequences  
<sup>31</sup> for each site following any requirements of the specific  
<sup>32</sup> Cas system (e.g. for Cas9, the sequence is designed  
<sup>33</sup> as the 20 bases located upstream of the PAM site).  
<sup>34</sup> Start and end positions for each guide, as well as the  
<sup>35</sup> cut-site for the Cas nuclease activity and a calculated  
<sup>36</sup> on-site score are then added as metadata to each guide  
<sup>37</sup> sequence.

<sup>38</sup> Each guide sequence is aligned against the entire  
<sup>39</sup> genome to identify potential sequences that target

more than a single location, however no penalty is ap-  
 plied to the score in positive cases. A tag is added to se-  
 quences that hit in multiple locations, and the user will  
 then decide whether these guides are of their interest  
 (e.g. for mutating all copies of a gene) or not. This re-  
 sults in a much larger number of guide sequences being  
 provided compared to the currently available alterna-  
 tives, and is especially useful for the complex genomes  
 of energy crops. A whole-genome analysis of a complex  
 genome, such as *Miscanthus sinensis* [14], can gener-  
 ate upwards of 200 million potential CRISPR targets  
 (table 3).

**Table 3 Number of gRNA sequences generated by CROPSR for *Miscanthus sinensis***

| Chromosome in<br><i>Miscanthus sinensis</i> | # of guide sequences<br>generated by CROPSR |
|---|---|
| Chr01                                       | 15,174,305                                  |
| Chr02                                       | 14,634,737                                  |
| Chr03                                       | 10,933,532                                  |
| Chr04                                       | 11,515,578                                  |
| Chr05                                       | 11,539,743                                  |
| Chr06                                       | 11,947,192                                  |
| Chr07                                       | 16,147,361                                  |
| Chr08                                       | 9,662,343                                   |
| Chr09                                       | 8,046,957                                   |
| Chr10                                       | 7,389,629                                   |
| Chr11                                       | 8,198,766                                   |
| Chr12                                       | 8,558,113                                   |
| Chr13                                       | 6,745,249                                   |
| Chr14                                       | 5,959,044                                   |
| Chr15                                       | 7,181,726                                   |
| Chr16                                       | 8,132,384                                   |
| Chr17                                       | 8,480,072                                   |
| Chr18                                       | 8,443,163                                   |
| Chr19                                       | 8,851,088                                   |
| Unplaced scaffolds                          | 19,158,456                                  |
| TOTAL                                       | 206,699,438                                 |

#### Run time, memory usage and disk space

The full analysis of the genome of *Miscanthus sinensis*  
 using CROPSR required a minimum of 16Gb of RAM,  
 and a quad-core processor. Memory usage peaked at  
 around 12 Gb. The total run time under these specs

<sup>1</sup>was 6 days, 7 hours, 17 minutes and 40 seconds. The  
<sup>2</sup>*.csv* file output is 26.94 Gb in size.

<sup>3</sup>

#### <sup>4</sup>Discussion

<sup>5</sup>As previously described, CROPSR is a tool developed  
<sup>6</sup>from the ground up as an open source Python appli-  
<sup>7</sup>cation to perform all steps required to design guide  
<sup>8</sup>and primer sequences for genome editing, with addi-  
<sup>9</sup>tional consideration paid to the complications of per-  
<sup>10</sup>forming CRISPR/Cas9 editing in complex, often poly-  
<sup>11</sup>ploid crop genomes, such as the need to target multiple  
<sup>12</sup>paralogs and the need for unique validation primers.

<sup>13</sup>

#### <sup>14</sup>CROPSR use cases

<sup>15</sup>The development of CROPSR was inspired by the lim-  
<sup>16</sup>itations of current gRNA sequence design tools for  
<sup>17</sup>A/T-rich regions of crop genomes. However, we quickly  
<sup>18</sup>realized that the polyploid genomes often found in  
<sup>19</sup>crops also impose limitations to algorithms that score  
<sup>20</sup>based on guide sequence uniqueness. This often causes  
<sup>21</sup>otherwise useful gRNA sequences that target multi-  
<sup>22</sup>ple homeologous copies of a gene to be filtered and  
<sup>23</sup>not accessible to the user. Finally, unique validation  
<sup>24</sup>primers for each target site are often the limiting fac-  
<sup>25</sup>tor in polyploid genomes or for high-copy genes, thus  
<sup>26</sup>the design of these is integrated into CROPSR. Al-  
<sup>27</sup>though CROPSR was developed with polyploid, A/T-  
<sup>28</sup>rich crop genomes in mind, it is the first tool for  
<sup>29</sup>genome-wide gRNA sequence design, and can be em-  
<sup>30</sup>ployed in any genome. Additionally, due to its modular  
<sup>31</sup>nature, individual modules can be utilized individually  
<sup>32</sup>to design PCR primers for a given sequence, or to origi-  
<sup>33</sup>inate gRNA sequences for a single gene rather than a  
<sup>34</sup>whole genome.

