



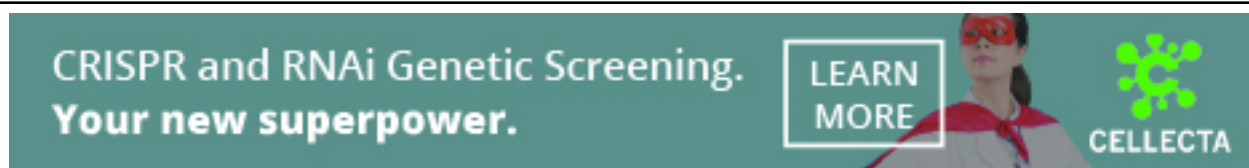
CRISPR-Cas9-based repeat depletion for the high-throughput genotyping of complex plant genomes

Marzia Rossato, Luca Marcolungo, Luca De Antoni, et al.

Genome Res. published online May 1, 2023

Access the most recent version at doi:[10.1101/gr.277628.122](https://doi.org/10.1101/gr.277628.122)

P<P	Published online May 1, 2023 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **CRISPR-Cas9-based repeat depletion for the high-throughput genotyping of**
2 **complex plant genomes**

3 Marzia Rossato^{1,2*}, Luca Marcolungo^{1*}, Luca De Antoni¹, Giulia Lopatriello¹, Elisa Bellucci³, Gaia
4 Cortinovis³, Giulia Frascarelli³, Laura Nanni³, Elena Bitocchi³, Valerio Di Vittori³, Leonardo Vincenzi¹, Filippo
5 Lucchini¹, Kirstin E. Bett⁴, Larissa Ramsay⁴, David James Konkin⁵, Massimo Delledonne^{1,2*} and Roberto
6 Papa^{3*}

7 *equal contribution

8 [§]corresponding authors

9 Marzia.rossato@univr.it; Luca.marcolungo@univr.it; Luca.deantoni@univr.it; giulia.lopatriello@univr.it;
10 e.bellucci@staff.univpm.it; gaia.cortinovis93@gmail.com; g.frascarelli@pm.univpm.it;
11 l.nanni@staff.univpm.it; e.bitocchi@staff.univpm.it; v.divittori@staff.univpm.it; leonardo.vincenzi@univr.it;
12 filippo.lucchini@univr.it; k.bett@usask.ca; l.ramsay@usask.ca; David.Konkin@nrc-cnrc.gc.ca;
13 massimo.delledonne@univr.it; r.papa@staff.univpm.it

14 ¹Department of Biotechnology, University of Verona, Strada Le Grazie 15, 37134, Verona, Italy

15 ²Genartis s.r.l., Via IV Novembre 24, 37126, Verona, Italy

16 ³Department of Agricultural, Food and Environmental Sciences, Polytechnic University of Marche, via Brecce
17 Bianche, 60131, Ancona, Italy

18 ⁴Department of Plant Sciences, University of Saskatchewan, 51 Campus Drive, Saskatoon, Saskatchewan
19 S7N 5A8, Canada

20 ⁵National Research Council Canada, 110 Gymnasium Place, Saskatoon, Ontario S7N 0W9

21

22 **Running Title:** Reduced genome representation by repeats exclusion

23

24

25

26

27 **ABSTRACT**

28 High-throughput genotyping enables the large-scale analysis of genetic diversity in population genomics and
29 genome-wide association studies that combine the genotypic and phenotypic characterization of large
30 collections of accessions. Sequencing-based approaches for genotyping are progressively replacing
31 traditional genotyping methods due to the lower ascertainment bias. However, genome-wide genotyping
32 based on sequencing becomes expensive in species with large genomes and a high proportion of repetitive
33 DNA. Here we describe the use of CRISPR-Cas9 technology to deplete repetitive elements in the 3.76-Gb
34 genome of lentil (*Lens culinaris*), 84% consisting of repeats, thus concentrating the sequencing data on
35 coding and regulatory regions (single-copy regions). We designed a custom set of 566,766 gRNAs targeting
36 2.9 Gbp of repeats and excluding repetitive regions overlapping annotated genes and putative regulatory
37 elements based on ATAC-seq data. The novel depletion method removed ~40% of reads mapping to
38 repeats, increasing those mapping to single-copy regions by ~2.6-fold. When analyzing 25 million fragments,
39 this repeat-to-single-copy shift in the sequencing data increased the number of genotyped bases of ~10-
40 fold compared to non-depleted libraries. In the same condition, we were also able to identify ~12-fold more
41 genetic variants in the single-copy regions and increased the genotyping accuracy by rescuing thousands of
42 heterozygous variants that otherwise would be missed due to low coverage. The method performed similarly
43 regardless of the multiplexing level, type of library or genotypes, including different cultivars and a closely-
44 related species (*L. orientalis*). Our results demonstrated that CRISPR-Cas9-driven repeat depletion focuses
45 sequencing data on meaningful genomic regions, thus improving high-density and genome-wide genotyping
46 in large and repetitive genomes.

47

48 **KEYWORDS:** repetitive elements, CRISPR-Cas9, sequencing-based genotyping, high-throughput
49 sequencing library

50

51

52

53

54

55

56 INTRODUCTION

57 The efficient and accurate determination of genotypes is necessary for large-scale projects investigating the
58 genetic composition of germplasm collections representing wild and domesticated species and inbred lines.
59 One example is the EU H2020 project INCREASE (www.pulsesincrease.eu) (Bellucci et al. 2021), which
60 focuses on four legume staples: chickpea, common bean, lentil and lupin. Such projects depend on large
61 cohorts of individuals to enable the comparative analysis of samples with sufficient statistical power. Cost-
62 effective high-throughput genotyping methods are therefore needed to increase the number of samples that
63 can be processed in an economically feasible manner (Bellucci et al. 2021). This can only be achieved by
64 reducing the fraction of each individual genome that is sequenced while ensuring that the same homologous
65 regions are examined in each individual (Peterson et al. 2012).

66 High-throughput low-cost genotyping has largely been achieved by the analysis of single-nucleotide
67 polymorphisms on microarray-based platforms (SNP arrays). These allow up to several thousand SNPs to
68 be tested simultaneously (Pavan et al. 2020). This approach considers a predefined set of markers, resulting
69 in fixed costs per individual regardless of the genome size and fraction of repetitive DNA. However, analysis
70 is restricted to known SNPs that are frequent in the population, while rare and unknown SNPs are ignored.
71 This is a drawback when analyzing diverse landraces and distant wild relatives, as required in the
72 germplasm characterization projects mentioned above (Lachance and Tishkoff 2013).

73 More recently, next generation sequencing (NGS) has provided an opportunity to discover genome-wide
74 variants in a less biased manner. Sequencing-based approaches for genotyping involves low coverage (5–
75 10×) whole-genome sequencing (lcWGS), allowing the characterization of several million variants (Tanaka et
76 al. 2021; Friel et al. 2021). To reduce costs enough to make WGS affordable even in large germplasm
77 collections, very low coverage (0.5–2×) WGS (ultra-lcWGS) can be combined with imputation to infer
78 positions that are not sequenced or genotyped (Deng et al. 2022; Zan et al. 2019; Wang et al. 2016).
79 Alternatively, sequencing costs are often minimized by reduced-representation sequencing, which comprises
80 methods such as genotyping by sequencing (GBS) (Elshire et al. 2011), restriction site-associated DNA

81 sequencing (RAD-Seq) (Davey et al. 2011; Baird et al. 2008) and double-digest RAD-Seq (ddRADseq)
82 (Truong et al. 2012; Peterson et al. 2012). These methods concentrate sequencing data on regions adjacent
83 to restriction sites by exploiting the specificity of restriction endonucleases. Reduced-representation
84 sequencing is suitable for large cohorts, but provides only low-resolution data, with a small fraction of
85 analysed and genotyped bases (Pavan et al. 2020) that may not provide sufficient marker density and depth
86 to confidently identify variants under selection in large genomes (Guerra-García et al. 2021). The resolution
87 can be increased without significantly greater costs in sample prep by using the Twist 96-Plex Library Prep
88 Kit (formerly iGenomX Riptide Kit) to generate multiplexed libraries, allowing 96–960 samples to be
89 processed simultaneously and resulting in the non-random sampling of millions of genomic positions
90 (Siddique et al. 2019).

