

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

Marzia Rossato^{1,2*}, Luca Marcolungo^{1*}, Luca De Antoni¹, Giulia Lopatriello¹, Elisa Bellucci³, Gaia Cortinovi³, Giulia Frascarelli³, Laura Nanni³, Elena Bitocchi³, Valerio Di Vittori³, Leonardo Vincenzi¹, Filippo Lucchini¹, Kirstin E. Bett⁴, Larissa Ramsay⁴, David James Konkin⁵, Massimo Delledonne^{1,2§} and Roberto Papa^{3§}

*equal contribution

§corresponding authors

Marzia.rossato@univr.it; Luca.marcolungo@univr.it; Luca.deantoni@univr.it; giulia.lopatriello@univr.it;
e.bellucci@staff.univpm.it; gaia.cortinovi93@gmail.com; g.frascarelli@pm.univpm.it; l.nanni@staff.univpm.it;
e.bitocchi@staff.univpm.it; v.divittori@staff.univpm.it; leonardo.vincenzi@univr.it; filippo.lucchini@univr.it;
k.bett@usask.ca; l.ramsay@usask.ca; David.Konkin@nrc-cnrc.gc.ca; massimo.delledonne@univr.it;
r.papa@staff.univpm.it

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

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

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

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

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

Running Title: Genome sequencing based on repeats exclusion

ABSTRACT

High-throughput genotyping enables the large-scale analysis of genetic diversity in population genomics and genome-wide association studies that combine the genotypic and phenotypic characterization of large collections of accessions. Genotyping by sequencing is progressively replacing traditional genotyping methods due to the lower ascertainment bias. However, genome-wide genotyping by sequencing becomes expensive in species with large genomes and a high proportion of repetitive DNA. Here we describe the use of CRISPR/Cas9 technology to deplete repetitive elements in the 3.76-Gb genome of lentil (*Lens culinaris*), 84% consisting of repeats, thus concentrating the sequencing data on coding and regulatory regions (unique regions). We designed a custom set of 566,722 gRNAs targeting 2.9 Gbp of repeats and excluding repetitive regions overlapping annotated genes and putative regulatory elements based on ATAC-Seq data. The novel depletion method removed 40% of reads mapping to repeats, increasing those mapping to unique regions by 2.6-fold. This repeat-to-unique shift in the sequencing data increased the number of genotyped bases by up to 17-fold compared to non-depleted libraries. We were also able to identify up to 18-fold more genetic variants in the unique regions and increased the genotyping accuracy by rescuing thousands of heterozygous variants that otherwise would be missed due to low coverage. The method performed similarly regardless of the multiplexing level, type of library or genotypes, including different cultivars and a closely-related species (*L. orientalis*). Our results demonstrated that CRISPR/Cas9-driven repeat depletion focuses sequencing data on meaningful genomic regions, thus improving high-density and genome-wide genotyping in large and repetitive genomes.

KEYWORDS: repetitive elements, CRISPR/Cas9, genotyping by sequencing, next-generation sequencing library

INTRODUCTION

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

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

More recently, next generation sequencing (NGS) has provided an opportunity to discover genome-wide variants in an unbiased manner. Genotyping by sequencing (GBS) involves low coverage (5–10x) whole-genome sequencing (lcWGS), allowing the characterization of several million variants (Tanaka et al. 2021; Friel et al. 2021). To reduce costs enough to make WGS affordable even in large germplasm collections, very low coverage (0.5–2x) WGS (ultra-lcWGS) can be combined with imputation to infer positions that are not sequenced or genotyped (Deng et al. 2022; Zan et al. 2019; Wang et al. 2016). Alternatively, sequencing costs are often minimized by reduced-representation sequencing, which comprises methods such as restriction site-associated DNA sequencing (RAD-Seq) (Davey et al. 2011) and double-digest RAD-Seq (ddRAD-seq) (Truong et al. 2012; Peterson et al. 2012). These methods concentrate sequencing data on regions adjacent to restriction sites by exploiting the specificity of restriction endonucleases. Reduced-representation sequencing is suitable for large cohorts, but provides only low-resolution data, with a small fraction of analysed and genotyped bases (Pavan et al. 2020). In large genomes, such that of *L. culinaris* (3.76 Gb), this may not provide sufficient marker density and depth to confidently identify variants under selection (Guerra-García et al. 2021). The resolution can be increased without significantly greater costs in sample prep by using the Twist 96-Plex Library Prep Kit (formerly iGenomX Riptide

Kit) to generate multiplexed libraries, allowing 96–960 samples to be processed simultaneously and resulting in the non-random sampling of millions of genomic positions (Siddique et al. 2019).

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

Ideally, lcWGS could be focused on the most complex parts of the genome, avoiding wasted effort on the sequencing of repetitive elements. This could be achieved by using enzymes that enable target enrichment by depleting unwanted sequences from NGS libraries. For example, the duplex-specific nuclease (DSN) selectively digests double-stranded DNA molecules, and can be used to eliminate highly abundant sequences in a controlled denaturation-reassociation reaction (Zhulidov et al. 2004). This method has been used in RNA-Seq analysis to remove abundant transcripts (Zhao et al. 2014; Miller et al. 2013) and, just occasionally, also to delete repetitive elements in DNA-Seq libraries generated from plant genomes (Ichida and Abe 2019; Matvienko et al. 2013). However, the use of this approach is limited by the need to define an optimal renaturation time for each specific genome, that can vary from some hours to a day (Matvienko et al. 2013). DSN also removes informative repetitive elements, such as the coding sequences of abundant gene families, which are particularly relevant in polyploid plants arising from whole genome duplication events (Matvienko et al. 2013). More recently, the CRISPR/Cas9 (clustered regularly interspaced short palindromic repeats and CRISPR-associated nuclease 9)

system has been used for the selective depletion of unwanted genome fractions from sequencing libraries (Gu et al. 2016). The Cas9 enzyme can be programmed to cut library fragments by designing specific guide-RNA sequences targeting the unwanted sequences. Subsequently, only intact fragments -retaining adapters at both ends- can be effectively amplified by PCR and generate productive clusters on a sequencing flow-cell. The DASH approach (depletion of abundant sequences by hybridization) involved the use of Cas9 to exclude ribosomal RNA (rRNA) sequences from RNA-Seq libraries and to remove DNA from common pathogens in order to detect rare pathogens in metagenomic samples (Gu et al. 2016). A similar technology has been recently commercialized under the name “CRISPRclean” by JumpCode Genomics (JumpCode Genomics 2021).

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

RESULTS

Depletion of *L. culinaris* repetitive DNA using CRISPR/Cas9

We designed a custom set of gRNAs to deplete the repetitive DNA content of the *L. culinaris* CDC Redberry genome (Ramsay et al.), targeting transposable elements (3.1 Gbp, 83%), simple sequence repeats (303 Mbp, 8%) and tandem repeats (13 Mbp, 0.4%) in the nuclear genome, as well as the entire mitochondrial genome (mtDNA, 489 kbp) and chloroplast genome (chlDNA, 118 kbp) (**Supplemental_Table_S1.xlsx**). We excluded repetitive DNA that overlapped with functional regions such as annotated genes (185 Mbp, 5%) and putative regulatory regions identified using ATAC-Seq data (78 Mbp, 2%) (**Supplemental_Table_S1.xlsx**). All nuclear DNA outside the gRNA target regions is hereafter defined as unique. The final design comprised 566,722 gRNAs with at least 25 recognition sites, potentially targeting 2.9 Gbp (77%) of the *L. culinaris* nuclear genome and 93.5% of its repetitive regions when using a sequencing library with 500-bp inserts (**Supplemental_Table_S2.xlsx**). An additional 2,366 gRNAs targeted the mitochondrial and chloroplast genomes (**Supplemental_Table_S1.xlsx**). The gRNAs were assigned to 11 pools based on their cutting frequency in the *L. culinaris* genome (**Supplemental_Table_S2.xlsx**).