<sup>35</sup> We have purposefully designed CROPSR as a set of  
<sup>36</sup>modular tools to facilitate the implementation of new  
<sup>37</sup>functionality, such as new scoring algorithms as they  
<sup>38</sup>become available, or new CRISPR systems as needed.  
<sup>39</sup>An additional advantage of this approach is that in-

dividual modules can be utilized separately from the<sup>1</sup>  
complete package. For example, the PCR primer de-<sup>2</sup>  
sign tool can be used as a stand-alone application for<sup>3</sup>  
any PCR experiment, or the CRISPR guide design tool<sup>4</sup>  
can be used for an isolated gene sequence (for example,<sup>5</sup>  
for an *in vitro*, controlled experiment where matches<sup>6</sup>  
elsewhere in the genome are not a consideration).<sup>7</sup>

<sup>8</sup>

<sup>9</sup>

#### Improved scoring model<sup>10</sup>

When developing the new scoring model for CROPSR,<sup>11</sup>  
we gave special consideration to avoiding imposing<sup>12</sup>  
penalties on guide sequences that match at more than<sup>13</sup>  
a single location on the target genome. Crop genomes<sup>14</sup>  
can possess multiple copies of important genes, and pe-<sup>15</sup>  
nalizing guides that match these genes at more than a<sup>16</sup>  
single locus often causes guide sequences designed to<sup>17</sup>  
target these genes to have poor scores. Such penalties<sup>18</sup>  
against repeated sequences are implemented in Chop-<sup>19</sup>  
chop and many other tools [11, 12, 24]. These scoring<sup>20</sup>  
algorithms have been developed to facilitate CRISPR<sup>21</sup>  
experiments in humans. In plant biology, it is some-<sup>22</sup>  
times desirable to target multiple paralogous sequences<sup>23</sup>  
simultaneously, and in the case of many genes in some<sup>24</sup>  
genomes, this is the only option. Therefore, removing<sup>25</sup>  
guides that target multiple sequences by default before<sup>26</sup>  
reporting results becomes an obstacle when attempt-<sup>27</sup>  
ing to design guide RNA sequences for crop genomes.<sup>28</sup>  
As existing optimizations frequently discard guide se-<sup>29</sup>  
quences with a low score at an early stage, and scores<sup>30</sup>  
are often heavily penalized for targeting multiple loca-<sup>31</sup>  
tions, such software tools can be unable to provide the<sup>32</sup>  
user with gRNA sequences for certain plant genes. The<sup>33</sup>  
approach adopted in the model utilized by CROPSR<sup>34</sup>  
does not apply penalties for sequences with multiple<sup>35</sup>  
hits even when scoring is performed utilizing the im-<sup>36</sup>  
plemented version of the Doench [11] algorithm. How-<sup>37</sup>  
ever, sequences with multiple hits can easily be filtered<sup>38</sup>  
from the reported results.<sup>39</sup>

<sup>1</sup> To further improve the reliability of the provided  
<sup>2</sup>gRNA sequences, a new algorithm based on a linear  
<sup>3</sup>support vector regression method was employed. This  
<sup>4</sup>choice was based primarily on the nature of the data  
<sup>5</sup>being analyzed. The methodology adopted Doench *et al.*  
<sup>6</sup>[11], which inspired our approach, initially converts  
<sup>7</sup>nucleotide sequences into one-hot matrices. A weight  
<sup>8</sup>matrix can then be obtained from a population of one-  
<sup>9</sup>hot matrices, based on the correlation between base  
<sup>10</sup>frequency at specific positions (a continuous variable)  
<sup>11</sup>and an effect state. This is where both methods di-  
<sup>12</sup>verge, however. Doench *et al.* set a threshold at an or-  
<sup>13</sup>dered rank data setting to create two discrete classes,  
<sup>14</sup>which then were used to train a classifier to segregate  
<sup>15</sup>the highest efficient sequences from the remaining pop-  
<sup>16</sup>ulation. The threshold was defined to maximize the  
<sup>17</sup>odds of a successful CRISPR experiment (in human  
<sup>18</sup>and mouse genomes), and simplifies the choice for the  
<sup>19</sup>user. We have instead opted for a different approach,  
<sup>20</sup>attempting to predict where in the ordered rank a  
<sup>21</sup>new data point would be situated. We opted for as-  
<sup>22</sup>sisting the user to make an informed decision about  
<sup>23</sup>whether a guide is suited for an experiment or not,  
<sup>24</sup>based on a combination of the ranking score (provided  
<sup>25</sup>by CROPSR) and field-specific information.

<sup>26</sup> Another important factor in repetitive genomes is  
<sup>27</sup>the ability to design primers to validate edits. Thus,  
<sup>28</sup>although gRNAs that target multiple genes can be de-  
<sup>29</sup>sirable, it is necessary to be able to design unique PCR  
<sup>30</sup>primers to amplify each target site and sequence the  
<sup>31</sup>edited region. For this reason, primer design is inte-  
<sup>32</sup>grated into the CROPSR software suite, and can be  
<sup>33</sup>used to filter CROPSR-designed gRNAs for those al-  
<sup>34</sup>lowing the design of unique primer sets.

### <sup>35</sup>New workflow for CRISPR experiments

<sup>37</sup>With CROPSR, the workflow for CRISPR experi-  
<sup>38</sup>ments in crop genomes can be revised and simpli-  
<sup>39</sup>fied (fig. 4). During the first utilization with a spe-

cific genome, a full analysis and genome-wide guide<sup>1</sup>  
sequence design will be performed. All gRNA se-<sup>2</sup>  
quences and their important metadata will be stored<sup>3</sup>  
in a database for ease of identification. This step, al-<sup>4</sup>  
though time-consuming, only needs to be performed<sup>5</sup>  
once per genome. On subsequent uses, a search to the<sup>6</sup>  
database will return all data required to perform the<sup>7</sup>  
experiment: the gRNA sequence, cut site location and<sup>8</sup>  
genome annotations matching that position, a pair of<sup>9</sup>  
unique PCR primers for experimental validation and<sup>10</sup>  
the melting temperatures (*Tms*) for the primer pair.<sup>11</sup>  
Grouping all the steps on a single platform provides<sup>12</sup>  
ease of use and helps minimize variations, either orig-<sup>13</sup>  
inated by running various different tools, or even be-<sup>14</sup>  
tween different people utilizing the same pipeline.<sup>15</sup>