91 Despite the advantages of sequencing-based genotyping over SNP arrays, one common disadvantage is
92 that sequencing methods generally do not distinguish between repetitive (low-complexity) and single-copy
93 (high-complexity) regions, the latter comprising coding and regulatory regions that are the main targets of
94 natural selection and thus the focus of most genotyping projects. In contrast, low-complexity regions of plant
95 genomes mainly comprise transposable elements, simple sequence repeats and tandem repeats.
96 Transposable elements play a key role in genome evolution, but the analysis of such regions is technically
97 challenging and largely uninformative in genotyping studies, unless dedicated analysis workflows are applied
98 (Yan et al. 2022). Mapping reads to transposable/repetitive elements can result in low-quality alignments that
99 hinder the calling of accurate genotypes, which is a consistent challenge particularly for those plant species
100 with large genomes, where repetitive elements account for up 90% of the total DNA. This includes many
101 domesticated crops such as corn (*Zea mays*), wheat (*Triticum* spp.), lentil (*Lens culinaris*) and onion (*Allium*
102 *cepa*) (Feuillet et al. 2011). One strategy to address this issue is whole exome sequencing (WES), which
103 selects coding regions for preferential sequencing (Hodges et al. 2007) as shown in lentil, wheat and barley
104 (Ogutcen et al. 2018; He et al. 2019). However, WES only focuses on coding sequences and thus overlooks
105 regulatory elements, which are equally important as sources of genetic diversity (Ricci et al. 2019; Wang et
106 al. 2019; Tian et al. 2020).

107 Ideally, lcWGS could be focused on the most complex parts of the genome, avoiding wasted effort on the
108 sequencing of repetitive elements. This could be achieved by using enzymes that enable target enrichment
109 by depleting unwanted sequences from NGS libraries. For example, the duplex-specific nuclease (DSN)
110 selectively digests double-stranded DNA molecules, and can be used to eliminate highly abundant

111 sequences in a controlled denaturation-reassociation reaction (Zhulidov et al. 2004). This method has been
112 used in RNA-seq analysis to remove abundant transcripts (Zhao et al. 2014; Miller et al. 2013) and, just
113 occasionally, also to delete repetitive elements in DNA-seq libraries generated from plant genomes (Ichida
114 and Abe 2019; Matvienko et al. 2013). However, (Matvienko et al. 2013)DSN can also remove informative
115 repetitive elements, such as the coding sequences of abundant gene families, which are particularly relevant
116 in polyploid plants arising from whole genome duplication events (Matvienko et al. 2013). More recently, the
117 CRISPR-Cas9 (clustered regularly interspaced short palindromic repeats and CRISPR-associated nuclease
118 9) system has been used for the selective depletion of unwanted genome fractions from sequencing libraries
119 (Gu et al. 2016). The Cas9 enzyme can be programmed to cut library fragments by designing specific guide-
120 RNA sequences targeting the unwanted sequences. Subsequently, only intact fragments -retaining adapters
121 at both ends- can be effectively amplified by PCR and generate productive clusters on a sequencing flow-
122 cell. The DASH approach (depletion of abundant sequences by hybridization) involved the use of Cas9 to
123 exclude ribosomal RNA (rRNA) sequences from RNA-seq libraries and to remove DNA from common
124 pathogens in order to detect rare pathogens in metagenomic samples (Gu et al. 2016). A similar technology
125 has been recently commercialized under the name “CRISPRclean” by JumpCode Genomics (JumpCode
126 Genomics 2021).

127 Here we determined whether CRISPRclean technology could be used to deplete the repetitive elements in
128 libraries prepared from the 3.76-Gbp genome of lentil (*L. culinaris*), 84% of which is repetitive DNA.
129 CRISPRclean technology was combined with Twist multiplexing libraries and we evaluated its performance,
130 focusing on the technical features required for genotyping. Our results will facilitate the large-scale genomic
131 analysis of lentil as well as other plant species with large and highly-repetitive genomes.

132

133 **RESULTS**

134 **Depletion of *L. culinaris* repetitive DNA using CRISPR-Cas9**

135 We designed a custom set of gRNAs to deplete the repetitive DNA content of the *L. culinaris* CDC Redberry
136 genome (Ramsay et al. 2021), targeting transposable elements (totaling 3.1 Gbp of sequence across the
137 whole genome, corresponding to 82.5% of its size), simple sequence repeats (58 Mbp, 1.5%) and tandem
138 repeats (13 Mbp, 0.4%) in the nuclear genome, as well as the entire mitochondrial genome (mtDNA, 489
139 kbp) and chloroplast genome (cpDNA, 118 kbp) (**Supplemental Table S1**). We excluded repetitive DNA that

140 overlapped with functional regions such as annotated genes (185 Mbp, 5%) and putative regulatory regions
141 identified using ATAC-seq data (78 Mbp, 2%) (**Supplemental Table S1**). All nuclear DNA outside the gRNA
142 target regions is hereafter defined as single-copy . The final design comprised 566,766 gRNAs with at least
143 25 recognition sites, potentially targeting 2.9 Gbp (77%) of the *L. culinaris* nuclear genome and 93.5% of its
144 repetitive regions when using a sequencing library with 500-bp inserts (**Supplemental Table S2** and
145 **Supplemental File S1**). An additional 2,366 gRNAs targeted the mitochondrial and chloroplast genomes
146 (**Supplemental Table S2**). The gRNAs were assigned to 11 pools based on their cutting frequency in the *L.*
147 *culinaris* genome (**Supplemental Table S2**).

148 The custom gRNAs were tested on three Twist 8-plex libraries, allowing the reproducible sampling of the
149 same genomic regions by random priming during first-strand DNA synthesis. Each 8-plex library comprised
150 replicates of three distinct *L. culinaris* samples (cv. Castelluccio) (**Supplemental Table S3**). Cas9/gRNAs
151 ribonucleoprotein (RNP) complexes (1:2.5 protein/gRNA ratio) were generated, and gRNAs with more target
152 sites in the genome were used at higher relative concentrations in the final reaction (**Supplemental Table**
153 **S2**). Depletion reactions in the presence of RNP complexes were conducted either using all gRNAs
154 simultaneously or by splitting the gRNA pools into three groups based on cutting frequency (**Supplemental**
155 **Table S2**) and using the groups sequentially, starting with the lowest cutting frequency. Depleted and non-
156 depleted libraries were sequenced, generating 91 million fragments on average (**Supplemental Table S4**).
157 The sequencing data were normalized at ~50 million fragments per library in order to compare the proportion
158 of reads mapping on repetitive and single-copy regions of the nuclear genome and on the organelle
159 genomes (**Figure 1**). The number of reads mapping to repetitive regions (total repeats, nuclear repeats,
160 mtDNA and cpDNA) was significantly lower in the depleted libraries compared to the non-depleted libraries,
161 with the sequential depletion strategy using three gRNA groups performing best and depleting 37.7% of the
162 repetitive DNA (**Figure 1A**). The results were similar when considering only the nuclear repetitive regions
163 (37.2% depletion) (**Figure 1B**). A small fraction of total reads mapped to the organelle genomes (~1%). Both
164 depletion strategies were similarly effective in the chloroplast genome, resulting in ~78.5% depletion (**Figure**
165 **1C**). In contrast, no significant depletion was observed in the mitochondrial genome (**Figure 1D**). In parallel,
166 the sequential depletion strategy achieved a 130% increase in the number of reads mapping to single-copy
167 regions, from 17.5 to 40.4 million (**Figure 1E**). Given that the concentration of Cas9 RNPs influences the
168 cutting efficiency (Gu et al. 2016), we repeated the sequential depletion strategy using double amount of
169 Cas9 and gRNAs. This modified the read distribution further, achieving 41.2% depletion of nuclear repeats

170 and a 160% increase in reads mapped to single-copy regions. The sequential depletion strategy with double
171 RNPs was therefore the most efficient, and was applied in all subsequent experiments. Overall, our results
172 demonstrated that the custom gRNA set and Cas9 effectively targeted fragments containing repetitive DNA
173 sequences and depleted them in the resulting sequencing libraries. **Figure 2A-B** shows the alignments of
174 reads at two representative genomic regions, confirming that less sequencing data was assigned to regions
175 of repetitive DNA and more reads were mapped to single-copy parts of the genome.

176 **Efficiency of CRISPR-Cas9-mediated depletion for different classes of nuclear repeats**

177 We next examined the depletion of different classes of repetitive sequences in the *L. culinaris* genome.
178 Reads mapping to the most abundant retroelements, namely the Ty3-Gypsy family (64% of the genome
179 (Ramsay et al. 2021), were reduced by 47% in the depleted libraries, whereas those mapping to the Ty3-
180 Copia family (15% of the genome) and other long terminal repeat (LTR) elements (3% of the genome) were
181 depleted by 3% and 38%, respectively (**Figure 3A** and **Supplemental Table S5**). In contrast, there was no
182 decrease in the abundance of other transposable elements (LINE, CACTA, Mu, hAT, Helitron, Harbinger,
183 mariner and Sine), each representing < 1% of the genome (**Supplemental Table S5**), and there was no
184 reduction in the number of reads mapping to tandem repeats (0.4% of the genome) or simple sequence
185 repeats (8% of the genome) (**Figure 3A**). We observed a significant correlation between the variation in
186 mapped reads after depletion and the abundance of these repeat classes in terms of overall repeat length
187 and occurrence in the genome (**Figure 3B-C** and **Supplemental Table S5**). Given that the number of
188 gRNAs targeting each repeat class increased proportionally with the repeat size and occurrence, the most
189 efficiently depleted repetitive elements also featured a higher density of gRNA targets (**Figure 3D**). There
190 was a significant correlation between the variation of mapping reads following depletion and the gRNA
191 density over the whole target region when considering each single cut site in the genome (**Supplemental**
192 **Figure S1**). In particular, target regions with a density > 8 gRNAs/kbp showed a read reduction in 85% of
193 cases (**Supplemental Figure S1**) whereas regions targeted by < 8 gRNAs/kbp usually showed limited or no
194 depletion (**Supplemental Figure S1**). We therefore concluded that the depletion efficiency across different
195 repeat classes was dependent on the density of gRNA targets.

