# A peer-reviewed version of this preprint was published in PeerJ on 18 November 2015.

<u>View the peer-reviewed version</u> (peerj.com/articles/cs-33), which is the preferred citable publication unless you specifically need to cite this preprint.

Roberson ED. 2015. Identification of high-efficiency 3'GG gRNA motifs in indexed FASTA files with ngg2. PeerJ Computer Science 1:e33 <a href="https://doi.org/10.7717/peerj-cs.33">https://doi.org/10.7717/peerj-cs.33</a>

# Identification of high-efficiency 3'GG gRNA motifs in indexed FASTA files with ngg2

- 3 Elisha D.O. Roberson<sup>1,2,\*</sup>
- <sup>1</sup>Department of Internal Medicine, Division of Rheumatology, Washington University,
- 5 St. Louis, MO, USA.
- 6 <sup>2</sup>Department of Genetics, Washington University, St. Louis, MO, USA.
- 7 \*Elisha D.O. Roberson, Ph.D.
- 8 Washington University
- 9 Depts. of Internal Medicine and Genetics,
- 10 Division of Rheumatology
- 11 660 South Euclid Ave.
- 12 Campus Box 8045
- 13 St. Louis, MO 63110
- 14 <u>eroberso@dom.wustl.edu</u>

#### Abstract

- 16 CRISPR/Cas9 is emerging as one of the most used methods of genome modification in
- organisms ranging from bacteria to human cells. However, the efficiency of editing
- 18 varies tremendously site-to-site. A recent report identified a novel motif, called the
- 19 3'GG motif, which substantially increases the efficiency of editing at all sites tested.
- 20 Furthermore, they highlighted that previously published gRNAs with high editing
- 21 efficiency also had this motif. I designed a python command-line tool, ngg2, to identify
- 22 3'GG gRNA sites from indexed FASTA files. As a proof-of-concept, I screened for these
- 23 motifs in six genomes: Saccharomyces cerevisiae, Caenorhabditis elegans, Drosophila
- 24 melanogaster, Danio rerio, Mus musculus, and Homo sapiens. I identified more than 24
- 25 million single match 3'GG motifs in these reference genomes. Greater than 87% of all
- 26 protein coding genes in the six reference genomes had at least one overlapping unique
- 27 3'GG gRNA site. In particular, more than 96% of mouse and 99% of human protein
- coding genes have at least one unique, overlapping 3'GG gRNA. These identified sites
- 29 can be used as a starting point in gRNA design, and the ngg2 tool provides an
- 30 important ability to identify high-efficiency editing sites in non-model species.

#### Introduction

- 2 Genome engineering allows for the targeted deletion or modification by homology
- 3 directed repair of a target locus. Currently, one of the most popular methods for
- 4 genome manipulation is the clustered regularly interspaced short palindromic repeat
- 5 (CRISPR) / CRISPR associated protein 9 (Cas9) system adapted from *Streptococcus*
- 6 pyogenes. The CRISPR/Cas system was initially thought to represent a novel DNA repair
- 7 mechanism, but was eventually found to provide heritable bacterial immunity to
- 8 invading exogenous DNA from sources, such as plasmids and bacteriophages
- 9 (Barrangou et al. 2007; Makarova et al. 2006). During endogenous CRISPR/Cas9
- 10 function, foreign DNA integrates into the CRISPR locus. The bacterial cell then
- 11 expresses the pre-CRISPR RNA (crRNA) and a trans-activating crRNA (tracrRNA) that
- pair to form a complex that is cleaved by RNAse III (Deltcheva et al. 2011). The
- resulting RNA is a hybrid of the pre-crRNA and the tracrRNA, and includes a 20 bp
- 14 guide RNA (gRNA) sequence. The gRNA is incorporated into Cas9 and can then guide
- 15 the cleavage of a complementary DNA sequence by the nuclease activity of the Cas9
- 16 protein. The topic of CRISPR-Cas genome editing has been reviewed extensively
- elsewhere (Doudna & Charpentier 2014; Hsu et al. 2014; Jiang & Doudna 2015; Mali et
- 18 al. 2013).
- 19 Codon-optimized versions of Cas9 are available for a wide range of organisms, and can
- 20 easily be synthesized if it is not already available. Transfecting cells with Cas9 plasmid
- 21 along with a fused crRNA-tracrRNA hybrid construct called a single-guide RNA
- 22 (sgRNA) allows for temporary activity of Cas9. Keeping a stock of plasmids with a
- 23 sgRNA backbone minus the gRNA site makes it easy to quickly generate new sgRNA
- 24 plasmids by site-directed mutagenesis. The Cas9 protein loaded with the sgRNA will
- 25 bind to sites complementary genomic loci, but will only cut it if a protospacer adjacent
- 26 motif (PAM) site immediately follows the complementary sequence (Mojica et al. 2009).
- 27 The PAM site for type-I CRISPR is an NGG. Therefore, a gRNA site can be defined as
- N<sub>20</sub>NGG. It is important to note that constitutively expressed sgRNAs typically use a U6
- 29 snRNA promoter that strongly prefers a G starting base. For U6 compatibility,
- 30 sequences starting with A, C, or T may be used if they are cloned into a sgRNA vector
- 31 with an appended G base, resulting in a 21 bp gRNA (Farboud & Meyer 2015; Ran et al.
- 32 2013b). The subset gRNA sites contain a starting G base (GN<sub>19</sub>NGG), and can be cloned
- as a 20 bp gRNA, which I will refer to as canonical gRNA sites.