### <sup>16</sup>CROPSR limitations, in comparison to other tools

<sup>17</sup>To utilize CROPSR's features fully, a whole-genome<sup>18</sup>  
run is recommended prior to performing the CRISPR<sup>19</sup>  
experiments. While this run only needs to be done<sup>20</sup>  
once for each organism, and the results can then be<sup>21</sup>  
stored in a database, the initial run requires signifi-<sup>22</sup>  
cant server-class compute resources. The time required<sup>23</sup>  
to perform this process should be kept in mind when<sup>24</sup>  
planning experiments utilizing this tool. This step is<sup>25</sup>  
required on the first use for a given genome, to generate<sup>26</sup>  
a database containing all potential gRNA sequences<sup>27</sup>  
and their associated metadata. On subsequent utiliza-<sup>28</sup>  
tion, database searches should provide the user with<sup>29</sup>  
the necessary guide and primer sequences with mini-<sup>30</sup>  
mum further calculation, and can be delivered through<sup>31</sup>  
a web interface.<sup>32</sup>

### <sup>33</sup>Future directions

<sup>34</sup>The release of CROPSR is the first step on a longer<sup>35</sup>  
path to provide an innovative, open source toolkit<sup>36</sup>  
to assist in the design of CRISPR experiments and<sup>37</sup>  
other manipulations of crop genomes. During devel-<sup>38</sup>  
opment and testing of CROPSR, the challenges that<sup>39</sup>

<sup>1</sup>emerged have helped develop improved strategies for  
<sup>2</sup>future work.

#### <sup>3</sup> <sup>4</sup>*Additional output formats*

<sup>5</sup>One of the consequences of the genome-wide approach  
<sup>6</sup>adopted in CROPSR is the dimension of the output  
<sup>7</sup>file. Processing large, complex genomes result in a  
<sup>8</sup>colossal number of guide sequences with associated  
<sup>9</sup>metadata, such as functional annotation and primers  
<sup>10</sup>for PCR validation. In an effort to make the output file  
<sup>11</sup>accessible to the majority of users, the options cur-  
<sup>12</sup>rently provided are CSV and and MongoDB. Future  
<sup>13</sup>releases may include support for other database for-  
<sup>14</sup>mat, targeted at optimizing storage size and search  
<sup>15</sup>functionality without compromising on the ability to  
<sup>16</sup>update the stored data with new information.

#### <sup>17</sup>*Novel scoring algorithms*

<sup>18</sup>The current implementation allows the utilization of  
<sup>19</sup>both the model developed by Doench *et al.* [11], as  
<sup>20</sup>well as the one developed for this work. Functional-  
<sup>21</sup>ity aimed at facilitating the addition of new models  
<sup>22</sup>will be a part of future releases. Additionally, novel  
<sup>23</sup>scoring algorithms based on different sets of features  
<sup>24</sup>besides positional frequency of nucleotides will be ex-  
<sup>25</sup>plored in an attempt to provide better information to  
<sup>26</sup>the user and further facilitate the design of CRISPR  
<sup>27</sup>experiments.

#### <sup>29</sup>*Hardware scaling optimization*

<sup>30</sup>The server-class compute time of a genome-wide anal-  
<sup>31</sup>ysis, such as the one performed by CROPSR, is very  
<sup>32</sup>significant. Due to the modular, open source nature of  
<sup>33</sup>CROPSR, new parameters and models that get added  
<sup>34</sup>in the future will cause the current compute times to  
<sup>35</sup>get increased further. To mitigate some of this effect,  
<sup>36</sup>as well as help minimize issues associated with the long  
<sup>37</sup>compute times, future plans include a revision of the  
<sup>38</sup>CROPSR code to provide better hardware scaling op-  
<sup>39</sup>timization. In our opinion the benefits to be gained

from having the tool available now outweigh those of<sup>1</sup>  
waiting for non-essential optimization updates.<sup>2</sup>

## **Conclusions**

We have developed CROPSR, the first open source<sup>5</sup>  
tool for genome-wide design and evaluation of gRNA<sup>6</sup>  
sequences for CRISPR. In an effort to provide the<sup>7</sup>  
scientific community with a tool aimed at facilitat-<sup>8</sup>  
ing the design of genomics experiments by minimiz-<sup>9</sup>  
ing the out-of-lab time, CROPSR is capable of cre-<sup>10</sup>  
ating complete databases containing genome-wide in-<sup>11</sup>  
formation needed for CRISPR experiments, including<sup>12</sup>  
unique primer pairs for validation through PCR. The<sup>13</sup>  
improved scoring model adopted by CROPSR repre-<sup>14</sup>  
sents a significant improvement over currently widely<sup>15</sup>  
utilized methods. Additionally, CROPSR provides the<sup>16</sup>  
user with all information required to decide whether<sup>17</sup>  
the generated gRNAs are fit for an experiment in-<sup>18</sup>  
stead of deciding for them. This change will greatly<sup>19</sup>  
benefit crop scientists, as previously available scoring<sup>20</sup>  
models could be unreliable for complex crop genomes.<sup>21</sup>  
CROPSR allows data output as CSV and MongoDB,<sup>22</sup>  
with more formats being planned for future addition,<sup>23</sup>  
together with optimizations to reduce compute time<sup>24</sup>  
and novel scoring algorithms.<sup>25</sup>

## **Abbreviations**

CRISPR: Clustered regularly interspaced short palindromic repeats; Cas9: CRISPR associated nuclease; DNA: Deoxyribonucleic acid; RNA: Ribonucleic acid; SVM: Support vector machine; SVR: Support vector regressor; PCR: Polymerase chain reaction; CSV: Comma-separated values;

## **Acknowledgements**

The authors would like to thank Ghana Challa for his assistance with the file structure of Phytozome, especially the location of the functional annotation file.