The custom gRNAs were tested on three Twist 8-plex libraries, allowing the reproducible sampling of the same genomic regions by random priming during first-strand DNA synthesis. Each multiplex library comprised eight replicates of three distinct *L. culinaris* samples (cv. Castelluccio) (**Supplemental_Table_S3.xlsx**). Cas9/gRNAs ribonucleoprotein (RNP) complexes (1:10 protein/gRNA ratio) were generated, and gRNAs with more target sites in the genome were used at higher relative concentrations in the final reaction (**Supplemental_Table_S2.xlsx**). Depletion reactions in the presence of RNP complexes were conducted either using all gRNAs simultaneously or by splitting the gRNA pools into three groups based on cutting frequency (**Supplemental_Table_S2.xlsx**) and using the groups sequentially, starting with the lowest cutting frequency. Depleted and non-depleted libraries were sequenced, generating 91 million fragments on average (**Supplemental_Table_S4.xlsx**). The sequencing data were normalized at ~50 million fragments per library in order to compare the proportion of reads mapping on repetitive and unique regions of the nuclear genome and on the organelle genomes (**Figure 1**). The number of reads mapping to repetitive regions (total repeats, nuclear repeats, mtDNA and chlDNA) was significantly lower in the depleted libraries compared to the non-depleted libraries, with the sequential depletion strategy using three gRNA groups performing best and depleting 37.7% of the repetitive DNA (**Figure 1A**). The results were similar when considering only the nuclear repetitive regions (37.2% depletion) (**Figure 1B**). A small fraction of total reads mapped to the organelle genomes (~1%). Both depletion strategies were similarly effective in the chloroplast genome, resulting in ~78.5% depletion (**Figure 1C**). In contrast, no significant depletion was observed in the mitochondrial genome (**Figure 1D**). In parallel, the sequential depletion strategy achieved a 130% increase in the number of reads mapping to unique regions, from 17.5 to 40.4 million (**Figure 1E**). Given that the concentration of Cas9 RNPs influences the cutting efficiency (Gu et al. 2016), we repeated the sequential depletion strategy using double amount of Cas9 and gRNA. This modified the read distribution further, achieving 41.2% depletion of nuclear repeats and a 160% increase in reads mapped to unique regions. This double sequential depletion strategy was therefore the most efficient, and was applied in all subsequent experiments. Overall, our results demonstrated that the custom gRNA set and Cas9 effectively targeted fragments containing repetitive DNA sequences and depleted them in the resulting sequencing libraries. Manual inspection of mapped reads confirmed that less sequencing data was assigned to regions of repetitive DNA and more reads were mapped to unique parts of the genome (**Figure 2A-B**).

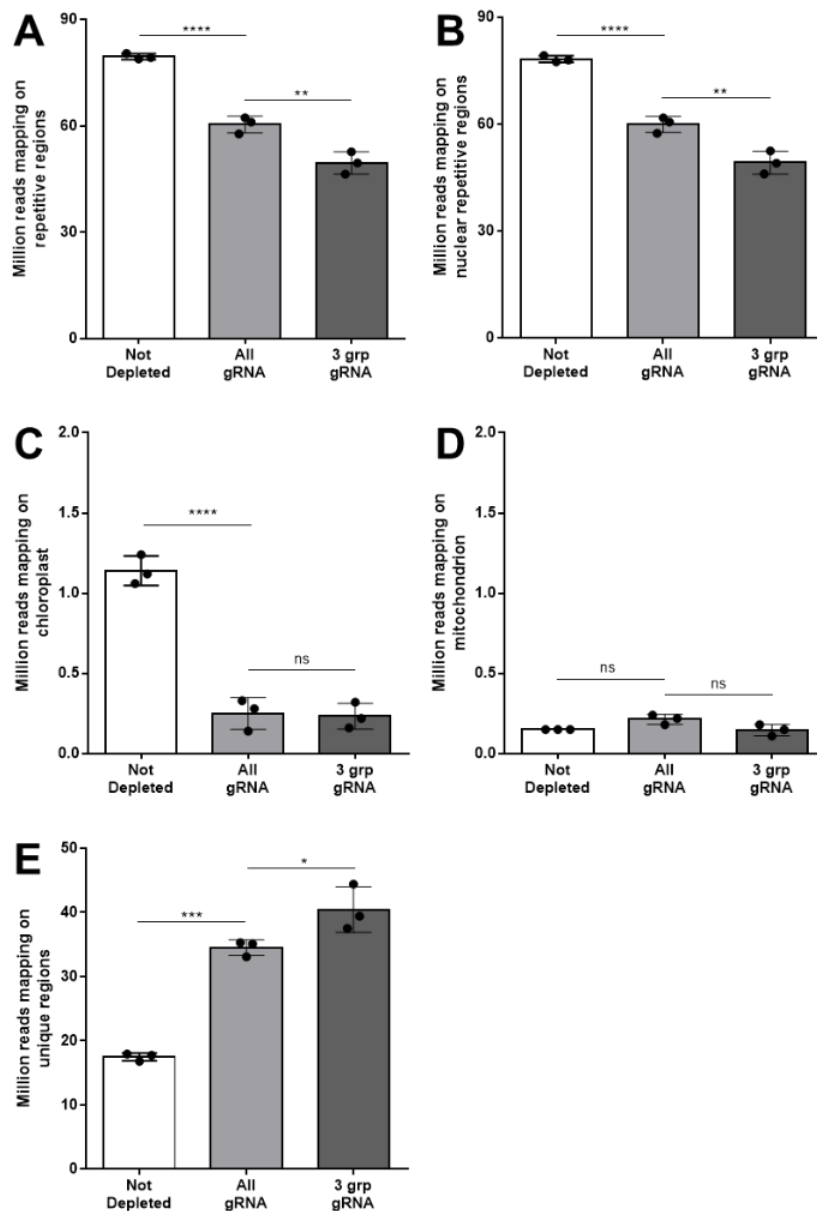


Figure 1. Distribution of mapped reads after CRISPR/Cas9-mediated repeat depletion. Libraries of *L. culinaris* cv. Castelluccio DNA (8-plex) were depleted using the custom gRNA set and Cas9. The gRNAs were used simultaneously (All) or were split into three groups that were used sequentially in order of increasing cutting frequency (3 grp). Bar graphs show the number (in millions) of reads mapping to repetitive DNA (A), to nuclear repetitive DNA (B), to the chloroplast genome (C), to the mitochondrial genome (D), and to the unique regions of the nuclear genome (E), starting from 100 million-normalized reads (50 million fragments). Data are means \pm SD ($n = 3$ for each condition; * p -adj < 0.05 , ** p -adj < 0.01 , **** p -adj < 0.0001 ; one-way ANOVA plus Tukey's multiple comparisons test; ns, not significant).