19

- 1 The rate of editing using the CRISPR/Cas9 system is far higher than homologous
- 2 recombination, but higher efficiency is still desirable. The introduction of a longer stem
- 3 in part the sgRNA stem-loop structure and the flip of a single A in a polyA track of a
- 4 separate sgRNA stem-loop, called the flip + extension (F+E) sgRNA design, resulted in
- 5 increased Cas9 editing efficiency (Chen et al. 2013). Recently, another improvement was
- 6 reported that increases efficiency. gRNA sites with a GG motif adjacent to the PAM site,
- 7 called 3'GG gRNAs, have far higher activity than equivalent gRNA sites in the same
- 8 region (Farboud & Meyer 2015). These sites take the form of N<sub>18</sub>GGNGG.
- 9 In this manuscript, I report a python command-line tool, ngg2, for identification of
- 10 high-efficiency 3'GG gRNA motifs from indexed FASTA files. Tools already exist to
- 11 identify Cas9 gRNA targets in common model organisms. However, support for less
- 12 common and non-model organisms is limited. This tool will enable the easy
- identification of high-efficiency gRNA sites in any genome. As a proof of concept, I
- 14 report all 3'GG gRNA motifs in 6 model species, identifying more than 35 million sites,
- of which more than 24 million are unique matches within the reference genome for that
- species. More than 90% of all protein coding genes in 5/6 species have at least one
- 17 unique 3'GG gRNA overlapping it for potential editing.

#### Materials & Methods

#### ngg2 Motif identification

- 20 I designed ngg2 using python with compiled regular expressions for the 3'GG gRNA
- 21 plus PAM motif. ngg2 identifies these sites on both the sense and antisense strands, and
- optionally can be restricted to only canonical sites starting with a G for 20 bp gRNAs.
- 23 Sites can be identified for a specified region, a whole contig, or all contigs in the input
- 24 FASTA file. ngg2 uses the FASTA index to directly seek the genomic target without
- 25 reading the entire file. This tool only identifies potential editing sites, and does not
- 26 report uniqueness of the gRNA. As part of this manuscript, I generated files listing the
- 27 uniqueness of each gRNA for each species, including gene overlaps. I recommend
- 28 checking any identified site with blast / blat to ensure the minimum number of potential
- off target sites. ngg2 output includes the contig name, start and end positions, the
- 30 gRNA sequence, the PAM sequence, and whether the site starts with a G.