196 **Impact of CRISPR-Cas9-mediated repeat depletion on genotyping accuracy**

197 Next, we investigated the impact of CRISPR-Cas9-mediated repeat depletion on the number of genomic
198 positions in the single-copy regions where a base can be reliably genotyped (PASS at a depth of ≥ 5 reads).

199 For this analysis, sequencing data generated from depleted and non-depleted *L. culinaris* cv. Castelluccio
200 samples were downsampled from 62 to 6 million fragments to mimic lcWGS and ultra-lcWGS. Consistently
201 more bases were genotyped within the single-copy regions of the depleted samples over the whole range
202 considered (from ~3.5 to ~13-fold), with the highest gains at the lowest amounts of sequencing data (6 to 25
203 million fragments) (**Figure 4A**). Because a genotyped position does not necessarily allow the variant to be
204 identified (this also depends on allele coverage), we also determined the impact of repeat depletion on
205 variant calling. Following depletion, the total number of variants identified in the single-copy regions
206 increased significantly from ~4.5 to ~18-fold, with a delta of ~25,000 to ~1 million more variants identified in
207 the depleted sample (**Figure 4B**). Also the number of heterozygous variants identified was increased,
208 although these constituted a minor fraction of total variants, as lentil is an autogamous species (**Figure 4C**).
209 This allowed us to identify from ~650 up to ~50,000 variant positions that would be erroneously classified as
210 reference without the depletion (**Figure 4D**). These false negative variants in the non-depleted sample were
211 not called due to allelic imbalance caused by the low coverage (**Supplemental Figure S2**). Consistently,
212 most of these false negative variants rescued by depletion were heterozygous (~97%, **Figure 4D**).

213 **Performance of CRISPR-Cas9-mediated repeat depletion on different samples and library types**

214 Finally, we assessed the performance of CRISPR-Cas9-mediated repeat depletion on different lentil
215 genotypes, multiplexing levels and library types (**Supplemental Table S3** and **Supplemental Table S4**).
216 Similar variations in the coverage of repetitive/single-copy regions and the number of mapped reads were
217 observed when depleting multiplex libraries generated from a different cultivar (RB, Redberry) or from the
218 closely-related species *L. orientalis* (**Figure 5A-B**) when compared to the original Castelluccio cultivar
219 (**Figure 1**). To assess the impact of depletion when comparing different samples, these data were
220 downsampled as described above, genotyped positions were identified and intersected with those of
221 Castelluccio samples. Depletion improved the genotyping reproducibility, as the number of genotyped
222 positions in common between all analyzed samples, within the single-copy regions, was consistently higher
223 than in the condition without depletion (**Figure 5C**). Finally, there was no significant difference in the
224 performance of CRISPR-Cas9-mediated repeat depletion when library multiplexing was increased from 8-
225 plex to 96-plex while maintaining the 1:2.5 Cas9:gRNA ratio and 1 ng of treated library per sample (**Figure**
226 **6A-B**), or when treating standard singleplex WGS libraries (**Figure 6C-D**). Overall, these results
227 demonstrated that CRISPR-Cas9-mediated repeat depletion using the same gRNA set is at least equally
228 effective when applied to a group of two lentil cultivars and one close wild relative, providing an improved

229 genotyping reproducibility and the possibility to process individual samples or multiple samples
230 simultaneously.

231

232 **DISCUSSION**

233 In a typical genotyping experiment based on sequencing, the data derived from repetitive DNA is directly
234 proportional to the repeat content of the genome, however these data are largely uninformative. The bigger
235 the genome, the more sequencing costs are therefore wasted on repeats. The traditional solution is the
236 capture and sequencing of coding regions (WES). However, we approached the problem from the opposite
237 perspective by depleting repetitive elements from the large genome of *L. culinaris* using CRISPR-Cas9
238 technology. Similar methods have been used for RNA-seq library normalization (Prezza et al. 2020),
239 metatranscriptomics (Gu et al. 2016), pathogen detection (Gu et al. 2016; Le et al. 2021) and single-cell
240 analysis (Homburger et al. 2023; Le et al. 2021). Here, the main challenge is that 84% of the 3.7-Gb *L.*
241 *culinaris* genome is repetitive DNA. The design of the gRNA array, therefore, required stringent multi-step
242 filtering resulting in a set of ~566,766 gRNAs. To our knowledge, this is the first time CRISPR-Cas9 has
243 been used with such a large number of gRNAs either *in vitro* or *in vivo*, representing a fundamental advance
244 in the technology platform beyond the specific goals of our project. Despite these technical challenges,
245 CRISPR-Cas9-mediated repeat depletion produced sequencing libraries with consistently lower proportions
246 of repetitive DNA (41.2% depletion) and enriched the single-copy regions (160% increase), thus allowing the
247 generation of more meaningful sequencing data. Equivalent results have been achieved not only in *L.*
248 *culinaris* cv. Redberry (source of the reference genome for gRNA design) but also in another *L. culinaris*
249 cultivar (Castelluccio) and in the close-related species *L. orientalis*, using the same gRNA set. The approach
250 demonstrated to be useful also for other species beyond lentil, namely in bread wheat (*Triticum aestivum*),
251 where a dedicated repeat-specific gRNA set showed similar depletion performances, both in Chinese Spring
252 and Jagger cultivars (Jumpcode Genomics personal communication).

253 Repetitive DNA in plant genomes can be divided into two broad categories: dispersed mobile elements and
254 tandem/simple repeats. Dispersed mobile elements are made up of DNA transposons and retrotransposons,
255 the most abundant of which are the LTR retrotransposons (Bennetzen and Wang 2014). Although mobile
256 elements are not under the same selection pressure as genes, the degree of conservation across multiple
257 copies of the same element is sufficiently high to allow the targeting of multiple copies with single gRNAs.

258 The most efficient depletion (47%) was achieved for the most abundant LTR retrotransposon family (Gypsy,
259 ~64% of the genome), followed by all the other LTR families (~15%). The cutting of LTR elements by Cas9
260 was responsible for almost the entire depletion observed at the genome-wide level, whereas the depletion of
261 other mobile elements was negligible. Given that some repetitive elements were not efficiently depleted, it is
262 likely that some gRNAs in the current lentil design are not yet optimal. Another factor influencing the
263 depletion performances was probably the lower repetition of such elements, which translated into a poor
264 cutting frequency in the final gRNA design, comprising only gRNAs with 25 targets. These targets usually
265 featured < 8 gRNA/kbp, namely a density associated with a low depletion rate. A similar observation was
266 reported for RNA-seq libraries, where a gRNA every 50–100 nucleotides (10–20 gRNAs/kbp) achieves
267 excellent depletion results (Gu et al. 2016). Simple and tandem repeats were also not depleted efficiently,
268 although the gRNA density was close to 8. In these cases, the highly repetitive motifs of such sequences
269 may have reduced the cutting efficiency of Cas9 (Müller Paul et al. 2022). Given that LTR retrotransposons
270 make up the majority of repeats in plant genomes (94% in *L. culinaris*) and are the principal cause of plant
271 genome size variation (Bennetzen and Wang 2014; Lee and Kim 2014), the design of gRNAs to target only
272 LTR sequences may be the most efficient strategy to reduce the genome size in sequencing experiments.
273 Future gRNA designs in other species and further optimization of the *L. culinaris* design should maximize the
274 gRNA number on these most abundant elements, instead of dispersing the effort across the remaining
275 repetitive fraction (< 10%). Another factor influencing the efficiency of CRISPR-Cas9-mediated repeat
276 depletion was the dose of Cas9 and gRNAs; doubling their dose indeed improved repeat depletion by ~10%,
277 albeit with a slight increase in overall costs. To maximize depletion, the concentration of most efficient
278 gRNAs could be also increased by excluding gRNA-Cas9 complexes with low cut performances and
279 preoccupying limited Cas9 molecules, while maintaining the final amount of gRNA and Cas9 enzyme.
280 Finally, also the order of gRNA addition was important, as higher depletion efficiency was achieved when
281 splitting gRNAs into three groups based on cutting frequency and using the groups sequentially. A possible
282 explanation of this phenomenon is that most abundant gRNAs, with the highest number of target sites in the
283 genome, could interfere with the “genome patrolling” of other RNPs targeting less abundant elements.
284 Therefore, the gRNA target density, RNP concentration and prioritization of gRNAs with less abundant
285 targets are factors that can improve depletion efficiency.

