- 1 Review Article
- 2 Proposed EU NGT legislation in light of plant genetic variation
- 3 Alan H. Schulman<sup>1,2\*</sup>, Frank Hartung<sup>3†</sup>, Marinus J.M. Smulders<sup>4†</sup>, Jens F. Sundström<sup>5†</sup>, Ralf Wilhelm<sup>3†</sup>,
- 4 Odd Arne Rognli<sup>6†</sup>, and Karin Metzlaff<sup>7</sup>
- <sup>5</sup> <sup>1</sup>HiLife Institute of Biotechnology and Viikki Plant Science Centre (ViPS), University of Helsinki,
- 6 Finland; alan.schulman@helsinki.fi
- 7 <sup>2</sup>Production Systems, Natural Resources Institute Finland (LUKE), Helsinki, Finland;
- 8 alan.schulman@luke.fi
- 9 <sup>3</sup>Julius Kuehn Institute (JKI) Federal Research Centre for Cultivated Plants, Institute for Biosafety in
- Plant Biotechnology, Quedlinburg, Saxony-Anhalt, Germany; frank.hartung@julius-kuehn.de;
   ralf.wilhelm@julius-kuehn.de
- <sup>4</sup>Plant Breeding, Wageningen University & Research, Wageningen, The Netherlands;
- 13 rene.smulders@wur.nl
- <sup>5</sup>Department of Plant Biology, Swedish University of Agricultural Science, The Linnean Centre for
- 15 Plant Biology, Box 7080, SE-75007 Uppsala, Sweden; jens.sundstrom@slu.se
- 16 <sup>6</sup>Faculty of Biosciences, Department of Plant Sciences, Norwegian University of Life Sciences (NMBU),
- 17 Ås, Norway; odd-arne.rognli@nmbu.no
- 18 <sup>7</sup>European Plant Science Organisation (EPSO), Brussels, Belgium; karin.metzlaff@epsomail.org
- 19 \*Correspondence: (Tel +358 407682242; email alan.schulman@helsinki.fi)
- 20 <sup>†</sup>These authors contributed equally to this work
- 21
- 22 Key words: gene editing; new genomic techniques (NGT); new breeding techniques (NBT); plant
- 23 genetic diversity; mutagenesis; CRISPR/Cas9; plant breeding

## 25 Summary

26 The European Commission (EC) proposal for New Genomic Techniques (NGTs) of July 2023 specifies 27 that Category 1 NGT (NGT1) plants, which are considered equivalent to conventional plants, i.e. 28 those obtainable by conventional plant breeding or mutagenesis, may differ from the recipient or 29 parental plant by no more than 20 insertions, which cannot be longer than 20 bp; deletions can be 30 no more than 20 but of any size. Here, we examine the proposed 20/20 NGT1 limit against the 31 background of the theoretical considerations and older data used to frame it and in light of recent 32 data from highly contiguous long-read assemblies for reference genomes and pangenomes. We find 33 that current genomic data indicate that natural variation in germplasm used by breeders is much 34 greater than earlier understood and that both conventional breeding and mutagenesis can introduce 35 genomic changes that are both more extensive in size and more frequent than the NGT Category 1 36 "20 insertions of maximum 20 bp" limit would allow. Furthermore, natural variation also scales with 37 genome size and complexity, a factor not considered in the EC proposal. We conclude that the 38 proposed cutoffs under which an NGT plant is considered equivalent to conventional plants do not 39 align with what is observed in nature, conventional breeding, and mutagenesis. Updating the 20/20 rule to broader limits would facilitate breeding for climate resilience, farming sustainability, and 40 41 nutritional security, while ensuring that NGT1 plants are equivalent to conventional ones.

## 42 Introduction

- 43 Annex 1 to the European Commission (EC) proposal (2023/0226) on New Genomic Techniques (NGTs)
- 44 specifies the number and types of changes that would be regarded as equivalent to the variation
- 45 found in conventional plants, i.e. those that "could have been produced through conventional plant
- 46 breeding or classical mutagenesis" (as stated in COM(2023) 411 final). Plants within these bounds are
- 47 regarded as NGT Category 1 (NGT1). With the rapid progress in genomic sequencing methods, our
- 48 understanding of plant genomic variation is improving quickly in parallel. Here, we consider what is
- 49 known about variation between plant genomes and to what extent the proposed NGT1 category
- reflects that. The EC also defines a Category 2 (NGT2), which includes all NGT modifications that
- 51 exceed the limits specified in NGT1. The traceability and labelling regime currently in place for GMOs
- 52 under Directive 2001/18/EC would be adapted to NGT2, likely reducing commercial interest in
- 53 introducing NGT2 products into the market. As a counterbalance, the EC proposes regulatory
- 54 incentives under Annex III for use of NGT2 for various sustainability traits, including disease
- resistance, abiotic stress tolerance, and nutrient and water use efficiency; herbicide tolerance is
- 56 explicitly excepted. Given that the line between NGT1 and NGT2 will crucially affect the use and
- 57 commercialization of NGTs, we will restrict our focus to the EC proposal and to the relevant research
- 58 underlying it (Figure 1).