#### Multi-species site identification

- 2 I used ngg2 to identify all 3' GG gRNA motifs 6 commonly studied organisms:
- 3 Saccharomyces cerevisiae, Caenorhabditis elegans, Drosophila melanogaster, Danio rerio, Mus
- 4 musculus, and Homo sapiens. I used a GNU Make script for genome downloads and site
- 5 identification to enable reproducibility. The Makefile downloads the top-level (Danio
- 6 rerio, Drosophila melanogaster, Caenorhabditis elegans, and Saccharomyces cerevisiae) or
- 7 primary assembly (Homo sapiens, Mus musculus) genomes from Ensembl Release 79,
- 8 indexes the FASTA files with samtools, downloads gene annotations for later
- 9 intersection, runs ngg2 on all contigs for each FASTA file, and calculates GC content for
- 10 each genome. I did not allow for 'N' bases in the gRNA sequence or PAM site. I based
- 11 the GC content on only non-N base content in each genome. I counted occurrences of
- 12 each gRNA sequence within each species, annotated the overlapping genes for each
- 13 gRNA, and determined the number of genes cut by any 3'GG gRNAs as well as by
- unique gRNAs using R (v3.1.2) with the GenomicRanges package (v1.18.4) (Lawrence et
- al. 2013; R Core Team 2014). I calculated all summary statistics, counts, and figures
- using RStudio (v0.98.1102) Markdown with knitr (Xie 2013). I used many accessory
- 17 functions from plyr, dplyr, tidry, stringr, and magrittr in the initial analysis, and I
- 18 generated all plots with ggplot2.

#### Results

19

#### 20 **3'GG gRNA sites are common in each species**

- 21 Overall, I identified greater than 35 million 3'GG gRNA sites in the six reference
- 22 genomes (Table 1). Some of these gRNA sequences were not unique in a genome,
- $\,$  leaving more than 24 million unique 3'GG sites. Approximately 6 million of the 24  $\,$
- 24 million unique sites were canonical G starting motifs. The sites identified in each
- 25 species with the gRNA sequence, PAM sequence, genome coordinates, annotated
- overlapping genes, and number of perfect genome matches are available for download
- 27 (Roberson 2015). The R scripts, python files, and Make files are also available in a public
- 28 repository for reproducibility.
- 29 The genomes I analyzed had vastly different sizes, ranging from approximately 12 Mb
- 30 for yeast to greater than 3 Gb for humans, and as a result had dramatically different
- 31 numbers of 3'GG gRNA sites per genome. Therefore, I also assessed the site density per
- 32 megabase of reference genome size (Table 2). Unique sites averaged a density of 3,091

- sites / Mb, or 1 unique site per 324 bp. *D. rerio* had the lowest density at 1,733 unique
- 2 sites / Mb, while *D. melanogaster* had the highest density at 4,062 unique sites / Mb. The
- 3 low density of unique sites in zebrafish may be due to genome complexity from
- 4 previous duplication events.

#### 5 Little strand bias observed for 3'GG gRNA sites

- 6 The strand of each gRNA site with respect to the reference was included in the ngg2
- 7 output files. I plotted the split of sense / antisense sites in each genome to visualize
- 8 strong deviations from an expection of no strand bias (Fig. 1). For each organism, I
- 9 considered every gRNA site as an independent Bernoulli trial with a 50% probability of
- success, and considered a "Sense" strand designation as a trial successful outcome
- 11 (Supp. Table 1). *C. elegans, D. melanogaster,* and *H. sapiens* all demonstrated a strand bias
- 12 significantly different from the expected ratio for all 3'GG sites. However, while the
- difference is significant, it may be unimportant. Wildtype Cas9 cleaves both DNA
  - strands simultaneously, and therefore the strand of the target sequence doesn't matter.
- 15 Strategies that employ dual nickases to reduce off target effects could be affected by
- such bias, as they require two separate gRNA sites on opposite strands (Ran et al.
- 17 2013a). The difference observed is only less than 0.4% different from expected 50% ratio,
- and whether this functionally affects the ability to choose paired 3'GG gRNAs remains
- 19 to be seen.

# 20 CGG & GGG PAM sites are underrepresented

- 21 I visualized the distribution of the four PAM sites (AGG, CGG, GGG, TGG) as a stacked
- 22 bar chart of each sites proportion of the total identified sites in each species (Fig. 2). In
- 23 general, the AGG and TGG sites represented the majority of 3'GG gRNA sites in all
- 24 species. I tested whether PAM site distribution differed from chance based on the GC
- 25 content of the reference genome. For each species, I considered each PAM site a
- 26 Bernoulli trial, and defined success as either CGG or GGG site identity. The probability
- of success was set equal to the estimated genome-wide GC content calculated from the
- 28 reference genome, excluding N (Supp. Table 2). Only D. melanogaster met the expected
- 29 GC success rate for 3'GG gRNA sites. The rate of picking a CGG or GGG PAM was less
- 30 than the genome GC content in *S. cerevisiae*, *D. rerio*, *M. musculus*, and *H. sapiens*. In
- 31 particular, the estimate for both *M. musculus* and *H. sapiens* is approximately 20% lower
- 32 than the average genome GC content. This is not necessarily unexpected. The CGG
- 33 PAM site includes a 5' CpG dinucleotide that is generally underrepresented due to the
- relatively high frequency of methyl-cytosine deamination to thymine in this context. *C.*

- 1 *elegans* was the exception, with CGG and GGG PAM selection greater than the expected
- 2 frequency. However, this species lacks DNA methylation and would not necessarily be
- at an advantage to limit CpG dinucleotides.