## **Funding**

This work was funded by the DOE Center for Advanced Bioenergy and Bioproducts Innovation (U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research under Award Number DE-SC0018420). Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the author(s) and do not necessarily reflect the views of the U.S. Department of Energy.

#### 1 Availability of data and materials

2 The CROPSR package, tutorial and links to the tools used in this project  
3 are hosted at <https://github.com/cabbi-bio/cropsr>

#### 4 Author's contributions

5 HMP designed CROPSR, conducted all data analysis, wrote the manuscript  
6 and prepared figures. MH conceived CROPSR, obtained funding, supervised  
7 data analysis and software design and development, and co-wrote the  
8 manuscript. JH and DDI contributed with improvements to the code  
9 structure and functionality. All authors edited and approved the manuscript.

#### 9 Ethics approval and consent to participate

10 Not applicable.

#### 11 Consent for publication

12 Not applicable.

#### 13 Competing interests

14 The authors declare that they have no competing interests.

#### 15 Author details

16 <sup>1</sup> Illinois Informatics Institute, University of Illinois at Urbana-Champaign,  
17 Urbana, IL 61820, US. <sup>2</sup> Urbana, IL 61820, US. <sup>3</sup> Department of Crop  
18 Sciences, University of Illinois at Urbana-Champaign, Urbana, IL 61820,  
19 US. <sup>4</sup> National Center for Supercomputer Applications, University of  
20 Illinois at Urbana-Champaign, Urbana, IL 61820, US.

#### 21 References

1. Jinek M, Chylinski K, Fonfara I, Hauer M, Doudna JA, Charpentier E. A Programmable Dual-RNA-Guided DNA Endonuclease in Adaptive Bacterial Immunity. *Science*. 2012 08;337(6096):816. Available from: <http://science.sciencemag.org/content/337/6096/816.abstract>.
2. Cong L, Ran FA, Cox D, Lin S, Barretto R, Habib N, et al. Multiplex genome engineering using CRISPR/Cas systems. *Science (New York, NY)*. 2013 02;339(6121):819–823. Available from: <https://pubmed.ncbi.nlm.nih.gov/23287718>.
3. Mali P, Yang L, Esvelt KM, Aach J, Guell M, DiCarlo JE, et al. RNA-Guided Human Genome Engineering via Cas9. *Science*. 2013 02;339(6121):823. Available from: <http://science.sciencemag.org/content/339/6121/823.abstract>.
4. Gasiunas G, Barrangou R, Horvath P, Siksnys V. Cas9-crRNA ribonucleoprotein complex mediates specific DNA cleavage for adaptive immunity in bacteria. *Proceedings of the National Academy of Sciences of the United States of America*. 2012 09;109(39):E2579–E2586. Available from: <https://pubmed.ncbi.nlm.nih.gov/22949671>.
5. Zetsche B, Gootenberg JS, Abudayyeh OO, Slaymaker IM, Makarova KS, Essletzbichler P, et al. Cpf1 Is a Single RNA-Guided Endonuclease of a Class 2 CRISPR-Cas System. *Cell*. 2015 2021/07/28;163(3):759–771. Available from: <https://doi.org/10.1016/j.cell.2015.09.038>.
6. Montague TG, Cruz J, Gagnon JA, Church GM, Valen E. CHOPCHOP: a CRISPR/Cas9 and TALEN web tool for genome editing. *Nucleic Acids Research*. 2014;42(W1):W401–W407.
7. Labun K, Montague TG, Gagnon JA, Thyme SB, Valen E. CHOPCHOP v2: a web tool for the next generation of CRISPR genome engineering. *Nucleic Acids Research*. 2016;44(W1):W272–W276.
8. Labun K, Montague TG, Krause M, Torres Cleuren YN, Tjeldnes H, Valen E. CHOPCHOP v3: expanding the CRISPR web toolbox beyond genome editing. *Nucleic Acids Research*. 2019;47(W1):W171–W174.
9. Lei Y, Lu L, Liu HY, Li S, Xing F, Chen LL. CRISPR-P: A Web Tool for Synthetic Single-Guide RNA Design of CRISPR-System in Plants. *Molecular Plant*. 2014;7(9):1494–1496.
10. Liu H, Ding Y, Zhou Y, Jin W, Xie K, Chen LL. CRISPR-P 2.0: An Improved CRISPR-Cas9 Tool for Genome Editing in Plants. *Molecular Plant*. 2017 2020/07/28;10(3):530–532.
11. Doench JG, Hartenian E, Graham DB, Tothova Z, Hegde M, Smith I, et al. Rational design of highly active sgRNAs for CRISPR-Cas9-mediated gene inactivation. *Nature Biotechnology*. 2014;32(12):1262–1267.
12. Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, et al. Optimized sgRNA design to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nature Biotechnology*. 2016;34(2):184–191.
13. Parry MAJ, Madgwick PJ, Bayon C, Tearall K, Hernandez-Lopez A, Baudo M, et al. Mutation discovery for crop improvement. *Journal of Experimental Botany*. 2009 06;60(10):2817–2825. Available from: <https://doi.org/10.1093/jxb/erp189>.
14. Mitros T, Session AM, James BT, Wu GA, Belaffif MB, Clark LV, et al. Genome biology of the paleotetraploid perennial biomass crop *Miscanthus*. *Nature Communications*. 2020;11(1):5442. Available from: <https://doi.org/10.1038/s41467-020-18923-6>.
15. Koressaar T, Remm M. Enhancements and modifications of primer design program Primer3. *Bioinformatics*. 2007;23(10):1289–1291.
16. Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, et al. Primer3–new capabilities and interfaces. *Nucleic acids research*. 2012 08;40(15):e115–e115.
17. Breslauer KJ, Frank R, Blöcker H, Marky LA. Predicting DNA duplex stability from the base sequence. *Proceedings of the National Academy of Sciences*. 1986;83(11):3746–3750.
18. Rychlik W, Spencer WJ, Rhoads RE. Optimization of the annealing temperature for DNA amplification in vitro. *Nucleic Acids Research*. 1990 11;18(21):6409–6412.
19. Breslauer KJ, Frank R, Blöcker H, Marky LA. Predicting DNA duplex stability from the base sequence. *Proceedings of the National Academy of Sciences*. 1986 06;83(11):3746. Available from: <http://www.pnas.org/content/83/11/3746.abstract>.
20. Owczarzy R, Vallone PM, Gallo FJ, Paner TM, Lane MJ, Benight AS. Predicting sequence-dependent melting stability of short duplex DNA oligomers. *Biopolymers*. 1997;44(3):217–239. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291097-0282%281997%2944%3A3%3C217%3A%3AAID-BIP3%3E3.0.CO%3B2-Y>.
21. SantaLucia J. A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics. *Proceedings*