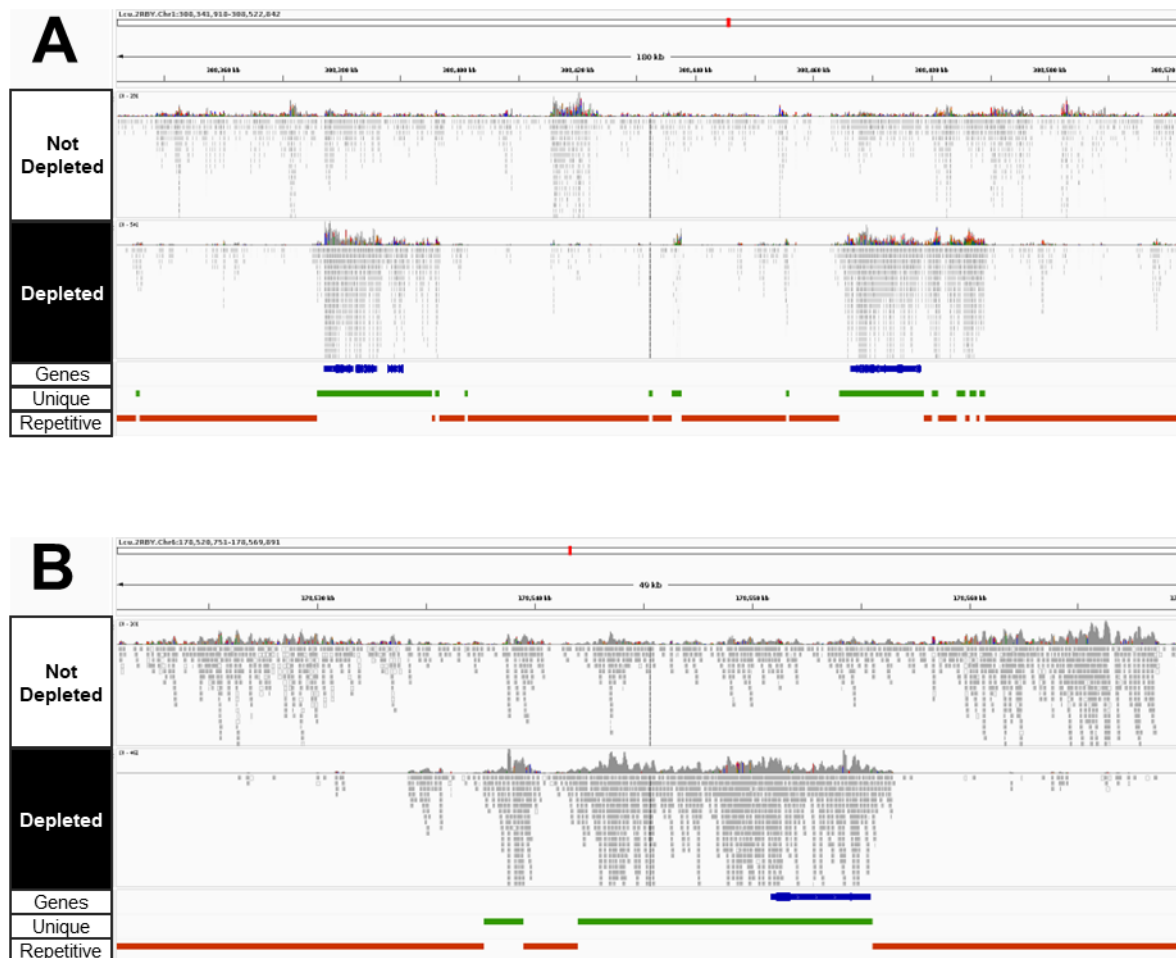


Figure 2. Reads mapped to repetitive or unique regions of the lentil genome before and after CRISPR/Cas9-mediated repeat depletion.

Integrative Genome Browser Visualization (IGV) of Illumina sequencing data mapped to two genomic sites of ~180 and ~50 kbp, before and after CRISPR/Cas9-mediated repeat depletion. Tracks in blue, green and red represent annotated genes, unique regions and repetitive regions, respectively.

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

We next examined the depletion of different classes of repetitive sequences in the *L. culinaris* genome. Reads mapping to the most abundant retroelements, namely the Ty3-Gypsy family (64% of the genome (Ramsay et al.), were reduced by 47% in the depleted libraries, whereas those mapping to the Ty3-Copia family (15% of the genome) and other long terminal repeat (LTR) elements (3% of the genome) were depleted by 3% and 38%, respectively (**Figure 3A** and **Supplemental_Table_S5.xlsx**). In contrast, there was no decrease in the abundance of other transposable elements (Line, CACTA, mu, hAT, Helitron, Harbinger, mariner and Sine), each representing < 1% of the genome

(Supplemental_Table_S5.xlsx), and there was no reduction in the number of reads mapping to tandem repeats (0.4% of the genome) or simple sequence repeats (8% of the genome) (Figure 3A). We observed a significant correlation between the variation in mapped reads after depletion and the abundance of these repeat classes in terms of overall repeat length and occurrence in the genome (Figure 3B-C and Supplemental_Table_S5.xlsx). Given that the number of gRNAs targeting each repeat class increased proportionally with the repeat size and occurrence, the most efficiently depleted repetitive elements also featured a higher density of gRNA targets (Figure 3D). There was a significant correlation between the variation of mapping reads following depletion and the gRNA density over the whole target region when considering each single cut site in the genome (Supplemental_Fig_S1.docx). In particular, target regions with a density > 8 gRNAs/kbp showed a negative read variation in 85% of cases (Supplemental_Fig_S1.docx) whereas regions targeted by < 8 gRNAs/kbp usually showed limited or no depletion (Supplemental_Fig_S1.docx). We therefore concluded that the depletion efficiency across different repeat classes was dependent on the density of gRNA targets.