#### 4 Most protein coding genes overlap at least one unique 3'GG gRNA

- 5 A common use of genome engineering is to knock out or otherwise modify the function
- 6 of a protein coding gene. The efficiency of such edits is critical, as just introducing
- 7 frame-shifting mutations can require screening a large number single-cell clones or
- 8 derived animals to identify a successful edit. As part of this study, I annotated for each
- 9 gRNA in the 6 species if there was any overlap with a gene. Conversely, I also annotate
- 10 a count of how many of each of the four classes (all sites, all unique sites, canonical
- sites, and unique canonical sites) overlap every gene. No less than 87% of any species'
- genes overlap at least one unique 3'GG gRNA (Fig. 3, Supp. Table 3). This catalog of
- potential sites demonstrates that most protein coding genes can be targeted by at least
- one 3'GG gRNA site to achieve high editing efficiency.

# **Discussion**

- 16 In this manuscript, I have described a new tool for identifying 3'GG gRNA sites and
- 17 presented a catalog of potential editing sites in 6 species. Importantly, many genomic
- 18 loci can be targeted by unique 3'GG gRNA sites for efficient genome modification. Blast
- 19 / blat searching of gRNA sequences is critical to ensure that there is not an abundance of
- 20 exact or near matches in the target genome. It is also important to consider the target
- 21 genome's specific genotypes when designing a gRNA. In particular, variants that alter
- 22 PAM sites away from NGG will not be cleaved by Cas9 even if the gRNA is an exact
- 23 match. Identification of potential gRNA sites from reference genomes is a starting point
- 24 for genome editing, but careful study of target sites is required before a final design is
- selected.
- 26 The accuracy of editing can be improved by using two gRNAs and a mutant Cas9
- 27 nickase. I observed some significant, but low-effect strand bias in these genomes. This
- 28 may lead to some loci not being compatible with paired 3'GG gRNA sites. When
- 29 possible, choosing paired 3'GG gRNA sites should be strongly considered. Efficiencies
- of less than 10% were increased to 50% efficiency or greater by using the 3'GG strategy
- 31 (Farboud & Meyer 2015). As such, using paired 3'GG gRNAs with a nickase may give
- 32 the best of both worlds with both high accuracy and high efficiency.

- 1 It is important to note that ngg2 will operate on any indexed FASTA file. Many gRNA
- 2 site finding tools are limited to catalogs of gRNA sites in model organisms. This tool
- 3 fills an important gap for individuals working outside of commonly used species,
- 4 allowing for rapid and accurate identification of high efficiency 3'GG gRNA sites. The
- 5 provided gRNA site survey and associated tool, ngg2, represent a valuable resource for
- 6 designing genomic modification strategies.

# Acknowledgments

- 8 This work was performed in the Human Genomics and Bioinformatics Facility of the
- 9 Rheumatic Disease Core Center at Washington University (P30 AR048335). I wish to
- 10 thank Dr. Li Cao for her helpful comments during the preparation of this manuscript.