- 1 of the National Academy of Sciences. 1998;95(4):1460–1465.
222. BOLTON ET, McCARTHY BJ. A general method for the isolation of  
 3 RNA complementary to DNA. Proceedings of the National Academy  
 4 of Sciences of the United States of America. 1962 08;48(8):1390–1397.  
 Available from: <https://pubmed.ncbi.nlm.nih.gov/13870855>.
523. Dong J, Zielinski RE, Hudson ME. t-SNAREs bind the Rhg1  $\alpha$ -SNAP  
 6 and mediate soybean cyst nematode resistance. The Plant Journal.  
 2020;104(2):318–331. Available from:  
 7 <https://onlinelibrary.wiley.com/doi/abs/10.1111/tpj.14923>.
824. Hsu PD, Scott DA, Weinstein JA, Ran FA, Konermann S, Agarwala V,  
 9 et al. DNA targeting specificity of RNA-guided Cas9 nucleases. Nature  
 10 biotechnology. 2013 09;31(9):827–832. Available from:  
<https://pubmed.ncbi.nlm.nih.gov/23873081>.

## 11 Figures

12

13 **Figure 1 Overview of a typical CRISPR/Cas9 gene editing**  
 14 **experiment.** (a) Overview of CRISPR/Cas9 mechanism used  
 15 to create deletions in crop genomes. (b) Diagram of a typical  
 16 knockout editing experiment in a crop plant, with associated  
 17 timeline. Improvements in the steps contained in gray blocks  
 are anticipated from the CROPSR software.