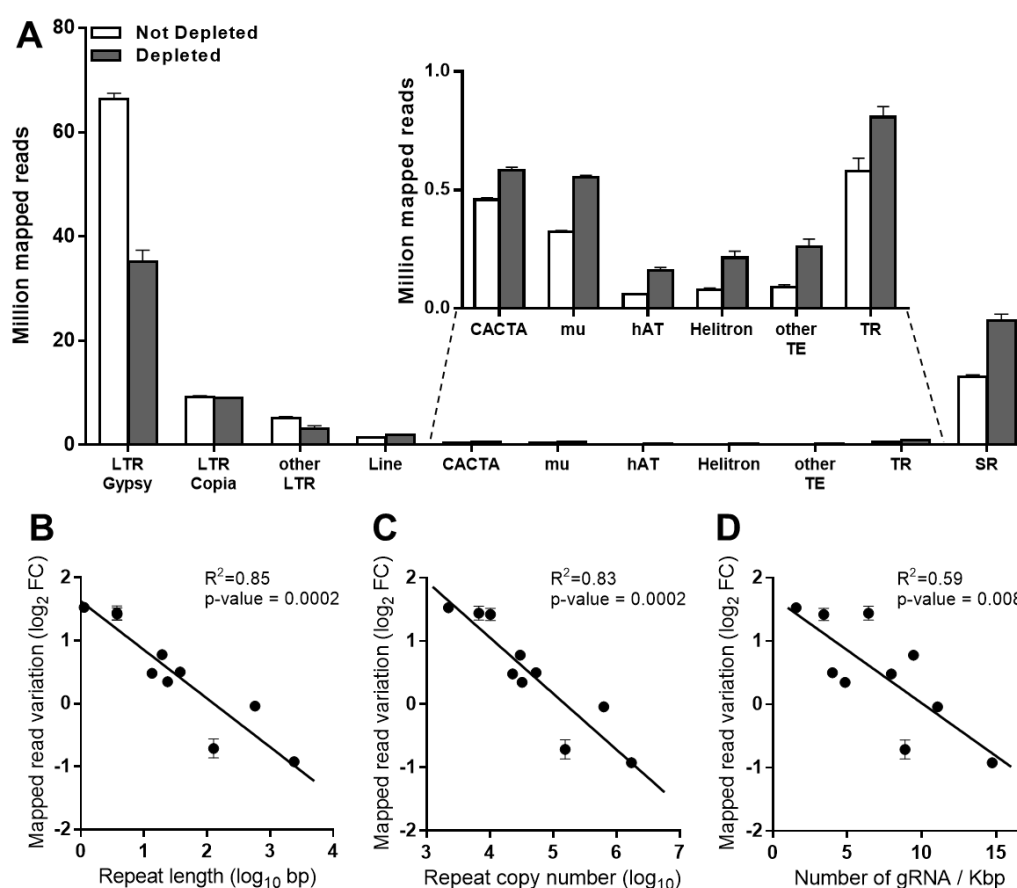


Figure 3. Sequencing data distribution after CRISPR/Cas9-mediated repeat depletion according to the class of nuclear repeat. (A) Number of reads (in millions) mapping to different nuclear repeat classes before and after CRISPR/Cas9-mediated repeat depletion: Ty3-Gypsy (LTR Gypsy), Ty1-Copia (LTR Copia), other LTR, Line, CACTA, mu, hAT, Helitron transposons, other transposable elements with abundance < 0.1% (Harbinger, mariner, Sine), tandem

repeats (TR) and simple repeats (SR). Correlation between the variation of mapped reads following CRISPR/Cas9-mediated repeat depletion on the same repeat classes versus the repeat length (B), repeat copy number (C) and density of gRNAs targeting each repeat class (D). Data are means \pm SE (n = 3). FC – fold change.

Impact of CRISPR/Cas9-mediated repeat depletion on genotyping accuracy

Next, we investigated the impact of CRISPR/Cas9-mediated repeat depletion on the number of genomic positions in the unique regions where a base can be reliably genotyped (% PASS at a depth of ≥ 5 reads). For this analysis, sequencing data generated from depleted and non-depleted *L. culinaris* cv. Castelluccio samples were downsampled from 62 to 6 million fragments to mimic lcWGS and ultra-lcWGS. Consistently more bases were genotyped within the unique regions of the depleted sample over the whole range considered (Figure 4A). Because a genotyped position does not necessarily allow the variant to be identified (this also depends on allele coverage), we also determined the impact of repeat depletion on variant calling. Following depletion, the total number of variants identified in the unique regions increased significantly by ~ 10 -fold on average, with a delta of $\sim 20,000$ to $\sim 300,000$ more variants identified in the depleted sample (Figure 4B). This allowed us to identify up to $\sim 50,000$ variants that would not be detected without the depletion strategy, despite occupying genotypable positions in the non-depleted sample (Figure 4C). These false negative variants in the non-depleted sample were not called due to allelic imbalance caused by the low coverage (Supplemental_Fig_S2.docx). Consistently, most of them ($\sim 97\%$) were heterozygous variants (Figure 4C).

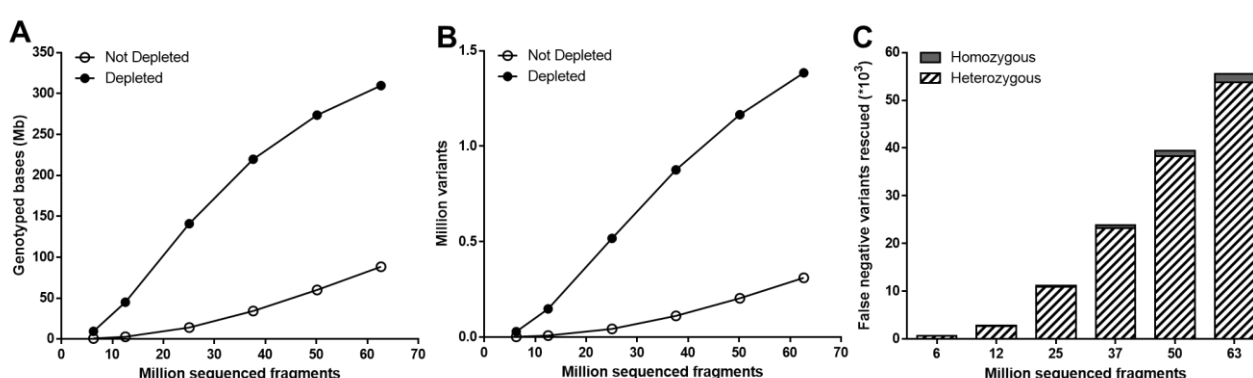


Figure 4. Genotyping performance on unique regions after CRISPR/Cas9-mediated repeat depletion. (A) Number of PASS positions within unique regions achieved with or without depletion, starting from different amounts of sequencing data. **(B)** Number of variants identified within unique regions with or without depletion, starting from different amounts of sequencing data. **(C)** Number of variants within unique regions that were located in a PASS position and identified in the depleted sample (1/1 or 1/0) but not in the non-depleted sample (0/0).

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

Finally, we assessed the performance of CRISPR/Cas9-mediated repeat depletion on different lentil genotypes, multiplexing levels and library types (**Supplemental_Table_S4.xlsx**). Similar variations in the coverage of repetitive/unique regions and the number of mapped reads were observed when depleting multiplex libraries generated from a different cultivar (RB, Redberry) or from the closely-related species *L. orientalis* (**Figure 5A-B**) when compared to the original Castelluccio cultivar (**Figure 1**). Furthermore, there was no significant difference in performance when library multiplexing was increased from 8-plex to 96-plex while maintaining the 1:10 Cas9:gRNA ratio and 1 ng of treated library per sample (**Figure 5C-D**), or when treating standard singleplex WGS libraries (**Figure 5E-F**). Overall, these results demonstrated that CRISPR/Cas9-mediated repeat depletion using the same gRNA set is at least equally effective when applied to a group of three lentil cultivars and one close wild relative, with the possibility to process individual samples or multiple samples simultaneously.