11

3

# 1 Figures

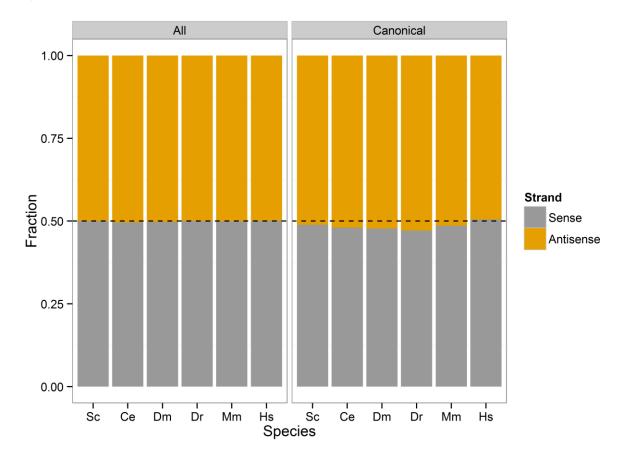


Fig. 1 - gRNA strand compared to reference genome

- 4 The fraction of gRNA sites sense and antisense to the reference are shown for each
- 5 species. The x-axis labels are the two-letter species abbreviations. When considering all
- $\,\,$  gRNA sites, there is not dramatic bias from a 50/50 strand ratio. The canonical sites are
- 7 more skewed toward being antisense to the reference.

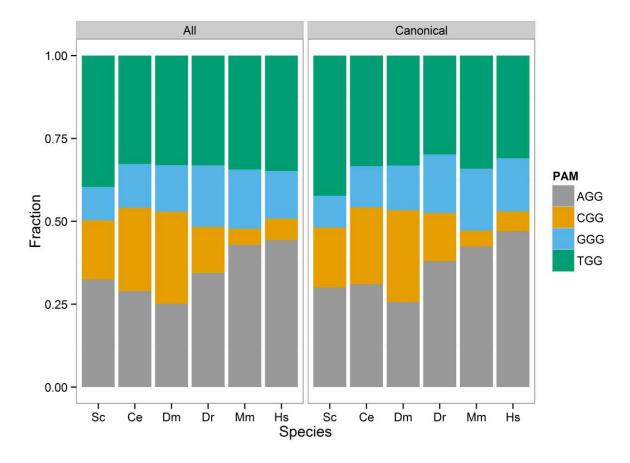


Fig. 2 - PAM site usage

- 3 Each species has four potential protospacer adjacent motifs (PAM) possible for
- 4 identified gRNA sites. The stacked bar chart shows the fraction of all PAM sites each
- 5 motif occupies. The CGG motif, that includes a CpG dinucleotide, is the least prevalent
- 6 motif in the zebrafish, mouse, and human genomes.

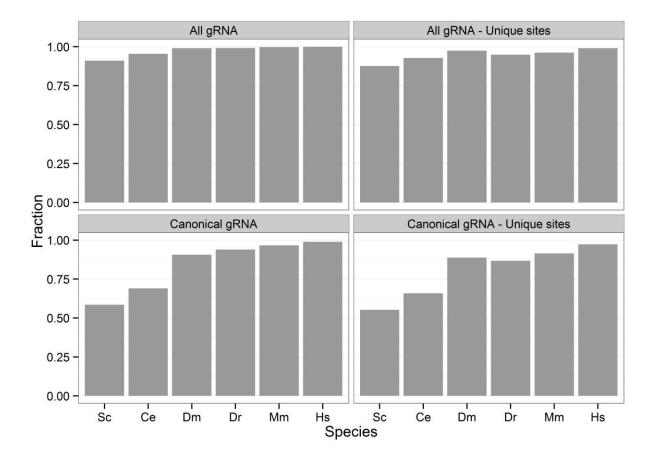


Fig. 3 - Overlap of protein coding genes and 3'GG gRNA sites

- 3 The fraction of protein coding genes overlapped by at least one gRNA of each class are
- 4 shown for each species. The majority of protein coding genes have at least one unique
- 5 overlapping 3'GG gRNA site.

#### 1 Tables

	All gl	RNAs	Canonical gRNAs		
	All	Unique	All	Unique	
Saccharomyces cerevisiae	38,430	35,733	8,328	8,148	
Caenorhabditis elegans	293,597	257,317	62,874	60,669	
Drosophila melanogaster	672,135	583,799	162,986	159,412	
Danio rerio	4,422,730	2,448,177	626,202	561,214	
Mus musculus	13,436,734	9,965,896	2,591,654	2,482,936	
Homo sapiens	16,454,683	11,145,670	2,934,283	2,797,025	
Total	35,318,309	24,436,592	6,386,327	6,069,404	

#### 2 Table 1 - Count of gRNA classes in each species

- 3 All N<sub>18</sub>GGNGG motifs are included in the All category for All gRNAs. However,
- 4 multiple edit sites in the genome for one gRNA are typically disadvantageous. Unique
- 5 gRNAs were only observed in a gRNA plus PAM context once in a given genome.