18

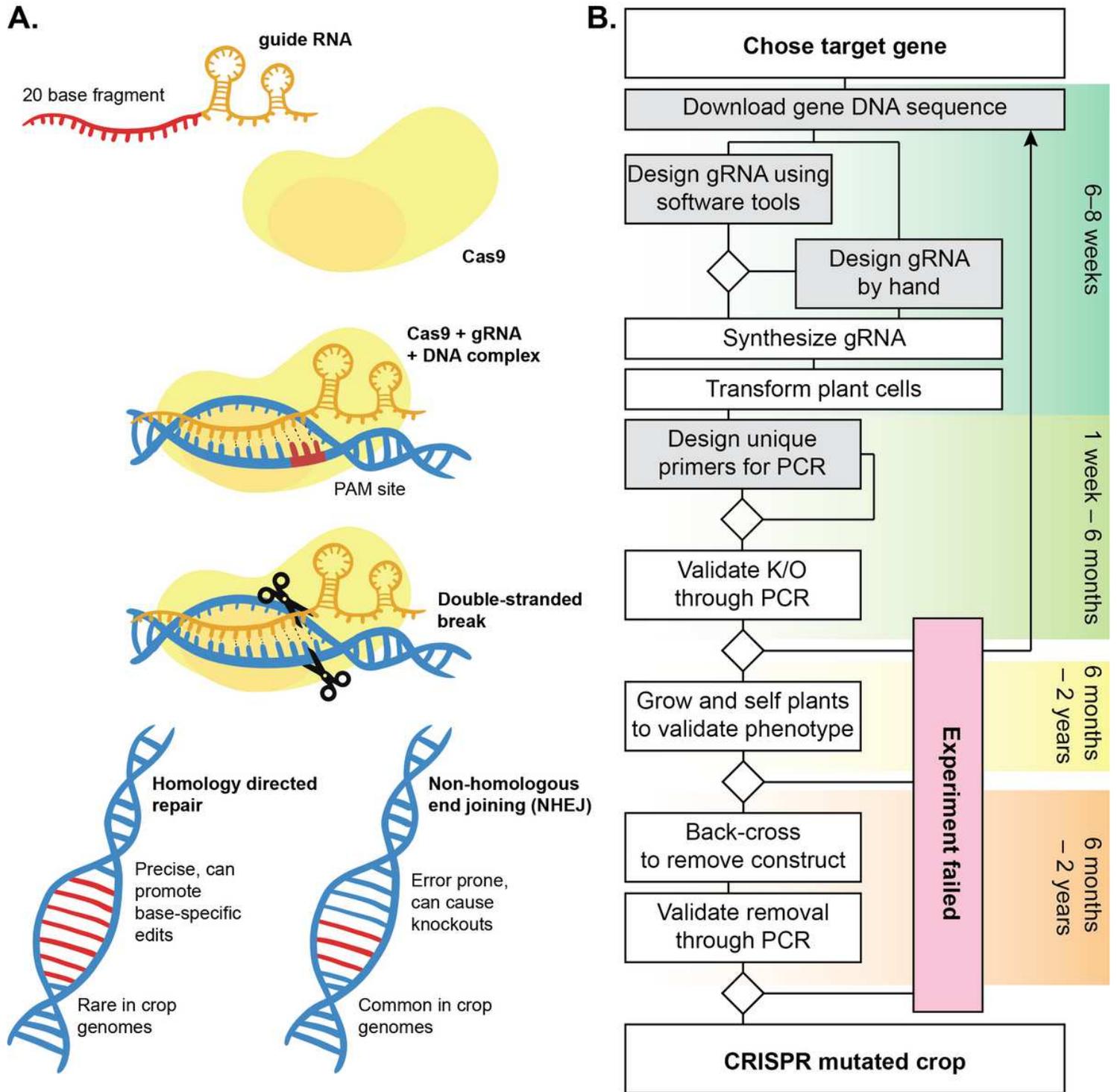
19

20 **Figure 2 Functional block diagram of CROPSR modules.** (a)  
 21 The different input data files (FASTA, GFF, Phytozome  
 22 annotation file) are imported and processed by multiple  
 23 modular programs within the CROPSR suite. The genome  
 24 sequence is submitted to the gRNA design program (shown in  
 25 detail in b), and the output is placed in a MongoDB database  
 26 (or optionally a CSV file). The GFF file, and Phytozome  
 27 annotation file when applicable, are processed by a separate  
 28 program, and then each entry in the database is updated with  
 29 functional annotation to be used for search queries. Unique  
 30 primer pairs are designed for each gRNA database entry. (b)  
 31 The gRNA module takes data from the file manager module  
 32 (which parses a FASTA input sequence file), and generates a  
 33 list of location pairs (5' – 3') for every PAM site match. The  
 34 sequence, strand, start and end positions and CRISPR system  
 for each guide are stored, and a score representing expected  
 35 performance of each potential gRNA is calculated utilizing one  
 of the available algorithms. Final data for each guide is then  
 added to the database to be associated with functional  
 annotation and PCR primers for validation.

7 **Figure 3 Comparison of the scoring performance of**  
 8 **CROPSR with the Chopchop algorithm.** (a) Density plot of  
 9 the score generated by the Chopchop scoring algorithm  
 10 against the "gene % rank", a ranking of  
 11 experimentally-determined relative performance of gRNAs on  
 12 a per-gene basis. (b) Density plot of the CROPSR scoring  
 13 algorithm against the gene % rank. (c) Binned scatter + box  
 14 plot of the Chopchop scoring algorithm against the gene %  
 15 rank. The gRNA targeted for experimental use by Chopchop  
 are those in the 80-100% bin. (d) Binned scatter + box plot of  
 the CROPSR scoring algorithm against the gene % rank.