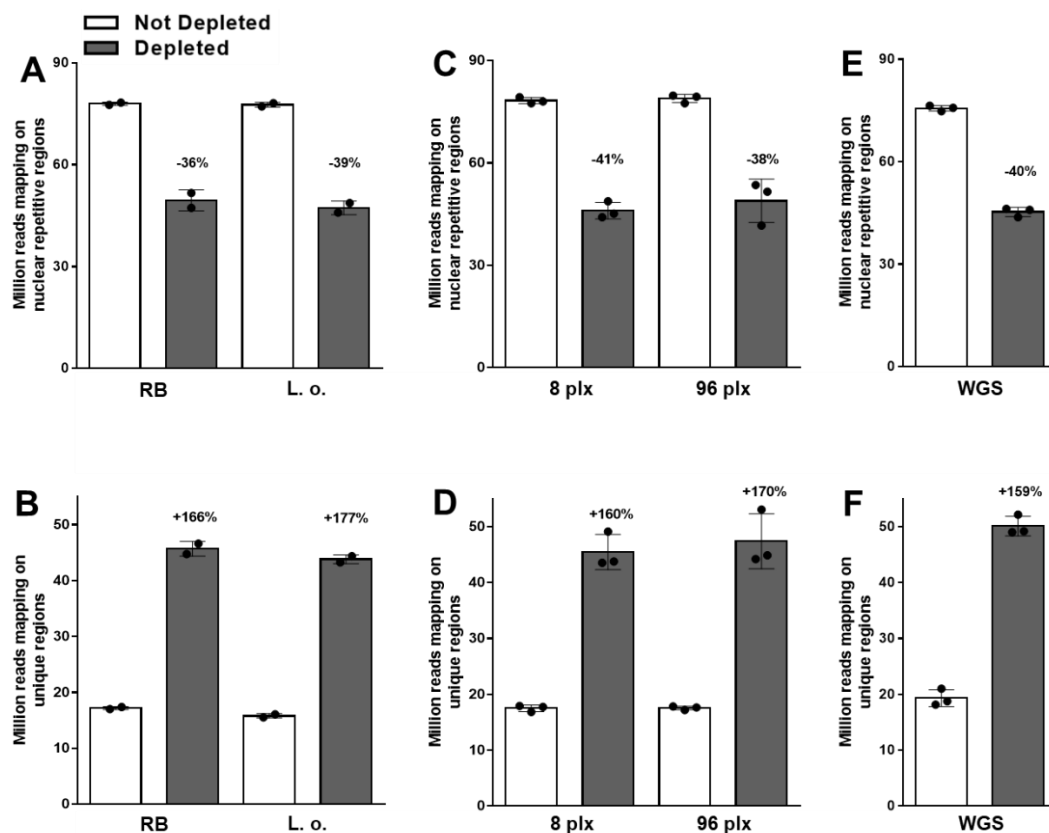


Figure 5. Variations of sequencing data distribution after CRISPR/Cas9-mediated repeat depletion in different lentil samples, multiplexing and library types. Sequencing reads mapping to the repetitive or unique regions before and after CRISPR/Cas9-mediated repeat depletion in **(A-B)** multiplex libraries generated from different lentil samples, namely *L. culinaris* cv. Redberry (RB) or *L. orientalis* (L. o.), **(C-D)** multiplex libraries containing 8 or 96 samples (plx) or **(E-F)** standard singleplex WGS libraries generated from the same samples. Data are means \pm SE ($n = 2-3$) normalized for the same sequencing input (50 million fragments). Variation percentages observed following CRISPR/Cas9-mediated repeat depletion are reported above each condition.

DISCUSSION

In a typical GBS experiment, the data derived from repetitive DNA is directly proportional to the repeat content of the genome, representing ~40% and ~80% for ddRAD-seq and lcWGS in *L. culinaris*, respectively. However, these data are largely uninformative. The bigger the genome, the more sequencing costs are therefore wasted on repeats. The traditional solution is the capture and sequencing of coding regions (WES). However, we approached the problem from the opposite perspective by depleting repetitive elements from the large genome of *L. culinaris* using CRISPR/Cas9 technology. Similar methods have been used for RNA-Seq library normalization (Prezza et al. 2020), metatranscriptomics (Gu et al. 2016), pathogen detection (Gu et al. 2016) and single-cell analysis (Parchman et al. 2018). Here, the main challenge is that 85% of the 3.7-Gb *L. culinaris* genome is repetitive DNA. The design of the gRNA array, therefore, required stringent multi-step filtering resulting in a set of ~600,000 gRNAs. To our knowledge, this is the first time CRISPR/Cas9 has been used with such a large number of gRNAs either *in vitro* or *in vivo*, representing a fundamental advance in the technology platform beyond the specific goals of our project. Despite these technical challenges, CRISPR/Cas9-mediated repeat depletion produced sequencing libraries with consistently lower proportions of repetitive DNA (~40% depletion) and enriched the unique regions, thus allowing the generation of more meaningful sequencing data. The depletion efficiency was comparable or even superior to the DSN method (Matvienko et al. 2013), but is also easier to standardize and allows a more specific targeting of repeats. CRISPR/Cas9-mediated repeat depletion allowed us to genotype up to 17-fold more bases on unique regions as compared to non-depleted libraries, and consequently we identified an average of 10-fold more genetic variants (up to 18-fold more in some cases). The shift in the distribution of sequencing data increased the coverage of unique regions by > 2.5-fold, and this was beneficial especially for the identification of heterozygous variants that would otherwise be missed due to unbalanced allelic sampling. Allele drop-

out is a well-known cause of errors in GBS experiments (Cooke et al. 2016). CRISPR/Cas9-mediated repeat depletion therefore improves the accuracy of genotyping experiments, especially in plants with highly heterozygous genomes.

Repetitive DNA in plant genomes can be divided into two broad categories: dispersed mobile elements and tandem/simple repeats. Dispersed mobile elements are made up of DNA transposons and retrotransposons, the most abundant of which are the LTR retrotransposons (Bennetzen and Wang 2014). Although mobile elements are not under the same selection pressure as genes, the degree of conservation across multiple copies of the same element is sufficiently high to allow the targeting of multiple copies with single gRNAs. The most efficient depletion (47%) was achieved for the most abundant LTR retrotransposon family (Gypsy, ~64% of the genome), followed by all the other LTR families (~15%). The cutting of LTR elements by Cas9 was responsible for almost the entire depletion observed at the genome-wide level, whereas the depletion of other mobile elements was negligible. This was probably due to the lower repetition of such elements, which translated into a poor cutting frequency in the final gRNA design, comprising only gRNAs with 25 targets. These targets usually featured < 8 gRNA/kbp, namely a density associated with a low depletion rate. A similar observation was reported for RNA-Seq libraries, where a gRNA every 50–100 nucleotides (10–20 gRNAs/kbp) achieves excellent depletion results (Gu et al. 2016). Simple and tandem repeats were also not depleted efficiently, although the gRNA density was close to 8. In these cases, the highly repetitive motifs of such sequences may have reduced the cutting efficiency of Cas9 (Müller Paul et al. 2022). Given that LTR retrotransposons make up the majority of repeats in plant genomes (94% in *L. culinaris*) and are the principal cause of plant genome size variation (Bennetzen and Wang 2014; Lee and Kim 2014), the design of gRNAs to target only LTR sequences may be the most efficient strategy to reduce the genome size in sequencing experiments. Future gRNA designs in other species and further optimization of the *L. culinaris* design should maximize the gRNA number on these most abundant elements, instead of dispersing the effort across the remaining repetitive fraction (< 10%). Another factor influencing the efficiency of CRISPR/Cas9-mediated repeat depletion was the dose of Cas9 and gRNAs; doubling their dose indeed improved repeat depletion by ~10%. Therefore, the gRNA target density and RNP concentration are two factors that can improve depletion efficiency, albeit with a slight increase in overall costs.

Organelle genomes often constitute a large fraction of DNA derived from plants (Sakamoto and Takami 2018). Although the mitochondrial and chloroplast genomes are smaller than the nuclear genome, they are present in multiple copies per cell, and they can represent > 20% of the total sequence data (Gargiulo et al. 2021; Ren et al. 2021). CRISPR/Cas9-mediated repeat depletion has been shown to reduce the fraction of sequencing libraries derived from organelle

genomes in ATAC-Seq experiments (Montefiori et al. 2017). Our method was efficient for the depletion of chloroplast DNA (by 67%) while the depletion of mitochondrial DNA was only marginal, possibly reflecting the different abundance of the two organelle genomes in the starting genomic sample. Still, in lentil, the fraction of sequencing data attributable to organelles was largely due to chloroplasts (88%), whose depletion was therefore sufficient to decrease the data mapping on organelles after Cas9 treatment (-65% overall). Although in the case of lentil the total sequencing data attributable to organelles was rather low (~1.3% in the not depleted libraries), the depletion of organelle DNA from sequencing libraries will be highly beneficial for organisms with a strongly unbalanced ratio of organelle vs nuclear DNA, such as *Cypripedium calceolus* (Gargiulo et al. 2021) and *Haematococcus pluvialis* (Ren et al. 2021).