# 59 <u>From Annex I:</u> 60 "An NGT plant is

- 60 "An NGT plant is considered equivalent to conventional plants when it differs from the recipient/parental plant by no more than [20] genetic modifications of the types referred to in 61 points 1 to 5, in predictable DNA sequences. A predictable DNA sequence is any DNA sequence 62 that shares sequence similarity with the targeted site." 63 (1) Substitution or insertion of no more than (20) nucleotides ; 64 65 (2) Deletion of any number of nucleotides; 66 (3) On the condition that the genetic modification does not result in an intragenic plant: 67 (a) Targeted insertion of a contiguous DNA sequence existing in the breeder's 68 gene pool; 69 (b) Targeted substitution of an endogenous DNA sequence with a contiguous DNA 70 sequence existing in the breeder's gene pool; 71 (4) Targeted inversion of a sequence of any number of nucleotides; 72 (5) Any other targeted modification of any size, on the condition that the resulting DNA
- requerces already occur (possibly with modifications as accepted under points 1
- 74 and/or 2) in a species from the breeders' gene pool.

- 76 Regulation of the European Parliament and of the Council on plants obtained by certain new genomic
- techniques and their food and feed, and amending Regulation (EU) 2017/625," of 5 July 2023.

<sup>75</sup> Figure 1 Excerpt from Annex I, "Criteria of equivalence of NGT plants to conventional plants," to "Proposal for a

Annex I sets a very specific standard for insertions within the NGT1 category. This can be interpreted 78 79 as being consistent with the original 2001/18/EC legislation on the deliberate release into the 80 environment of genetically modified organisms, where Article 2(2) specifies a GMO as one in which 81 the genetic material has been altered in a way that "does not occur naturally by mating and/or 82 natural recombination." Likewise, the proposed restriction on the number of changes appears to 83 respond to the 2018 ECJ Curia judgement (ECLI:EU:C:2018:583), which states according to the 84 referring court that, "the new techniques of mutagenesis allows the production of modifications ... at 85 a rate out of all proportion to the modifications likely to occur naturally or randomly...", implying a 86 resulting safety risk.

87 Standards of "naturalness" and "conventional" beg the question of what is found in nature. In the 88 European Commission's document 14204/23 (Commission Services, 2023), "Regulation on new 89 genomic techniques (NGT) – Technical paper on the rationale for the equivalence criteria in Annex I", 90 the criteria are based on a literature analysis of 90 scientific, peer-reviewed original studies, which 91 are cited in the annex of that document. The cited EFSA study on site-directed mutagenesis, 92 however, is from 2012 (EFSA GMO Panel, 2012), which is well before the advent of the current state 93 of the art. The EFSA risk assessment studies in 2020 (EFSA GMO Panel et al., 2020) and 2022 (EFSA 94 GMO Panel et al., 2022) did not revisit the state of knowledge of genome structural variation, either 95 natural or that induced by conventional mutagenesis. Virtually all of the 90 papers that support the 96 proposed standards were based on research from before long-read sequencing. This recently available approach has greatly increased the contiguity and completeness of genome assemblies -97 98 akin to reproduction of manuscripts without missing punctuations or words, or misplaced sentences 99 and paragraphs – and has thereby improved the detection of insertions and deletions ("indels", when 100 taken together), chromosomal rearrangements, and both presence-absence and copy-number 101 variations in gene families.

102 Critically, the true dynamic nature of the genome could not be resolved by the methods available 103 before 2012, and in fact not before the advent of the PacBio HiFi long-read sequencing method in 104 2022, complemented today with e.g. Nanopore technology. Indeed, the 14204/23 document 105 anticipates its own obsolescence, stating that "...improvement of detection methods (i.e. long-read 106 sequencing) has started to unveil higher rates than previously estimated" for genomic changes larger 107 than single-nucleotide polymorphisms. The PacBio HiFi method, for example, produces read lengths 108 of 10 to 25 kb at accuracies of 99.5% or higher (Hon et al., 2020). Long reads span retrotransposons, 109 which can be 5 to 15 kb and even 30 kb in length, whose high copy number and conservation 110 interfere with unique assignment and assembly of reads into genomes (Jayakodi et al., 2023). To aid 111 in correct orientation and assembly of long reads into chromosome-scale scaffolds that bridge the