286 Organelle genomes often constitute a large fraction of DNA derived from plants (Sakamoto and Takami
287 2018). Although the mitochondrial and chloroplast genomes are smaller than the nuclear genome, they are

288 present in multiple copies per cell, and they can represent > 20% of the total sequence data (Gargiulo et al.
289 2021; Ren et al. 2021). CRISPR-Cas9-mediated repeat depletion has been shown to reduce the fraction of
290 sequencing libraries derived from organelle genomes in ATAC-seq experiments (Montefiori et al. 2017). Our
291 method was efficient for the depletion of chloroplast DNA (by 67%) while the depletion of mitochondrial DNA
292 was only marginal, possibly reflecting the different abundance of the two organelle genomes in the starting
293 genomic sample. Still, in lentil, the fraction of sequencing data attributable to organelles was largely due to
294 chloroplasts (88%), whose depletion was therefore sufficient to decrease the data mapping on organelles
295 after Cas9 treatment (-65% overall). Although in the case of lentil the total sequencing data attributable to
296 organelles was rather low (~1.3% in the not depleted libraries), the depletion of organelle DNA from
297 sequencing libraries will be highly beneficial for organisms with a strongly unbalanced ratio of organelle vs
298 nuclear DNA, such as *Cypridium calceolus* (Gargiulo et al. 2021) and *Haematococcus pluvialis* (Ren et al.
299 2021). Although this has not been investigated in the present work, it is plausible that gRNAs designed for
300 organelle's genomes could also target nuclear integrants of plastid/mitochondrial DNA (NUPTs and NUMTs),
301 that are homologous to the cpDNA and mtDNA (Zhang et al. 2020; Sloan et al. 2018). Although the efficiency
302 of CRISPR-Cas9-mediated repeat depletion could be improved, the current set of gRNAs allowed us to
303 genotype consistently more bases on single-copy as compared to not depleted libraries, and consequently to
304 identify more genetic variants. We observed that coupling the depletion with 25 million sequencing fragments
305 provided the best balance between costs and fraction of genotyped bases. In this condition, depleted
306 samples reached an average coverage of ~5 fold on single-copy regions, with gains of 10- and 12-fold in the
307 number of genotyped positions and identified variants, respectively, as compared to not depleted libraries.
308 We estimated that one would need 3.5-4 times more sequencing data to achieve the same performances
309 when using not depleted libraries. By targeting the majority of genomic length (84%), the depletion allowed to
310 pour a large fraction of sequencing data on the single-copy regions, that are instead just a small fraction
311 (16%). For this reason, the CRISPR-Cas9-repeat depletion was more effective to improve genotyping
312 performances than increasing the overall sequencing coverage, being also beneficial for the identification of
313 heterozygous variants that would otherwise be missed due to unbalanced allelic sampling. Although this
314 involved only ~3% of total variants identified in lentil, as this is an autogamous diploid species, allelic
315 imbalance is a well-known cause of errors in genotyping experiments based on sequencing (Cooke et al.
316 2016). CRISPR-Cas9-mediated repeat depletion can therefore improve the accuracy of genotyping
317 experiments, especially in plants with highly heterozygous genomes and/or in polyploids.

318 Given that CRISPR-Cas9 repeat depletion allows to concentrate the sequencing data on the desired regions,
319 the amount of positions that were genotyped in common between multiple samples was consistently higher
320 in the depleted ones. In population studies, the shift in the distribution of sequencing data may contribute to
321 reduce the number of missing data, thereby detecting a larger number of differences between samples.
322 Furthermore, this approach was successful in different cultivars of *L. culinaris*, and also in the closely-related
323 species *L. orientalis*. This is important because genotyping experiments typically include distant/wild relatives
324 and related species, from which it is possible to develop evolutionary studies and plan breeding experiments,
325 including the introgression of characters of interest. For example, the INCREASE project features a
326 collection of 2000 lentil accessions that includes both cultivated varieties and local landraces (Guerra-García
327 et al. 2021). Further experiments could determine whether the gRNAs designed in this study are also
328 suitable for the depletion of repeats in other closely related leguminous species (Fabaceae) with very large
329 and repetitive genomes, such as pea (*P. sativum*, 3.92 Gb, 83% repetitive)(Kreplak et al. 2019) and faba
330 bean (*Vicia faba* L., 12 Gb, 79% transposon-derived repeats)(Jayakodi et al. 2023).

331 The cost of NGS library preparation for a genotyping project can easily exceed the cost of sequencing in the
332 case of small genomes and/or ultra-lcWGS, especially given the steadily falling price of sequencing. More
333 recent library-preparation kits circumvent several lengthy steps that require expensive reagents, and allow
334 large sample sets to be processed in multiplex reactions. We used the Twist 96-Plex Library Prep Kit
335 (formerly iGenomX Riptide) that constructs Illumina NGS libraries by polymerase-mediated extension of
336 barcoded random primers. This type of library is beneficial for genotyping in general because random
337 priming reduces uniform genome coverage but allows more reproducible sampling of the same sites across
338 multiple samples (Siddique et al. 2019). Most importantly, the kit is designed to process large numbers of
339 samples (up to 96 simultaneously) at low costs and without advanced equipment (just a multichannel
340 pipette). To design, filter and synthesize the gRNA set targeting *L. culinaris* repeats required 10 weeks and
341 cost ~20,000 USD. The latter comprised gRNAs and depletion reagents sufficient for 30 reactions, each for a
342 maximum of 96 samples treated in multiplex, corresponding to approximately 10 USD per sample. We
343 estimated that the net cost to achieve ~5-fold average coverage of the single-copy regions of the lentil
344 genome by combining CRISPR-Cas9-mediated repeat depletion with a 96-plex Twist library is approximately
345 US\$75, made up of US\$15 for library preparation, US\$10 for depletion and US\$40 for sequencing on a
346 NovaSeq 6000 S4 flowcell, generating 25 million fragments per sample. As such, our results demonstrated
347 that depleted libraries are more informative than standard ones when normalized for the amount of

348 sequencing data. Alternatively, CRISPR-Cas9-mediated repeat depletion can be used to reduce sequencing
349 costs (by ~75% in lentil), because the same number of genotyped bases (or detected variants) in single-copy
350 regions can be detected with much less sequencing data. CRISPR-Cas9-mediated repeat depletion
351 combined with Twist multiplex libraries is therefore an effective strategy for genotyping projects involving
352 hundreds or thousands of samples. Dealing with less-repetitive datasets can also reduce the complexity of
353 the genotyping analysis and the computational resources required.

354 The method therefore has the potential to increase our genetic knowledge of plant species that are currently
355 difficult to analyze without a significant economic investment due to the large genome size and high
356 proportion of repetitive DNA. Population studies, eQTL analysis, GWAS and pre-breeding programs are just
357 some of the approaches that can benefit from CRISPR-Cas9-mediated repeat depletion.

358

359 **METHODS**

360 **Multiplex-library preparation and sequencing.** We prepared 8-plex and 96-plex multiplex libraries
361 according to the Twist 96-Plex Library Preparation Kit protocol (Twist Bioscience, South San Francisco, CA,
362 USA) with the following modifications. For each sample, we denatured 100 ng of genomic DNA (25 ng/ μ l) at
363 98 °C for 1 min. Ultra-low (30%) GC random primer set A was used for the extension and termination
364 reaction (Reaction A) followed by 8 and 9 cycles of PCR amplification for the 96-plex and 8-plex libraries,
365 respectively. Final libraries were purified using Twist DNA Purification Beads (0.65 \times volume) and a second
366 round of purification was applied to the supernatant using 10 μ l of beads to achieve a median insert size of
367 500 bp. Libraries were quantified using the Qubit BR DNA kit and a Qubit device (Thermo Fisher Scientific,
368 Waltham, MA, USA) and size distributions were assessed using a Tape Station System (Agilent
369 Technologies, Santa Clara, CA, USA). Non-depleted libraries were pooled at equimolar concentrations and
370 sequenced on a NovaSeq 6000 instrument (Illumina, San Diego, CA, USA) to generate 150-bp paired-end
371 reads.

372 **WGS library preparation and sequencing.** Genomic DNA samples were fragmented using a Covaris
373 sonicator to achieve an average size of 400 bp, and Illumina PCR-free libraries were prepared from 700 ng
374 DNA using the KAPA Hyper prep kit and unique dual-indexed adapters (5 μ L of a 15 μ M stock) according to
375 the supplier's protocol (Roche, Basel, Switzerland). The library concentration and size distribution were
376 assessed on a Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). Non-depleted WGS libraries were

377 pooled at equimolar concentrations and sequenced on a NovaSeq 6000 instrument (Illumina, San Diego,
378 CA, USA) to generate 150-bp paired-end reads.