Although the efficiency of CRISPR/Cas9-mediated repeat depletion could be improved, our results demonstrated that depleted libraries were more informative than standard ones when normalized for the amount of sequencing data. The coverage of unique regions increased, allowing us to genotype more bases and to discover more variants. Alternatively, CRISPR/Cas9-mediated repeat depletion can be used to reduce sequencing costs (by ~75% in lentil), because the same number of genotyped bases (or detected variants) in unique regions can be detected with much less sequencing data. The investment required to produce the gRNA set targeting *L. culinaris* repeats was similar to that necessary for WES custom probe synthesis, but the depletion approach is not restricted to coding regions and can also seek variants in regulatory regions, which are excluded by WES.

The cost of library preparation for a GBS project can easily exceed the cost of sequencing in the case of small genomes and/or ultra-lcWGS, especially given the steadily falling price of sequencing. More recent library-preparation kits circumvent several lengthy steps that require expensive reagents, and allow large sample sets to be processed in multiplex reactions. We used the Twist 96-Plex Library Prep Kit (formerly iGenomX Riptide) that constructs Illumina NGS libraries by polymerase-mediated extension of barcoded random primers. This type of library is beneficial for GBS in general because random priming reduces uniform genome coverage but allows more reproducible sampling of the same sites across multiple samples (Siddique et al. 2019). Most importantly, the kit is designed to process large numbers of samples (up to 96 simultaneously) at low costs and without advanced equipment (just a multichannel pipette). We estimated that the net cost to achieve 5-fold average coverage of the unique regions of the lentil genome by combining CRISPR/Cas9-mediated repeat depletion with a 96-plex Twist library is approximately US\$90, made up of US\$15 for library preparation, US\$10 for depletion and US\$65 for sequencing on a NovaSeq6000 S4 flowcell. This is similar to standard ddRAD-Seq costs (US\$ 40–60) but can genotype at least an order of magnitude more bases by focusing on unique regions.

CRISPR/Cas9-mediated repeat depletion combined with Twist multiplex libraries is therefore an effective strategy for genotyping projects involving hundreds or thousands of samples.

The repeat-to-unique shift in the profile of sequencing data allows to concentrate data on the same -relevant- portions across different libraries, thereby detecting a larger number of differences between samples. Overall this improves genotyping accuracy and reduces the fraction of missing data at the population level. This approach was successful in different cultivars of *L. culinaris*, and also in the closely-related species *L. orientalis*. This is important because genotyping experiments typically include distant/wild relatives and related species, from which it is possible to develop evolutionary studies and plan breeding experiments, including the introgression of characters of interest. For example, the INCREASE project features a collection of 2000 lentil accessions that includes both cultivated varieties and local landraces (Guerra-García et al. 2021). Further experiments could determine whether the gRNAs designed in this study are also suitable for the depletion of repeats in other closely related leguminous species (Fabaceae) with very large and repetitive genomes, such as pea (*P. sativum*, 3.92 Gb, 83% repetitive)(Kreplak et al. 2019) and faba bean (*Vicia faba* L., 12 Gb, 79% transposon-derived repeats)(Jayakodi et al.).

The novel depletion method based on CRISPR/Cas9 demonstrated efficient in decreasing the sequencing reads mapping on repetitive elements, while increasing those mapping to unique and functional regions. This can improve the high-density and genome-wide genotyping in large and repetitive genomes or, alternatively, reduce sequencing costs to achieve the same performances of non-depleted libraries. The method therefore has the potential to increase our genetic knowledge of plant species that are currently difficult to analyze without a significant economic investment due to the large genome size and high proportion of repetitive DNA. Population studies, eQTL analysis, GWAS and pre-breeding programs are just some of the approaches that can benefit from CRISPR/Cas9-mediated repeat depletion.

METHODS

Multiplex-library preparation and sequencing. We prepared 8-plex and 96-plex multiplex libraries according to the Twist 96-Plex Library Preparation Kit protocol (Twist Bioscience, South San Francisco, CA, USA) with the following modifications. For each sample, we denatured 100 ng of genomic DNA (25 ng/μl) at 98 °C for 1 min. Ultra-low (30%) GC random primer set A was used for the extension and termination reaction (Reaction A) followed by 8 and 9 cycles of PCR amplification for the 96-plex and 8-plex libraries, respectively. Final libraries were purified using Twist DNA Purification Beads (0.65x

volume) and a second round of purification was applied to the supernatant using 10 µl of beads to achieve a median insert size of 500 bp. Libraries were quantified using the Qubit BR DNA kit and a Qubit device (Thermo Fisher Scientific, Waltham, MA, USA) and size distributions were assessed using a Tape Station System (Agilent Technologies, Santa Clara, CA, USA). Non-depleted libraries were pooled at equimolar concentrations and sequenced on a NovaSeq6000 instrument (Illumina, San Diego, CA, USA) to generate 150-bp paired-end reads.

WGS library preparation and sequencing. Genomic DNA samples were fragmented using a Covaris sonicator to achieve an average size of 400 bp, and Illumina PCR-free libraries were prepared from 700 ng DNA using the KAPA Hyper prep kit and unique dual-indexed adapters (5 µL of a 15 µM stock) according to the supplier's protocol (Roche, Basel, Switzerland). The library concentration and size distribution were assessed on a Bioanalyzer (Agilent Technologies, Santa Clara, CA, USA). Non-depleted WGS libraries were pooled at equimolar concentrations and sequenced on a NovaSeq6000 instrument (Illumina, San Diego, CA, USA) to generate 150-bp paired-end reads.

Design of gRNAs. The gRNA set was designed by JumpCode Genomics (San Diego, CA, USA) against the repetitive regions of the *L. culinaris* CDC Redberry v2.0 reference genome (Ramsay et al.) (<https://knowpulse.usask.ca/genome-assembly/Lcu.2RBY>). The available repeat annotation (transposable elements and tandem repeats) was integrated with the annotation of simple and tandem repeats identified by RepeatMasker v4.0.6 and Tandem Repeat Finder v4.9 using 2 7 7 80 10 50 2000 -d -h parameters to identify intervals for gRNA design (**Supplemental_Table_S1.xlsx**). Adjacent or overlapping intervals were collapsed into single intervals before design. As a first step, all 20 nt sequences with adjacent PAM sites for Cas9 (NGG) were identified in the target intervals. Second, the guides were filtered to exclude secondary structure, high and low GC content, homopolymers, dinucleotide repeats and low *in vitro* cleavage efficiency prediction scores (Azimuth algorithm; (Doench et al. 2016)). Third, the resulting guides were filtered to minimize off-target cleavage in unique regions of the genome by excluding guides that have complementary sites in genomic regions corresponding to genes and open-chromatin regions identified by ATAC-Seq (PRJNA912311) (allowing for up to 3 mismatches). As a final step, and to reduce the number of guides in the set, guides were selected to have no fewer than 25 cleavage sites each and to maintain an inter-guide spacing of at least 500 bp. The final guide set, comprising 569,088 unique guides, was split into 11 pools for the purpose of synthesis. The number of copies of each guide varied and reflected the number of on-target cleavage sites for each guide. DNA oligonucleotides containing the target-specific 20 nt gRNA sequence and invariant single gRNA sequence were synthesized, after which pools of oligonucleotides were amplified by PCR and converted to RNA by *in vitro* transcription. The products of transcription were treated with DNase I and column purified

to generate the final gRNA material. Pools 1–3, 5–8 and 10 contain only gRNAs targeting the nuclear genome. Pools 9 and 11 contain both nuclear and chloroplast genome gRNAs, and pool 4 is the only pool containing gRNAs that target the nuclear, chloroplast and mitochondrial genomes (**Supplemental_Table_S2.xlsx**).