	All gRNAs		<b>Canonical gRNAs</b>		
<b>Species</b>	All	Unique	All	Unique	
Saccharomyces cerevisiae	3,161.11	2,939.27	685.03	670.23	
Caenorhabditis elegans	2,927.59	2,565.82	626.94	604.96	
Drosophila melanogaster	4,676.50	4,061.89	1,134.01	1,109.14	
Danio rerio	3,131.21	1,733.27	443.34	397.33	
Mus musculus	4,920.31	3,649.35	949.02	909.21	
Homo sapiens	5,308.39	3,595.67	946.62	902.34	

# 7 Table 2 - 3'GG gRNA Sites per Megabase Genome Size

- 8 Reference genome size was determined from the FASTA index. The number of unique
- 9 3'GG gRNA sites in the genome is encouraging, with an average across all species of
- 10 one site per 324 bp.

# 1 Supplementary

		All gRNAs		Canonical gRNAs			
Species	estimate	p.value	p.adj	estimate	p.value	p.adj	
Saccharomyces cerevisiae	0.500	9.55E-01	1.00E+00	0.489	3.90E-02	1.56E-01	
Caenorhabditis elegans	<u>0.496</u>	<u>1.44E-05</u>	<u>8.66E-05</u>	<u>0.481</u>	2.29E-24	<u>1.60E-23</u>	
Drosophila melanogaster	<u>0.498</u>	<u>1.82E-04</u>	<u>9.12E-04</u>	0.478	2.78E-79	2.78E-78	
Danio rerio	0.500	1.41E-01	4.19E-01	0.471	3.95E-323	4.74E-322	
Mus musculus	0.500	8.50E-01	1.00E+00	0.486	6.92E-323	7.61E-322	
Homo sapiens	<u>0.501</u>	3.59E-25	2.87E-24	0.505	6.85E-70	6.16E-69	

- 2 Supp. Table 1 Strand preference for gRNA sites
- The estimate is the binomial estimate of choosing the sense strand for a given gRNA,
- 4 and the p-values are the significance of the strand bias. Half of the examined genomes
- 5 show statistically significant, but low effect strand bias.

		All gRNAs			Canonical gRNAs			
<b>Species</b>	GC	estimate	p.value	p.adj	estimate	p.value	p.adj	
Saccharomyces cerevisiae	0.382	<u>0.279</u>	<u>4.94E-324</u>	<u>5.93E-323</u>	<u>0.277</u>	<u>9.77E-96</u>	<u>3.91E-95</u>	
Caenorhabditis elegans	0.354	<u>0.384</u>	<u>7.99E-236</u>	<u>4.00E-235</u>	0.356	3.40E-01	3.40E-01	
Drosophila melanogaster	0.417	0.418	1.16E-01	2.33E-01	<u>0.412</u>	<u>1.56E-05</u>	<u>4.69E-05</u>	
Danio rerio	0.366	<u>0.325</u>	<u>7.41E-323</u>	<u>6.92E-322</u>	<u>0.322</u>	3.46E-323	3.80E-322	
Mus musculus	0.416	<u>0.228</u>	<u>1.33E-322</u>	<u>9.34E-322</u>	<u>0.235</u>	<u>6.92E-323</u>	<u>6.92E-322</u>	
Homo sapiens	0.407	<u>0.209</u>	1.48E-322	9.34E-322	<u>0.220</u>	6.92E-323	6.92E-322	

- 1 Supp. Table 2 PAM site frequency compared to genome GC content
- 2 The average genome GC content and the estimated chance of picking a GC PAM site
- 3 (CGG or GGG) is shown for each species. GC content was calculated from the
- 4 downloaded reference files. For all 3'GG gRNA sites, only *D. melanogaster* had no strand
- 5 bias. Most genomes show significantly fewer CGG and GGG PAM sites than expected
- 6 based on genome GC content.