379 **Design of gRNAs.** The gRNA set was designed by JumpCode Genomics (San Diego, CA, USA) against the
380 repetitive regions of the *L. culinaris* CDC Redberry v2.0 reference genome (Ramsay et al. 2021)
381 (<https://knowpulse.usask.ca/genome-assembly/Lcu.2RBY>). The available repeat annotation (transposable
382 elements and tandem repeats) was integrated with the annotation of simple and tandem repeats identified by
383 RepeatMasker v4.0.6 and Tandem Repeat Finder v4.9 using 2 7 7 80 10 50 2000 -d -h parameters to
384 identify intervals for gRNA design (**Supplemental Table S1**). Adjacent or overlapping intervals were
385 collapsed into single intervals before design. As a first step, all 20 nt sequences with adjacent PAM sites for
386 Cas9 (NGG) were identified in the target intervals. Second, the guides were filtered to exclude secondary
387 structure, high and low GC content, homopolymers, dinucleotide repeats and low *in vitro* cleavage efficiency
388 prediction scores (Azimuth algorithm; (Doench et al. 2016)). Third, the resulting guides were filtered to
389 minimize off-target cleavage in single-copy regions of the genome by excluding guides that have
390 complementary sites in genomic regions corresponding to genes and open-chromatin regions identified by
391 ATAC-seq (PRJNA912311) (allowing for up to 3 mismatches). As a final step, and to reduce the number of
392 guides in the set, guides were selected to have no fewer than 25 cleavage sites each and to maintain an
393 inter-guide spacing of at least 500 bp. The final guide set, comprising 569,088 unique guides, was split into
394 11 pools for the purpose of synthesis. The number of copies of each guide varied and reflected the number
395 of on-target cleavage sites for each guide. DNA oligonucleotides containing the target-specific 20 nt gRNA
396 sequence and invariant single gRNA sequence were synthesized, after which pools of oligonucleotides were
397 amplified by PCR and converted to RNA by *in vitro* transcription. The products of transcription were treated
398 with DNase I and column purified to generate the final gRNA material. Pools 1–3, 5–8 and 10 contain only
399 gRNAs targeting the nuclear genome. Pools 9 and 11 contain both nuclear and chloroplast genome gRNAs,
400 and pool 4 is the only pool containing gRNAs that target the nuclear, chloroplast and mitochondrial genomes
401 (**Supplemental Table S2** and **Supplemental File S1**). The number of gRNAs targeting each repeat class
402 are reported in **Supplemental Table S5**.

403 **Repeat depletion with JumpCode CRISPRclean.** Repetitive regions were depleted using the Cas9 protein
404 and the custom gRNA set described above, according to the Jumpcode CRISPRclean Ribosomal RNA
405 Depletion from Human RNA-seq Libraries for Illumina Sequencing protocol (Jumpcode Genomics, San
406 Diego, CA, USA) with the following modifications. The input was 10 and 100 ng for the 8-plex and 96-plex

407 libraries, respectively. Depletion was carried out either using all gRNA simultaneously or by splitting the
408 gRNA pools into three groups based on cutting frequency, which were used sequentially in order of
409 increasing cutting frequency (**Supplemental Table S2**). The sequential depletion strategy was also
410 conducted using the double amounts of gRNAs and Cas9. The reaction volume was 20 μ l when using all
411 gRNAs simultaneously or 26 μ l for the sequential and double sequential protocols. The quantity of each
412 gRNA pool per reaction is shown in **Supplemental Table S2** and amounted to 620 ng in the simultaneous
413 and sequential depletion reactions or 1240 ng in the double sequential depletion reaction. The Cas9 enzyme
414 was diluted 1:5 in 1 \times Cas9 Buffer and 0.0029 μ l was used per ng gRNA. The reactions were incubated at 37
415 $^{\circ}$ C and libraries were treated in the presence of gRNAs for a total of 1 h (simultaneous depletion protocol) or
416 3 h (sequential and double sequential depletion protocols, with the sequential gRNA pools added at 1-h
417 intervals). The depleted samples were then size selected using 0.6 \times volume of AMPure XP Beads (Beckman
418 Coulter, Brea, CA, USA). Libraries were amplified with 10 and 6 PCR cycles for the 8-plex and 96-plex
419 libraries, respectively before final purification with 60 μ l (0.6 \times volume) AMPure XP Beads. The concentrations
420 of depleted libraries were measured using the Qubit system and size distributions were assessed on a Tape
421 Station System as described above. Depleted libraries were pooled at equimolar concentrations and
422 sequenced on a NovaSeq 6000 instrument to generate 150-bp paired-end reads.

423 **Data analysis and variant calling.** Raw read quality was assessed using FastQC
424 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and the multiplex libraries were demultiplexed
425 using fgbio v1.3.0 DemuxFastqs (<http://fulcrumgenomics.github.io/fgbio/>), assigning fragments by exploiting
426 the unique sample identifier included during first-strand synthesis. The raw reads were then quality filtered
427 and the Illumina sequencing adapters removed using scythe v0.991 (<https://github.com/vsbuffalo/scythe>) and
428 sickle v1.33 (<https://github.com/najoshi/sickle>), respectively. Filtered reads were aligned to the *L. culinaris*
429 v2.0 reference genome using BWA-MEM v2.2.1 (Md et al. 2019) and the resulting alignments were
430 converted to BAM files and sorted using SAMtools v1.13 (Danecek et al. 2021). PCR-derived duplicates
431 were removed using the GATK MarkDuplicates tool v4.1.7.0 (Van der Auwera and O'Connor 2020) and
432 overlapping portions of the paired-end reads were clipped using the fgbio v1.3.0 ClipBam tool
433 (<http://fulcrumgenomics.github.io/fgbio/>). The resulting BAM files were used to calculate coverage depth,
434 breadth and fraction of PASS bases (at $\geq 5\times$) using BEDTools v2.30.0 genomecov (Quinlan and Hall 2010)
435 and GATK v3.8 CallableLoci, respectively (Van der Auwera and O'Connor 2020). The number of reads
436 aligning to the reference genome and to different regions of interest was calculated using SAMtools v1.13

437 (Danecek et al. 2021) with option -c to discard reads with a 2308 sam flag in order to consider only the
 438 primary alignment, thus omitting repetitive counts of the same multimapping reads. When necessary,
 439 sequencing data were normalized to a pre-defined number of input fragments using seqtk sample v1.3
 440 (<https://github.com/lh3/seqtk>).

441 The variation of mapped reads in depleted vs not depleted libraries was calculated using the formula:

442 $Variation\ of\ mapped\ reads = \frac{Mapped\ reads\ (depleted) - Mapped\ reads\ (not\ depleted)}{Mapped\ reads\ (not\ depleted)}$ To achieve a normal
 443 distribution of variation, in Supplemental Figure S1 the variation of mapped read coverage between the
 444 depleted and not depleted libraries was calculated using the following formula:

$$Variation\ of\ mapped\ read\ coverage = \log_2 \frac{Mean\ coverage\ (depleted)}{Mean\ coverage\ (not\ depleted)}$$

445 Regions with zero coverage in either depleted or not depleted conditions (6,097 and 6,319, respectively,
 446 over a total of 237,136 repetitive segments, of which 4,662 in common) were excluded from the calculation,
 447 as they represented a negligible fraction of total (3.3%) and showed a minimal coverage also in the opposite
 448 condition (< 0.2 fold). The plot in Supplemental Figure S1 was generated using the ggplot package in R
 449 (Hadley Wickham 2016; R Core Team 2020).

450 Genomic variants were identified using GATK HaplotypeCaller v4.1.7.0 with the parameters "--min-base-
 451 quality-score 20 -ERC GVCF" (Van der Auwera and O'Connor 2020). Individual gVCF files were merged
 452 using GATK GenomicsDBImport v4.1.7.0 and the final VCF file was generated using GATK GenotypeGVCFs
 453 v4.1.7.0 (Van der Auwera and O'Connor 2020). Variant filtration was achieved using GATK hard filters
 454 (<https://gatk.broadinstitute.org/hc/en-us/articles/360037499012?id=3225>).

455

456 **DATA ACCESS.** The sequencing datasets generated and analysed in this study have been submitted to the
 457 NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number
 458 PRJNA915594 and PRJNA912311.

459 **COMPETING INTEREST STATEMENT.** Authors MR and MD are partners of Genartis srl. The remaining
 460 authors declare that the research was conducted in the absence of any commercial or financial relationships
 461 that could be construed as a potential conflict of interest.

462 **ACKNOWLEDGMENTS.** This research was supported by the European Union's Horizon 2020 research and
463 innovation program, through the project INCREASE (www.pulsesincrease.eu) (grant agreement No.
464 862862). The founding sponsors had no role in the design of the study; in the collection, analyses, or
465 interpretation of data; in the writing of the manuscript; and in the decision to publish the results. We
466 acknowledge Matteo De Biasi for the support in bioinformatic analysis, the Twist Bioscience and Jumpcode
467 Genomics teams for the excellent technical support.