Repeat depletion with JumpCode CRISPRclean. Repetitive regions were depleted using the Cas9 protein and the custom gRNA set described above, according to the Jumpcode CRISPRclean Ribosomal RNA Depletion from Human RNA-Seq Libraries for Illumina Sequencing protocol (Jumpcode Genomics, San Diego, CA, USA) with the following modifications. The input was 10 and 100 ng for the 8-plex and 96-plex libraries, respectively. Depletion was carried out either using all gRNA simultaneously or by splitting the gRNA pools into three groups based on cutting frequency, which were used sequentially in order of increasing cutting frequency (**Supplemental_Table_S2.xlsx**). The sequential depletion strategy was also conducted using the double amounts of gRNAs and Cas9. The reaction volume was 20 µl when using all gRNAs simultaneously or 26 µl for the sequential and double sequential protocols. The quantity of each gRNA pool per reaction is shown in **Supplemental_Table_S2.xlsx** and amounted to 620 ng in the simultaneous and sequential depletion reactions or 1240 ng in the double sequential depletion reaction. The Cas9 enzyme was diluted 1:5 in 1x Cas9 Buffer and 0.0029 µl was used per ng gRNA. The reactions were incubated at 37 °C and libraries were treated in the presence of gRNAs for a total of 1 h (simultaneous depletion protocol) or 3 h (sequential and double sequential depletion protocols, with the sequential gRNA pools added at 1-h intervals). The depleted samples were then size selected using 0.6x volume of AMPure XP Beads (Beckman Coulter, Brea, CA, USA). Libraries were amplified with 10 and 6 PCR cycles for the 8-plex and 96-plex libraries, respectively before final purification with 60µl (0.6x volume) AMPure XP Beads. The concentrations of depleted libraries were measured using the Qubit system and size distributions were assessed on a Tape Station System as described above. Depleted libraries were pooled at equimolar concentrations and sequenced on a NovaSeq6000 instrument to generate 150-bp paired-end reads.

Data analysis and variant calling. Raw read quality was assessed using FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) and the multiplex libraries were demultiplexed using fgbio v1.3.0 DemuxFastqs, assigning fragments by exploiting the unique sample identifier included during first-strand synthesis. The raw reads were then quality filtered and the Illumina sequencing adapters removed using scythe v0.991 and syckle v1.33, respectively. Filtered reads were aligned to the *L. culinaris* v2.0 reference genome using bwa-mem v2.2.1 and the resulting alignments were converted to bam files and sorted using samtools v1.13. PCR-derived duplicates were removed using the GATK MarkDuplicates tool v4.1.7.0 and overlapping portions of the paired-end reads were clipped using the

fgbio v1.3.0 ClipBam tool. The cleaned bam files were used to calculate coverage depth, breadth and fraction of PASS bases (at $\geq 5\times$) using bedtools v2.30.0 genomecov and GATK v3.8 CallableLoci, respectively. The number of reads aligning to the reference genome and to different regions of interest was calculated using samtools v1.13 with option -c to discard reads with a 2308 sam flag in order to consider only the primary alignment, thus omitting repetitive counts of the same multimapping reads. When necessary, sequencing data were normalized to a pre-defined number of input fragments using seqtk sample v1.3.

The coverage variation between the depleted and non-depleted libraries was calculated for target regions with at least one target site using the following formula:

$$Coverage_variation = \log_2 \frac{Mean_cov(depleted)}{Mean_cov(not\ depleted)}$$

In those regions with a mean coverage of 0 in the depleted libraries, the value was set to the lowest calculated value of -12.1325. The coverage variation plot was generated using the ggplot package in R.

Genomic variants were identified using GATK HaplotypeCaller v4.1.7.0 with the parameters “--min-base-quality-score 20 -ERC GVCF”. Individual gVCF files were merged using GATK GenomicsDBImport v4.1.7.0 and the final VCF file was generated using GATK GenotypeGVCFs v4.1.7.0. Variant filtration was achieved using GATK hard filters (<https://gatk.broadinstitute.org/hc/en-us/articles/360037499012?id=3225>).

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

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

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

REFERENCES

- Bellucci E, Mario Aguilar O, Alseekh S, Bett K, Brezeanu C, Cook D, de la Rosa L, Delledonne M, Dostatny DF, Ferreira JJ, et al. 2021. The INCREASE project: Intelligent Collections of food-legume genetic resources for European agrofood systems. *Plant Journal* **108**: 646–660.
- Bennetzen JL, Wang H. 2014. The contributions of transposable elements to the structure, function, and evolution of plant genomes. *Annu Rev Plant Biol* **65**: 505–530.
- Cooke TF, Yee MC, Muzzio M, Sockell A, Bell R, Cornejo OE, Kelley JL, Bailliet G, Bravi CM, Bustamante CD, et al. 2016. GBStools: A Statistical Method for Estimating Allelic Dropout in Reduced Representation Sequencing Data. *PLoS Genet* **12**.
- Davey JW, Hohenlohe PA, Etter PD, Boone JQ, Catchen JM, Blaxter ML. 2011. Genome-wide genetic marker discovery and genotyping using next-generation sequencing. *Nat Rev Genet* **12**: 499–510.
- Deng T, Zhang P, Garrick D, Gao H, Wang L, Zhao F. 2022. Comparison of Genotype Imputation for SNP Array and Low-Coverage Whole-Genome Sequencing Data. *Front Genet* **12**.
- Doench JG, Fusi N, Sullender M, Hegde M, Vaimberg EW, Donovan KF, Smith I, Tothova Z, Wilen C, Orchard R, et al. 2016. Optimized sgRNA design to maximize activity and minimize off-target effects of CRISPR-Cas9. *Nat Biotechnol* **34**: 184–191.
- Feuillet C, Leach JE, Rogers J, Schnable PS, Eversole K. 2011. Crop genome sequencing: Lessons and rationales. *Trends Plant Sci* **16**: 77–88.
- Friel J, Bombarely A, Fornell CD, Luque F, Fernández-Ocaña AM. 2021. Comparative analysis of genotyping by sequencing and whole-genome sequencing methods in diversity studies of olea europaea l. *Plants* **10**.
- Gargiulo R, Kull T, Fay MF. 2021. Effective double-digest RAD sequencing and genotyping despite large genome size. *Mol Ecol Resour* **21**: 1037–1055.
- Gu W, Crawford ED, O'Donovan BD, Wilson MR, Chow ED, Retallack H, DeRisi JL. 2016. Depletion of Abundant Sequences by Hybridization (DASH): Using Cas9 to remove unwanted high-abundance species in sequencing libraries and molecular counting applications. *Genome Biol* **17**.
- Guerra-García A, Gioia T, von Wettberg E, Logozzo G, Papa R, Bitocchi E, Bett KE. 2021. Intelligent Characterization of Lentil Genetic Resources: Evolutionary History, Genetic Diversity of Germplasm, and the Need for Well-Represented Collections. *Curr Protoc* **1**.