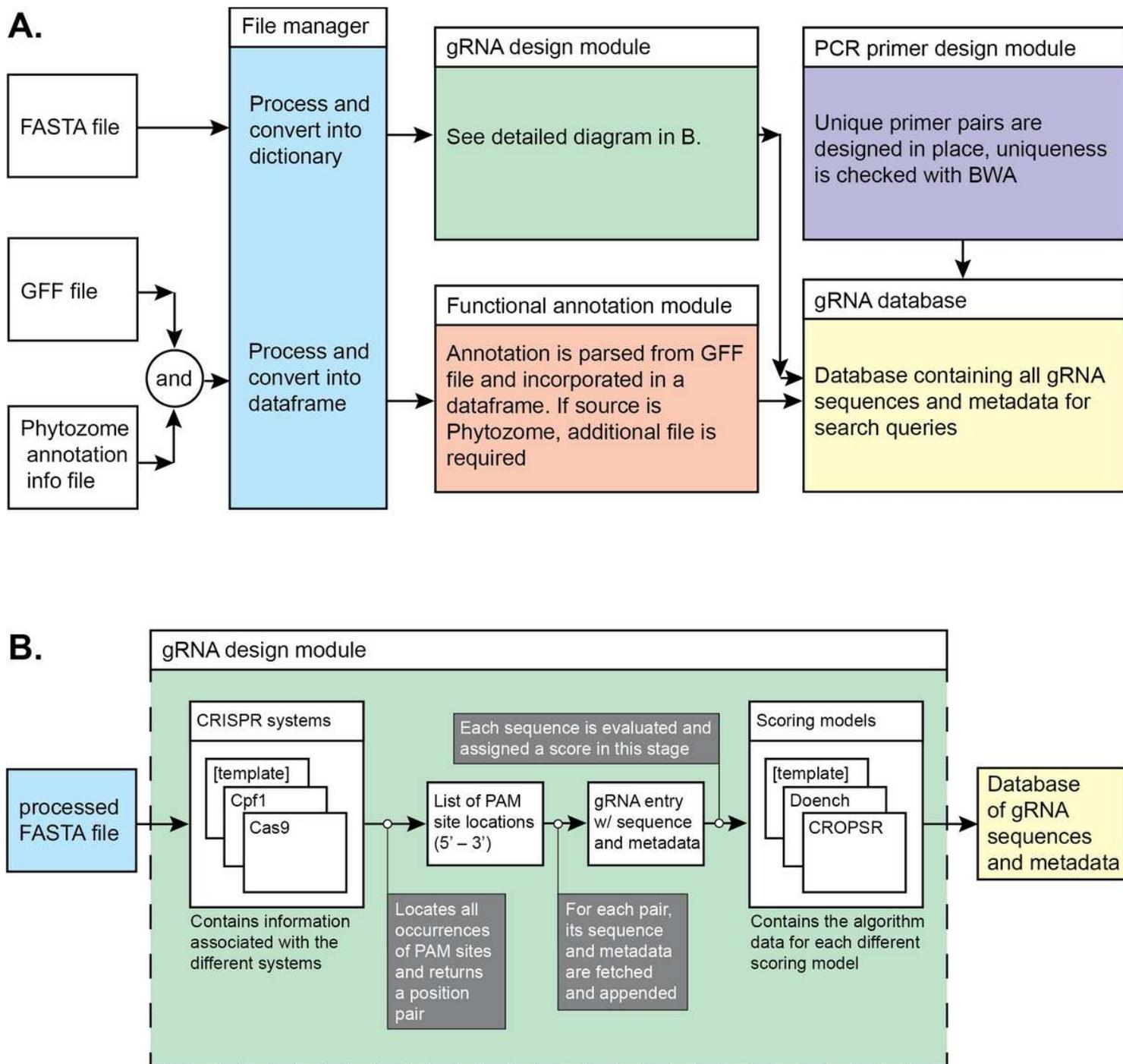
28 **Figure 4 Overview of a CRISPR experiment using CROPSR**  
 29 **Timeline and steps of a typical CRISPR/Cas9 knockout**  
 30 **experiment in a crop plant genome, utilizing CROPSR.** Steps  
 31 contained in gray blocks represent steps that only need to be  
 32 done once per genome, at the first utilization of CROPSR  
 33 (database generation). Consecutive uses on the same genome  
 require only a database search, as shown.

# Figures



**Figure 1**

Overview of a typical CRISPR/Cas9 gene editing experiment. (a) Overview of CRISPR/Cas9 mechanism used to create deletions in crop genomes. (b) Diagram of a typical knockout editing experiment in a crop plant, with associated timeline. Improvements in the steps contained in gray blocks are anticipated from the CROPSR software



**Figure 2**

Functional block diagram of CROPSR modules. (a) The different input data files (FASTA, GFF, Phytozome annotation file) are imported and processed by multiple modular programs within the CROPSR suite. The genome sequence is submitted to the gRNA design program (shown in detail in b), and the output is placed in a MongoDB database (or optionally a CSV file). The GFF file, and Phytozome annotation file when applicable, are processed by a separate program, and then each entry in the database is updated with functional annotation to be used for search queries. Unique primer pairs are designed for each gRNA database entry. (b) The gRNA module takes data from the file manager module (which parses a FASTA input sequence file), and generates a list of location pairs (5' - 3') for every PAM site match. The

sequence, strand, start and end positions and CRISPR system for each guide are stored, and a score representing expected performance of each potential gRNA is calculated utilizing one of the available algorithms. Final data for each guide is then added to the database to be associated with functional annotation and PCR primers for validation.