468 **Authors' contributions.** Conceptualization MR, MD and RP; Methodology MR, LM and MD; Software LM
469 and GL; Investigation LDA, EBE, GC and FL; Formal analysis MR, LM and LDA; Validation GF, LN and LV;
470 Resources KB, LR and DJK; Data Curation LM and GL; Writing - Original Draft MR; Writing - Review &
471 Editing LM, LDA and MD; Visualization MR and LM; Supervision MR and MD; Project administration MR;
472 Funding acquisition EBI, MD and RP.

473

474 REFERENCES

- 475 Baird NA, Etter PD, Atwood TS, Currey MC, Shiver AL, Lewis ZA, Selker EU, Cresko WA, Johnson EA. 2008.
476 Rapid SNP discovery and genetic mapping using sequenced RAD markers. *PLoS One* **3**.
- 477 Bellucci E, Mario Aguilar O, Alseekh S, Bett K, Brezeanu C, Cook D, De la Rosa L, Delledonne M, Dostatny DF,
478 Ferreira JJ, et al. 2021. The INCREASE project: Intelligent Collections of food-legume genetic resources
479 for European agrofood systems. *Plant Journal* **108**: 646–660.
- 480 Bennetzen JL, Wang H. 2014. The contributions of transposable elements to the structure, function, and
481 evolution of plant genomes. *Annu Rev Plant Biol* **65**: 505–530.
- 482 Cooke TF, Yee MC, Muzzio M, Sockell A, Bell R, Cornejo OE, Kelley JL, Bailliet G, Bravi CM, Bustamante CD,
483 et al. 2016. GBStools: A Statistical Method for Estimating Allelic Dropout in Reduced Representation
484 Sequencing Data. *PLoS Genet* **12**.
- 485 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T, McCarthy SA,
486 Davies RM. 2021. Twelve years of SAMtools and BCFtools. *Gigascience* **10**.
- 487 Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML. 2011. Genome-wide genetic marker
488 discovery and genotyping using next-generation sequencing. *Nat Rev Genet* **12**: 499–510.
- 489 Deng T, Zhang P, Garrick D, Gao H, Wang L, Zhao F. 2022. Comparison of Genotype Imputation for SNP
490 Array and Low-Coverage Whole-Genome Sequencing Data. *Front Genet* **12**.
- 491 Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, Smith I, Tothova Z, Wilen C, Orchard
492 R, et al. 2016. Optimized sgRNA design to maximize activity and minimize off-target effects of CRISPR-
493 Cas9. *Nat Biotechnol* **34**: 184–191.

- 494 Elshire RJ, Glaubitz JC, Sun Q, Poland JA, Kawamoto K, Buckler ES, Mitchell SE. 2011. A robust, simple
495 genotyping-by-sequencing (GBS) approach for high diversity species. *PLoS One* **6**.
- 496 Feuillet C, Leach JE, Rogers J, Schnable PS, Eversole K. 2011. Crop genome sequencing: Lessons and
497 rationales. *Trends Plant Sci* **16**: 77–88.
- 498 Friel J, Bombarely A, Fornell CD, Luque F, Fernández-Ocaña AM. 2021. Comparative analysis of genotyping
499 by sequencing and whole-genome sequencing methods in diversity studies of *olea europaea* l. *Plants*
500 **10**.
- 501 Gargiulo R, Kull T, Fay MF. 2021. Effective double-digest RAD sequencing and genotyping despite large
502 genome size. *Mol Ecol Resour* **21**: 1037–1055.
- 503 Gu W, Crawford ED, O'Donovan BD, Wilson MR, Chow ED, Retallack H, DeRisi JL. 2016. Depletion of
504 Abundant Sequences by Hybridization (DASH): Using Cas9 to remove unwanted high-abundance
505 species in sequencing libraries and molecular counting applications. *Genome Biol* **17**.
- 506 Guerra-García A, Gioia T, von Wettberg E, Logozzo G, Papa R, Bitocchi E, Bett KE. 2021. Intelligent
507 Characterization of Lentil Genetic Resources: Evolutionary History, Genetic Diversity of Germplasm,
508 and the Need for Well-Represented Collections. *Curr Protoc* **1**.
- 509 Hadley W. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York
510 <https://ggplot2.tidyverse.org>.
- 511 He F, Pasam R, Shi F, Kant S, Keeble-Gagnere G, Kay P, Forrest K, Fritz A, Hucl P, Wiebe K, et al. 2019. Exome
512 sequencing highlights the role of wild-relative introgression in shaping the adaptive landscape of the
513 wheat genome. *Nat Genet* **51**: 896–904.
- 514 Hodges E, Xuan Z, Balija V, Kramer M, Molla MN, Smith SW, Middle CM, Rodesch MJ, Albert TJ, Hannon GJ,
515 et al. 2007. Genome-wide in situ exon capture for selective resequencing. *Nat Genet* **39**: 1522–1527.
- 516 Homberger C, Hayward RJ, Barquist L, Vogel J. 2023. Improved Bacterial Single-Cell RNA-Seq through
517 Automated MATQ-Seq and Cas9-Based Removal of rRNA Reads. *mBio*.
- 518 Ichida H, Abe T. 2019. An improved and robust method to efficiently deplete repetitive elements from
519 complex plant genomes. *Plant Science* **280**: 455–460.
- 520 Jayakodi M, Golicz AA, Kreplak J, Fehete LI, Angra D, Bednář P, Bornhofen E, Zhang H, Bousageon R, Kaur
521 S, et al. 2023. The giant diploid faba genome unlocks variation in a global protein crop. *Nature*.
- 522 JumpCode Genomics. 2021. Technology Overview Version 1.2 Harnessing CRISPR to boost NGS sensitivity
523 with CRISPRclean™. [https://www.jumpcodegenomics.com/wp-content/uploads/2021/07/jumpcode-](https://www.jumpcodegenomics.com/wp-content/uploads/2021/07/jumpcode-technical-overview-20210521_v1-1_F.pdf)
524 [technical-overview-20210521_v1-1_F.pdf](https://www.jumpcodegenomics.com/wp-content/uploads/2021/07/jumpcode-technical-overview-20210521_v1-1_F.pdf) (Accessed December 24, 2022).
- 525 Kreplak J, Madoui MA, Cápál P, Novák P, Labadie K, Aubert G, Bayer PE, Gali KK, Syme RA, Main D, et al.
526 2019. A reference genome for pea provides insight into legume genome evolution. *Nat Genet* **51**:
527 1411–1422.
- 528 Lachance J, Tishkoff SA. 2013. SNP ascertainment bias in population genetic analyses: Why it is important,
529 and how to correct it. *BioEssays* **35**: 780–786.

- 530 Le C, Liu Y, López-Orozco J, Joyce MA, Le XC, Tyrrell DL. 2021. CRISPR Technique Incorporated with Single-
531 Cell RNA Sequencing for Studying Hepatitis B Infection. *Anal Chem* **93**: 10756–10761.
- 532 Lee S-I, Kim N-S. 2014. Transposable Elements and Genome Size Variations in Plants. *Genomics Inform* **12**:
533 87.
- 534 Matvienko M, Kozik A, Froenicke L, Lavelle D, Martineau B, Perroud B, Michelmore R. 2013. Consequences
535 of Normalizing Transcriptomic and Genomic Libraries of Plant Genomes Using a Duplex-Specific
536 Nuclease and Tetramethylammonium Chloride. *PLoS One* **8**.
- 537 Md V, Misra S, Li H, Aluru S. 2019. Efficient architecture-aware acceleration of BWA-MEM for multicore
538 systems. In *Proceedings - 2019 IEEE 33rd International Parallel and Distributed Processing Symposium,*
539 *IPDPS 2019*, pp. 314–324, Institute of Electrical and Electronics Engineers Inc.
- 540 Miller DFB, Yan PS, Buechlein A, Rodriguez BA, Yilmaz AS, Goel S, Lin H, Collins-Burow B, Rhodes L V., Braun
541 C, et al. 2013. A new method for stranded whole transcriptome RNA-seq. *Methods* **63**: 126–134.
- 542 Montefiori L, Hernandez L, Zhang Z, Gilad Y, Ober C, Crawford G, Nobrega M, Sakabe NJ. 2017. Reducing
543 mitochondrial reads in ATAC-seq using CRISPR/Cas9. *Sci Rep* **7**.
- 544 Müller Paul H, Istanto DD, Heldenbrand J, Hudson ME. 2022. CROPSR: an automated platform for complex
545 genome-wide CRISPR gRNA design and validation. *BMC Bioinformatics* **23**.
- 546 Ogutcen E, Ramsay L, von Wettberg EB, Bett KE. 2018. Capturing variation in Lens (Fabaceae): Development
547 and utility of an exome capture array for lentil. *Appl Plant Sci* **6**.
- 548 Pavan S, Delvento C, Ricciardi L, Lotti C, Ciani E, D'Agostino N. 2020. Recommendations for Choosing the
549 Genotyping Method and Best Practices for Quality Control in Crop Genome-Wide Association Studies.
550 *Front Genet* **11**.
- 551 Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: An inexpensive
552 method for de novo SNP discovery and genotyping in model and non-model species. *PLoS One* **7**.
- 553 Prezza G, Heckel T, Dietrich S, Homberger C, Westermann AJ, Vogel J. 2020. Improved bacterial RNA-seq by
554 Cas9-based depletion of ribosomal RNA reads. <http://www.rnajournal.org/cgi/doi/10.1261/rna>.
- 555 Quinlan AR, Hall IM. 2010. BEDTools: A flexible suite of utilities for comparing genomic features.
556 *Bioinformatics* **26**: 841–842.
- 557 R Core Team 2020. R: A language and environment for statistical computing. R Foundation for Statistical
558 Computing, Vienna, Austria. <https://www.r-project.org/>.
- 559 Ramsay L, Koh CS, Kagale S, Gao D, Kaur S, Haile T, Gela TS, Chen L-A, Cao Z, Konkin DJ, et al. 2021. Genomic
560 rearrangements have consequences for introgression breeding as revealed by genome assemblies of
561 wild and cultivated lentil species. <https://doi.org/10.1101/2021.07.23.453237>.
- 562 Ren Q, Wang Y, Lin Y, Zhen Z, Cui Y, Qin S. 2021. The extremely large chloroplast genome of the green
563 alga *Haematococcus pluvialis*: Genome structure, and comparative analysis. *Algal Res* **56**.