- 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 sequencing highlights the role of wild-relative introgression in shaping the adaptive landscape of the wheat genome. *Nat Genet* **51**: 896–904.
- Hodges E, Xuan Z, Balija V, Kramer M, Molla MN, Smith SW, Middle CM, Rodesch MJ, Albert TJ, Hannon GJ, et al. 2007. Genome-wide in situ exon capture for selective resequencing. *Nat Genet* **39**: 1522–1527.
- Ichida H, Abe T. 2019. An improved and robust method to efficiently deplete repetitive elements from complex plant genomes. *Plant Science* **280**: 455–460.
- Jayakodi M, Golicz AA, Kreplak J, Fechete LI, Angra D, Bednář P, Bornhofen E, Zhang H, Boussageon R, Kaur S, et al. The giant diploid faba genome unlocks variation in a global protein crop. <https://doi.org/10.1101/2022.09.23.509015>.
- JumpCode Genomics. 2021. Technology Overview Version 1.2 Harnessing CRISPR to boost NGS sensitivity with CRISPRclean™. https://www.jumpcodegenomics.com/wp-content/uploads/2021/07/jumpcode-technical-overview-20210521_v1-1_F.pdf.
- Kreplak J, Madoui MA, Cápál P, Novák P, Labadie K, Aubert G, Bayer PE, Gali KK, Syme RA, Main D, et al. 2019. A reference genome for pea provides insight into legume genome evolution. *Nat Genet* **51**: 1411–1422.
- Lachance J, Tishkoff SA. 2013. SNP ascertainment bias in population genetic analyses: Why it is important, and how to correct it. *BioEssays* **35**: 780–786.
- Lee S-I, Kim N-S. 2014. Transposable Elements and Genome Size Variations in Plants. *Genomics Inform* **12**: 87.
- Matvienko M, Kozik A, Froenicke L, Lavelle D, Martineau B, Perroud B, Michelmore R. 2013. Consequences of Normalizing Transcriptomic and Genomic Libraries of Plant Genomes Using a Duplex-Specific Nuclease and Tetramethylammonium Chloride. *PLoS One* **8**.
- Miller DFB, Yan PS, Buechlein A, Rodriguez BA, Yilmaz AS, Goel S, Lin H, Collins-Burow B, Rhodes L v., Braun C, et al. 2013. A new method for stranded whole transcriptome RNA-seq. *Methods* **63**: 126–134.
- Montefiori L, Hernandez L, Zhang Z, Gilad Y, Ober C, Crawford G, Nobrega M, Sakabe NJ. 2017. Reducing mitochondrial reads in ATAC-seq using CRISPR/Cas9. *Sci Rep* **7**.
- Müller Paul H, Istanto DD, Heldenbrand J, Hudson ME. 2022. CROPSR: an automated platform for complex genome-wide CRISPR gRNA design and validation. *BMC Bioinformatics* **23**.
- Ogutcen E, Ramsay L, von Wettberg EB, Bett KE. 2018. Capturing variation in Lens (Fabaceae): Development and utility of an exome capture array for lentil. *Appl Plant Sci* **6**.
- Parchman TL, Jahner JP, Uckele KA, Galland LM, Eckert AJ. 2018. RADseq approaches and applications for forest tree genetics. *Tree Genet Genomes* **14**.
- Pavan S, Delvento C, Ricciardi L, Lotti C, Ciani E, D’Agostino N. 2020. Recommendations for Choosing the Genotyping Method and Best Practices for Quality Control in Crop Genome-Wide Association Studies. *Front Genet* **11**.

- Peterson BK, Weber JN, Kay EH, Fisher HS, Hoekstra HE. 2012. Double digest RADseq: An inexpensive method for de novo SNP discovery and genotyping in model and non-model species. *PLoS One* **7**.
- Prezza G, Heckel T, Dietrich S, Homberger C, Westermann AJ, Vogel J. 2020. Improved bacterial RNA-seq by Cas9-based depletion of ribosomal RNA reads. <http://www.rnajournal.org/cgi/doi/10.1261/rna>.
- Ramsay L, Koh CS, Kagale S, Gao D, Kaur S, Haile T, Gela TS, Chen L-A, Cao Z, Konkin DJ, et al. Genomic rearrangements have consequences for introgression breeding as revealed by genome assemblies of wild and cultivated lentil species. <https://doi.org/10.1101/2021.07.23.453237>.
- Ren Q, Wang Y, Lin Y, Zhen Z, Cui Y, Qin S. 2021. The extremely large chloroplast genome of the green alga *Haematococcus pluvialis*: Genome structure, and comparative analysis. *Algal Res* **56**.
- Ricci WA, Lu Z, Ji L, Marand AP, Ethridge CL, Murphy NG, Noshay JM, Galli M, Mejía-Guerra MK, Colomé-Tatché M, et al. 2019. Widespread long-range cis-regulatory elements in the maize genome. *Nat Plants* **5**: 1237–1249.
- Sakamoto W, Takami T. 2018. Chloroplast DNA dynamics: Copy number, quality control and degradation. *Plant Cell Physiol* **59**: 1120–1127.
- Siddique A, Suckow G, Ordoukhanian P, Head S, Homer N, Hernandez A, Brown K, Glick L, Baruch K, Doran P, et al. 2019. RipTide High Throughput NGS Library Prep for Genotyping in Populations. *J Biomol Tech.*; **30**(Suppl):S35-S36.
- Tanaka N, Shenton M, Kawahara Y, Kumagai M, Sakai H, Kanamori H, Yonemaru JI, Fukuoka S, Sugimoto K, Ishimoto M, et al. 2021. Investigation of the Genetic Diversity of a Rice Core Collection of Japanese Landraces using Whole-Genome Sequencing. *Plant Cell Physiol* **61**: 2087–2096.
- Tian F, Yang DC, Meng YQ, Jin J, Gao G. 2020. PlantRegMap: Charting functional regulatory maps in plants. *Nucleic Acids Res* **48**: D1104–D1113.
- Truong HT, Ramos AM, Yalcin F, de Ruiter M, van der Poel HJA, Huvenaars KHJ, Hogers RCJ, van Enckevort LJG, Janssen A, van Orsouw NJ, et al. 2012. Sequence-based genotyping for marker discovery and co-dominant scoring in germplasm and populations. *PLoS One* **7**.
- Wang J, Sun G, Ren X, Li C, Liu L, Wang Q, Du B, Sun D. 2016. QTL underlying some agronomic traits in barley detected by SNP markers. *BMC Genet* **17**.
- Wang P, Xiong Y, Gong R, Yang Y, Fan K, Yu S. 2019. A key variant in the cis-regulatory element of flowering gene *Ghd8* associated with cold tolerance in rice. *Sci Rep* **9**.
- Yan H, Haak DC, Li S, Huang L, Bombarely A. 2022. Exploring transposable element-based markers to identify allelic variations underlying agronomic traits in rice. *Plant Commun* **3**.
- Zan Y, Payen T, Lillie M, Honaker CF, Siegel PB, Carlborg Ö. 2019. Genotyping by low-coverage whole-genome sequencing in intercross pedigrees from outbred founders: A cost-efficient approach. *Genetics Selection Evolution* **51**.

Zhao W, He X, Hoadley KA, Parker JS, Hayes DN, Perou CM. 2014. *Comparison of RNA-Seq by poly (A) capture, ribosomal RNA depletion, and DNA microarray for expression profiling*.
<http://www.biomedcentral.com/1471-2164/15/419>.

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