4

- remaining sequence gaps, the Chromosome Conformation Capture (Hi-C) sequencing method has
- 113 come into use (Himmelbach et al., 2018). The Hi-C method ligates together segments of DNA that are
- on proximal DNA strands in chromatin, thereby helping to establish the linear order of long sequence
- reads in scaffolds and also enabling phasing of haplotypes (Sun *et al.*, 2025). Hi-C can be
- supplemented as well by optical mapping (Paajanen et al., 2019; Jayakodi et al., 2023) to further
- 117 improve long-read assemblies.
- 118 Insertion and deletion sizes in plant genomes vs the 20 bp limit
- A key restriction in 14204/23 is the 20 bp insertion limit for NGT1. A likely rationale for the limit is to
- 120 distinguish short, random repair-type insertions from long insertions that can be identified as unique
- 121 or specific genomic constituents, i.e., equivalent to cisgenes. There are two components to this
- rationale: first, the assumption that natural "random" insertions are short; second, that insertions
- 123 longer than 20 bp can be uniquely identified as pre-existing in the genome (i.e., in the "breeders"
- gene pool"). As is stated in the report, "Insertions of more random sequences are typically of a length
- 125 of less than ten nucleotides but have been observed to extend to approximately fifty nucleotides" and
- 126 that "...a threshold of twenty nucleotides in criterion 1 for substitutions and insertions was set since it
- 127 fits with the sizes observed in the scientific analysis."
- 128 The first question one can raise is: What is the actual size distribution of spontaneous insertions and
- deletions in conventional plants, compared with the 20 bp limit of NGT1? The answer is that recent
- advances in genome sequencing show that natural variations extend from 1 to 1 million bp. Although
- 131 the proposed regulations distinguish between insertions and deletions, in practice it is seldom
- possible to determine the initial state in accessions of cultivars or wild materials, i.e. whether it was a
- 133 spontaneous insertion or deletion that occurred to distinguish the versions of a sequence. Hence,
- the term "indel" is used collectively for insertions and deletions; very often a particular "complex"
- 135 indel will contain a combination of both. While the Commission draft proposal distinguishes between
- the legal status of insertions (of fixed number and size for NGT1) and deletions (of any size or
- 137 number), the data from the actual natural world does not support this distinction (Figures 1—3). For
- 138 indels found in sequenced genomes, even early (2013) data from long genome assemblies showed
- no bimodal distribution expected for short, random insertions and longer gene-like (or transposon-
- 140 like) insertions within eukaryotes overall (Figure 2).



Figure 2 For each indel size class (x-axis), the number of simple (total = 901) and complex (total = 3,806) indels are indicated by the red and blue bars, respectively. 501 indels (10 simple indels and 491 complex indels) longer than 50 amino acid residues are not shown. Simple indels" occur in only two states, present or absent, and are potentially the result of a single indel event, while "complex indels" occur in two or more states and represent multiple compounded indel events. Modified from Ajawatanawong and Baldauf (2013).

147

148 "Complex" indels (Figure 2), which are the likely result of multiple, nested events over time, show no

sharp decline and, as compound events, are *ipso facto* not structurally equivalent to cisgenes. In rice,

indel markers varied from 3 to 39 bp, with 88.2% 6—25 bp,  $6.2\% \le 5$  bp, and 5.6% were  $\ge 26$  bp

151 (Zeng *et al.*, 2013). Work from 2015 in soybean (Figure 3) with the older short-read technologies

152 indicated a rapid drop-off in indel length, consistent with the earlier EFSA studies. Nevertheless,

analyses of individual gene families (e.g., *RPB2* in barley; Sun *et al.* (2009)), where alignments were

154 carefully constructed for the genes), indicated that indels of 20 – 100 bp are quite common; MITE

155 transposons, which are abundant, are 90 – 100 bp.

156



Figure 3 Frequency of length distribution of indels between soybean cultivars JS-335 and UPSM-534. From
Yadav *et al.* (2015).

160 Critically, it is has become clear that the apparent indel length distribution can be influenced by the 161 limits of alignment and assembly of short-read ("Illumina") sequences; a better picture is now 162 emerging from long-read (PacBio HiFi and Nanopore) sequencing approaches as anticipated but not yet documented by research by EFSA in the 14204/23 technical paper. A striking example is the 163 164 distribution of indels and presence-absence variations (PAVs) between two well-assembled barley 165 cultivars (Figure 4). Another example was published in 2024 for lentil, Lens culinaris (Shivaprasad et 166 al., 2024). These researchers compared the genomes of a lentil parental line with recombinant 167 inbred bulks, finding almost 735 000 indels, of which almost 16 000 were longer than 20 bp, 3600 168 greater than 40 bp, and 1200 greater than 50 bp.