	All g	RNAs	Canonical gRNAs		
<b>Species</b>	All	Unique	All	Unique	
Saccharomyces cerevisiae	0.911	0.876	0.585	0.552	
Caenorhabditis elegans	0.954	0.928	0.690	0.658	
Drosophila melanogaster	0.991	0.975	0.906	0.888	
Danio rerio	0.992	0.949	0.939	0.867	
Mus musculus	0.998	0.962	0.967	0.915	
Homo sapiens	0.999	0.990	0.990	0.973	

#### 1 Supp. Table 3 - Fraction protein coding genes with overlapping gRNAs

- 2 The proportion of protein-coding genes with at least one gRNA of the given type
- 3 overlapping the annotated transcript start and end are shown for each species. Both *C*.
- 4 elegans and S. cerevisiase have a lower rate of overlap with unique gRNAs. However,
- 5 even the lowest fraction of unique overlap (yeast) still has >87% of protein-coding genes
  - with at least one unique 3'GG gRNA site.

#### References

- 3 Barrangou R, Fremaux C, Deveau H, Richards M, Boyaval P, Moineau S, Romero DA,
- 4 and Horvath P. 2007. CRISPR Provides Acquired Resistance Against Viruses in
- 5 Prokaryotes. *Science* 315:1709-1712.
- 6 Chen B, Gilbert Luke A, Cimini Beth A, Schnitzbauer J, Zhang W, Li G-W, Park J,
- 7 Blackburn Elizabeth H, Weissman Jonathan S, Qi Lei S, and Huang B. 2013. Dynamic
- 8 Imaging of Genomic Loci in Living Human Cells by an Optimized CRISPR/Cas System.
- 9 *Cell* 155:1479-1491.
- 10 Deltcheva E, Chylinski K, Sharma CM, Gonzales K, Chao Y, Pirzada ZA, Eckert MR,
- 11 Vogel J, and Charpentier E. 2011. CRISPR RNA maturation by trans-encoded small
- 12 RNA and host factor RNase III. *Nature* 471:602-607.
- 13 Doudna JA, and Charpentier E. 2014. The new frontier of genome engineering with
- 14 CRISPR-Cas9. Science 346:1258096.
- 15 Farboud B, and Meyer BJ. 2015. Dramatic Enhancement of Genome Editing by
- 16 CRISPR/Cas9 Through Improved Guide RNA Design. *Genetics*
- 17 10.1534/genetics.115.175166.
- 18 Hsu PD, Lander ES, and Zhang F. 2014. Development and applications of CRISPR-Cas9
- 19 for genome engineering. Cell 157:1262-1278.
- 20 Jiang F, and Doudna JA. 2015. The structural biology of CRISPR-Cas systems. Current
- 21 Opinion in Structural Biology 30:100-111.
- 22 Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, Morgan MT,
- 23 and Carey VJ. 2013. Software for Computing and Annotating Genomic Ranges. PLoS
- 24 Comput Biol 9:e1003118.
- 25 Makarova K, Grishin N, Shabalina S, Wolf Y, and Koonin E. 2006. A putative RNA-
- 26 interference-based immune system in prokaryotes: computational analysis of the
- 27 predicted enzymatic machinery, functional analogies with eukaryotic RNAi, and
- 28 hypothetical mechanisms of action. *Biology Direct* 1:7.

- 1 Mali P, Esvelt KM, and Church GM. 2013. Cas9 as a versatile tool for engineering
- 2 biology. *Nat Meth* 10:957-963.
- 3 Mojica FJM, Díez-Villaseñor C, García-Martínez J, and Almendros C. 2009. Short motif
- 4 sequences determine the targets of the prokaryotic CRISPR defence system. *Microbiology*
- 5 155:733-740.
- 6 R Core Team. 2014. R: A Language and Environment for Statistical Computing. 3.1.2 ed:
- 7 R Foundation for Statistical Computing.
- 8 Ran FA, Hsu Patrick D, Lin C-Y, Gootenberg Jonathan S, Konermann S, Trevino AE,
- 9 Scott David A, Inoue A, Matoba S, Zhang Y, and Zhang F. 2013a. Double Nicking by
- 10 RNA-Guided CRISPR Cas9 for Enhanced Genome Editing Specificity. Cell 154:1380-
- 11 1389.
- 12 Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, and Zhang F. 2013b. Genome
- engineering using the CRISPR-Cas9 system. *Nat Protocols* 8:2281-2308.
- Roberson E. 2015. Survey of 3'GG gRNA sites in 6 genomes
- 15 <a href="http://dx.doi.org/10.6084/m9.figshare.1371077">http://dx.doi.org/10.6084/m9.figshare.1371077</a>.
- 16 Xie Y. 2013. *Dynamic Documents with R and knitr*: Chapman and Hall/CRC.