- 564 Ricci WA, Lu Z, Ji L, Marand AP, Ethridge CL, Murphy NG, Noshay JM, Galli M, Mejía-Guerra MK, Colomé-
565 Tatché M, et al. 2019. Widespread long-range cis-regulatory elements in the maize genome. *Nat*
566 *Plants* **5**: 1237–1249.
- 567 Sakamoto W, Takami T. 2018. Chloroplast DNA dynamics: Copy number, quality control and degradation.
568 *Plant Cell Physiol* **59**: 1120–1127.
- 569 Siddique A, Suckow G, Ordoukhanian P, Head S, Homer N, Hernandez A, Brown K, Glick L, Baruch K, Doran
570 P, et al. 2019. RipTide High Throughput NGS Library Prep for Genotyping in Populations. *J Biomol*
571 *Tech.*; **30**(Suppl):S35-S36.
- 572 Sloan DB, Warren JM, Williams AM, Wu Z, Abdel-Ghany SE, Chicco AJ, Havird JC. 2018. Cytonuclear
573 integration and co-evolution. *Nat Rev Genet* **19**: 635–648.
- 574 Tanaka N, Shenton M, Kawahara Y, Kumagai M, Sakai H, Kanamori H, Yonemaru JI, Fukuoka S, Sugimoto K,
575 Ishimoto M, et al. 2021. Investigation of the Genetic Diversity of a Rice Core Collection of Japanese
576 Landraces using Whole-Genome Sequencing. *Plant Cell Physiol* **61**: 2087–2096.
- 577 Tian F, Yang DC, Meng YQ, Jin J, Gao G. 2020. PlantRegMap: Charting functional regulatory maps in plants.
578 *Nucleic Acids Res* **48**: D1104–D1113.
- 579 Truong HT, Ramos AM, Yalcin F, de Ruiter M, van der Poel HJA, Huvenaars KHJ, Hogers RCJ, van Enckevort
580 LJG, Janssen A, van Orsouw NJ, et al. 2012. Sequence-based genotyping for marker discovery and co-
581 dominant scoring in germplasm and populations. *PLoS One* **7**.
- 582 Van der Auwera GA, O'Connor BD. 2020. Genomics in the Cloud: Using Docker, GATK, and WDL in Terra (1st
583 Edition). O'Reilly Media.
- 584 Wang J, Sun G, Ren X, Li C, Liu L, Wang Q, Du B, Sun D. 2016. QTL underlying some agronomic traits in
585 barley detected by SNP markers. *BMC Genet* **17**.
- 586 Wang P, Xiong Y, Gong R, Yang Y, Fan K, Yu S. 2019. A key variant in the cis-regulatory element of flowering
587 gene *Ghd8* associated with cold tolerance in rice. *Sci Rep* **9**.
- 588 Yan H, Haak DC, Li S, Huang L, Bombarely A. 2022. Exploring transposable element-based markers to
589 identify allelic variations underlying agronomic traits in rice. *Plant Commun* **3**.
- 590 Zan Y, Payen T, Lillie M, Honaker CF, Siegel PB, Carlborg Ö. 2019. Genotyping by low-coverage whole-
591 genome sequencing in intercross pedigrees from outbred founders: A cost-efficient approach.
592 *Genetics Selection Evolution* **51**.
- 593 Zhang GJ, Dong R, Lan LN, Li SF, Gao WJ, Niu HX. 2020. Nuclear integrants of organellar DNA contribute to
594 genome structure and evolution in plants. *Int J Mol Sci* **21**.
- 595 Zhao W, He X, Hoadley KA, Parker JS, Hayes DN, Perou CM. 2014. *Comparison of RNA-Seq by poly (A)*
596 *capture, ribosomal RNA depletion, and DNA microarray for expression profiling.*
597 <http://www.biomedcentral.com/1471-2164/15/419>.

598 Zhulidov PA, Bogdanova EA, Shcheglov AS, Vagner LL, Khaspekov GL, Kozhemyako VB, Matz M V.,
 599 Meleshkevitch E, Moroz LL, Lukyanov SA, et al. 2004. Simple cDNA normalization using kamchatka
 600 crab duplex-specific nuclease. *Nucleic Acids Res* **32**.

601

602

603

604 **FIGURE LEGEND**

605 **Figure 1. Distribution of mapped reads after CRISPR-Cas9-mediated repeat depletion.** Libraries of *L.*
 606 *culinaris* cv. Castelluccio DNA (8-plex) were depleted using the custom gRNA set and Cas9. The gRNAs
 607 were used simultaneously (All) or were split into three groups that were used sequentially in order of
 608 increasing cutting frequency (3 grp). Bar graphs show the number (in millions) of reads mapping to repetitive
 609 DNA (**A**), to nuclear repetitive DNA (**B**), to the chloroplast genome (**C**), to the mitochondrial genome (**D**), and
 610 to the single-copy regions of the nuclear genome (**E**). Data are means \pm SD (n = 3 for each condition; *p-adj
 611 < 0.05, **p-adj < 0.01, ****p-adj < 0.0001; one-way ANOVA plus Tukey's multiple comparisons test; ns, not
 612 significant).

613 **Figure 2. Reads mapped to repetitive or single-copy regions of the lentil genome with or without**
 614 **CRISPR-Cas9-mediated repeat depletion.** Integrative Genome Browser Visualization (IGV) of Illumina
 615 sequencing data mapped to two representative genomic sites of ~180 and ~50 kbp. Tracks in blue, green and
 616 red represent annotated genes, single-copy regions and repetitive regions, respectively.