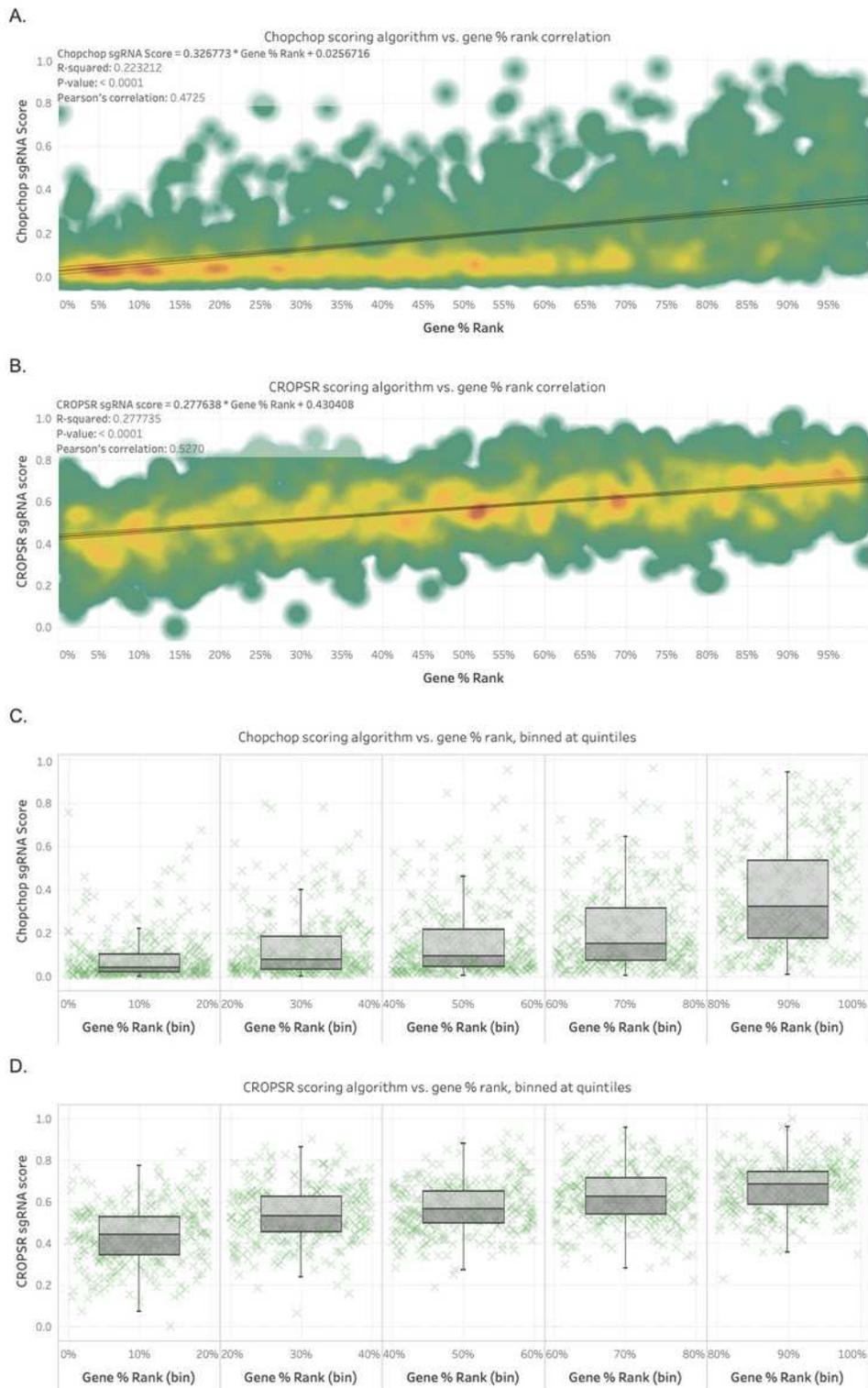
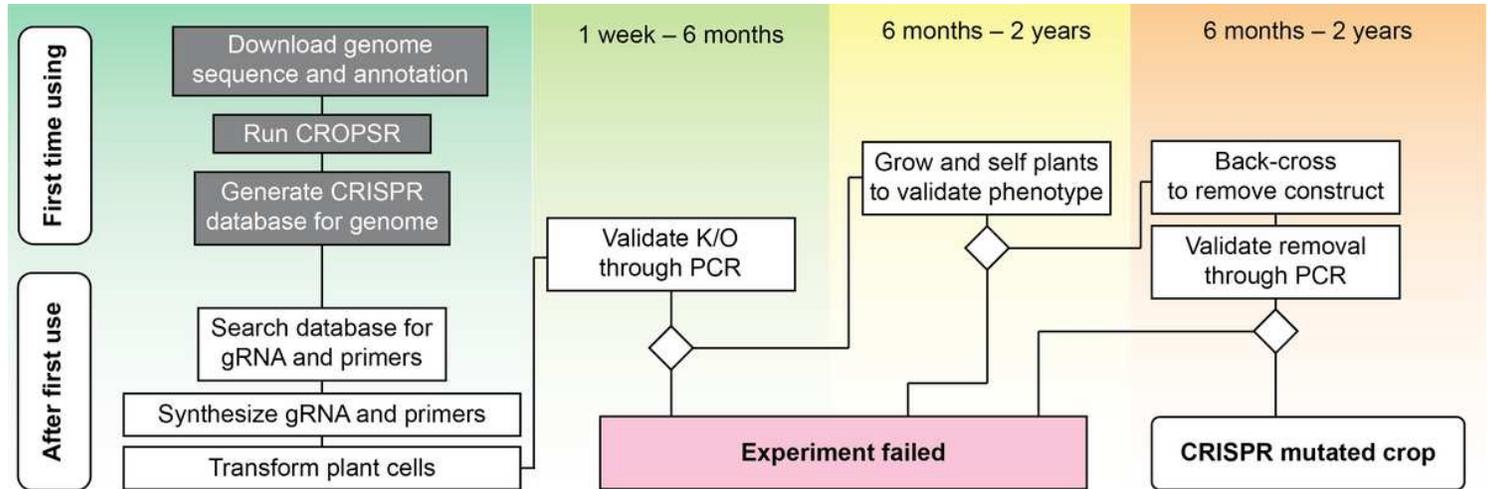


Figure 3

Comparison of the scoring performance of CROPSR with the Chopchop algorithm. (a) Density plot of the score generated by the Chopchop scoring algorithm against the "gene % rank", a ranking of experimentally-determined relative performance of gRNAs on a per-gene basis. (b) Density plot of the CROPSR scoring algorithm against the gene % rank. (c) Binned scatter + box plot of the Chopchop scoring algorithm against the gene % rank. The gRNA targeted for experimental use by Chopchop are those in the 80-100% bin. (d) Binned scatter + box plot of the CROPSR scoring algorithm against the gene % rank.



**Figure 4**

Overview of a CRISPR experiment using CROPSR Timeline and steps of a typical CRISPR/Cas9 knockout experiment in a crop plant genome, utilizing CROPSR. Steps contained in gray blocks represent steps that only need to be done once per genome, at the first utilization of CROPSR (database generation). Consecutive uses on the same genome require only a database search, as shown.