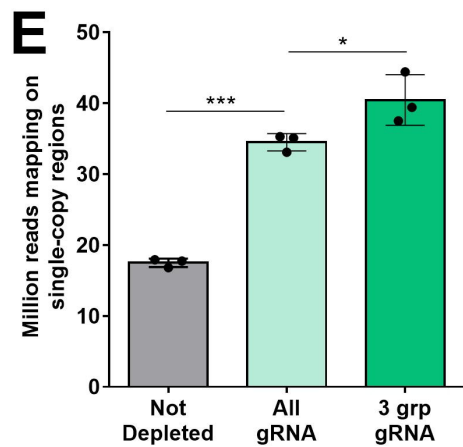
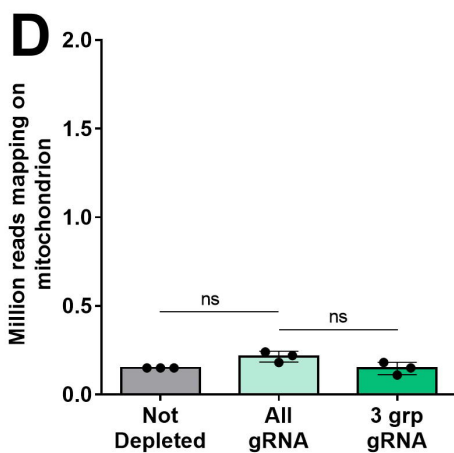
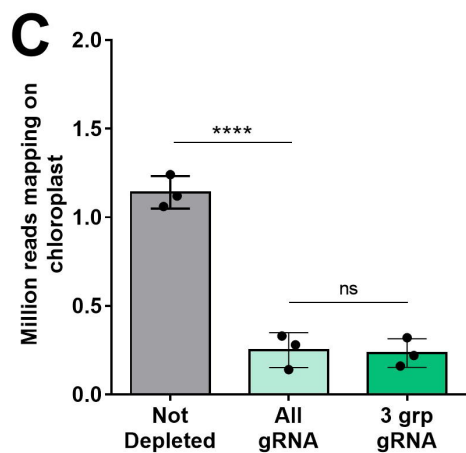
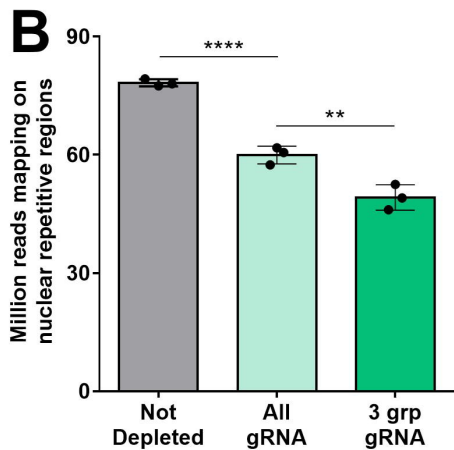
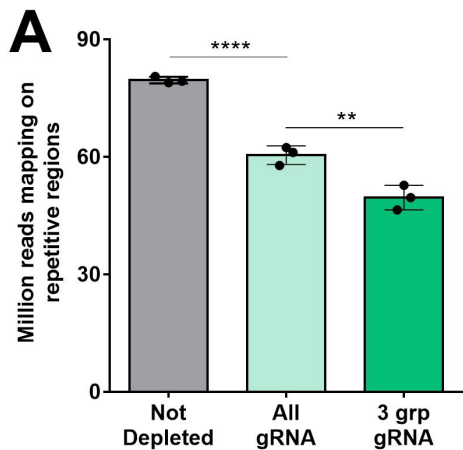
617 **Figure 3. Sequencing data distribution after CRISPR-Cas9-mediated repeat depletion according to**
 618 **the class of nuclear repeat.** (**A**) Number of reads (in millions) mapping to different nuclear repeat classes
 619 with or without CRISPR-Cas9-mediated repeat depletion: Ty3-Gypsy (LTR Gypsy), Ty1-Copia (LTR Copia),
 620 other LTR, LINE, CACTA, Mu, hAT, Helitron transposons, other transposable elements with abundance <
 621 0.1% (Harbinger, mariner, Sine), tandem repeats (TR) and simple repeats (SR). Correlation between the
 622 variation of mapped reads following CRISPR-Cas9-mediated repeat depletion on the same repeat classes
 623 versus the repeat length (**B**), repeat copy number (**C**) and density of gRNAs targeting each repeat class (**D**).
 624 Data are means \pm SE (n = 3). FC – fold change.

625 **Figure 4. Genotyping performance on single-copy regions with or without CRISPR-Cas9-mediated**
626 **repeat depletion starting from different amounts of sequencing data. (A)** Number of genotyped
627 positions . **(B)** Number of total variants identified . **(C)** Number of heterozygous variants identified. **(D)**
628 Number of variants that were genotyped as reference (0/0) in the non-depleted sample and identified as
629 homozygous (1/1) or heterozygous (1/0) alternative in the depleted sample . N=3 at 6, 12 and 25 million
630 fragments, N=2 at 37, 50 and 63 million fragments.

631 **Figure 5. Performances of CRISPR-Cas9-mediated repeat depletion in different lentil samples. (A-B)**
632 Sequencing reads mapping to the repetitive or single-copy regions with or without CRISPR-Cas9-mediated
633 repeat depletion in multiplex libraries generated from different lentil samples, namely *L. culinaris* cv.
634 Redberry (*L. c.* RB) or *L. orientalis* (*L. o.*). Data are means \pm SE (n = 2) normalized for the same sequencing
635 input (50 million fragments). Variation percentages observed following CRISPR-Cas9-mediated repeat
636 depletion are reported above each condition. **(C)** Number of genotyped positions in common between all
637 samples analyzed in the study (3 distinct samples of *L. culinaris* cv. Castelluccio, 1 sample of *L. culinaris* cv.
638 Redberry, and 1 samples of *L. orientalis*), within the single-copy regions. N=5 at 6, 12 and 25 million
639 fragments, N=4 at 37, 50 and 63 million fragments.

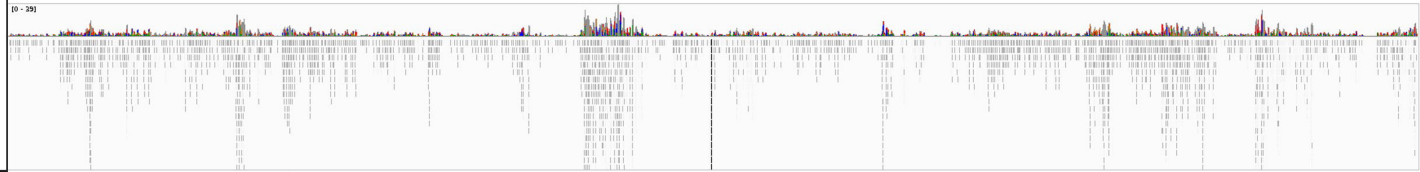
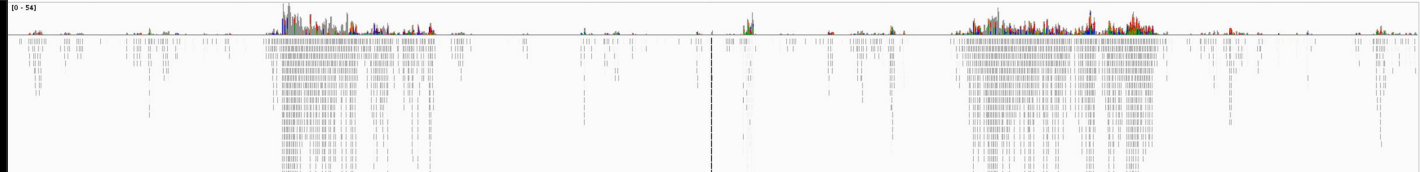
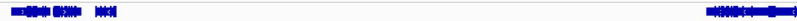
640 **Figure 6. Performances of CRISPR-Cas9-mediated repeat depletion at different multiplexing level and**
641 **library types.** Sequencing reads mapping to the repetitive or single-copy regions with or without CRISPR-
642 Cas9-mediated repeat depletion in **(A-B)** multiplex libraries containing 8 or 96 samples (plx) or **(E-F)**
643 standard singleplex WGS libraries generated from one *L. culinaris* cv. Castelluccio, one *L. culinaris* cv.
644 Redberry and one *L. orientalis* sample. Data are means \pm SE (n = 3) normalized for the same sequencing
645 input (50 million fragments). Variation percentages observed following CRISPR-Cas9-mediated repeat
646 depletion are reported above each condition.

647

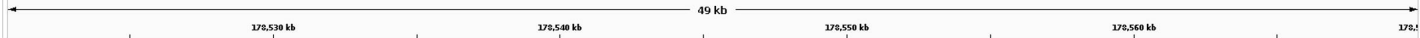
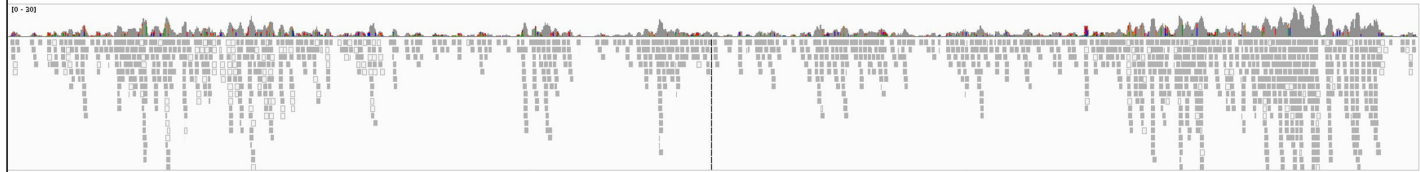
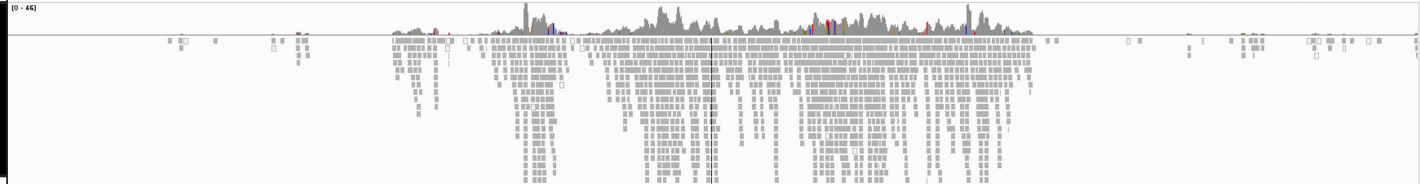


A

Lcu.2RBV.Chr1:308,341,918-308,522,842

**Not Depleted****Depleted****Genes****Single-copy****Repetitive****B**

Lcu.2RBV.Chr6:178,520,751-178,569,891

**Not Depleted****Depleted****Genes****Single-copy****Repetitive**