173 Distribution of indels generated by break repair following intentional mutagenesis

174 The question of indel size distribution is relevant because the criteria of naturalness and pre-2001 175 methodology are used to set the outer limits for GE acceptable as "conventional-like," i.e. NGT1. 176 Conventional breeding methods – including mutagenesis methods taken into use pre-2001 – are not 177 subject to the 2001/18 regulatory regime and hence are worth comparing with the outcomes of NGT 178 methods, which are subject to 2001/18 and considered in 14204/23. Ion-beam mutagenesis is one 179 frequently used method (Guo et al., 2024). A recent study in Arabidopsis demonstrated that 180 insertions generated by repair of the double-strand DNA breaks induced by ion-beam irradiation of 181 seedlings ranged from less than 5 bp to over 100 bp, with an average of ~ 12 bp (Kitamura et al., 182 2024). In contrast, data shows that when the early CRISPR/Cas9 method (SDN-1), which causes 183 double-strand breaks, is used to knock-out gene function, the breaks are precisely repaired 36-41% 184 of the time, the remainder not (Ben-Tov et al., 2024). In a study of 361 CRISPR/Cas9-mutated plants, 185 the imprecise break repairs were predominantly short insertions or deletions; 87% of the induced 186 indels were smaller than 10 bp (Zhang et al., 2020). Insertions comprised 30% of the total; 73% of 187 the insertions were only 1 bp in length, 2% were 2-50 bp, and 6% > 50 bp.

188 In many cases, it will be necessary to introduce changes at the mutation site that preserve gene

189 function rather than knocking it out by a deletion or by the repair process of the cell that can

190 generate small random insertions. The currently most popular targeted (NGT) mutagenesis method,

191 CRISPR/*Cas9*, generates distinctly smaller break-repair insertions (~1—10 bp) when used for

192 knockouts than does either conventional mutagenesis or natural processes. Hence, while the Curia

193 judgement of 2018 viewed the genetic changed wrought by new mutagenic techniques as far in

excess of those occurring naturally or by earlier-established methods, the available data shows that

the opposite is the case: NGT methods are therefore considerably gentler in their genomic impact

196 than traditional breeding approaches, whether crossing or random mutagenesis.

197 Minimum length needed to specify a unique sequence in the genome

To address the need for practical monitoring under NGT regulatory regimes, an alternative approach for defining a maximum insertion length acceptable as NGT1 is that it should be below the minimum identifiable unique sequence in a genome, hence it should be one that could result from a random

201 process. Report 14204/23 posits that, "...when considering genome diversity, the JRC calculated that

202 the theoretical probability that a random sequence is unique in the genome of various crops boils

203 *down to a consistent relatively narrow size range between 19 and 21 bases."* As justification for this

- claim, it cites Broothaerts *et al.* (2021). However, this publication (section 4.4, p. 20), has no
   explanation given for the claim; only undescribed and unpublished results from rice are cited.
- 206 One conceivable explanation is that 20 bp length is based on a mathematical calculation. Assuming
- 207 that the four nucleotides (A, C, G, T) occur at equal frequency, the likelihood of occurrence of any
- arbitrary nucleotide sequence of length bp in a genome of random nucleotides is 1/4<sup>bp</sup>. Hence, a
- 209 19mer would have a frequency of 3.6 x 10<sup>-12</sup> bp; a 20mer, 0.9 x 10<sup>-12</sup> bp; a 21mer, 2.3 x 10<sup>-13</sup>. One of
- 210 the largest crop genomes known is that of faba bean (*Vicia faba*), where the basic set of
- chromosomes (monoploid genome) comprises 13 x10<sup>12</sup> bp. An arbitrary 19mer would be expected,
- based on fully random sequence, to occur by chance three or four times in *V. faba* monoploid
- 213 genome (seven times in the diploid, i.e, in all cells except for pollen and the egg cell), whereas a
- 214 20mer would be found only once in the monoploid and two or three times in the diploid; a 21mer,
- 215 once or less.
- 216 Critically, this is an inaccurate estimate of the true frequency of oligomers, and therefore uniqueness,
- 217 in a likely target crop for NGT1. First, the four nucleotides do not occur at the same frequency
- 218 (usually, GC < AT) and plant genome sequences are far from random. This is due both to the
- 219 functional importance of sequence information in genes (both regulatory and coding regions) and
- especially to the high percentage, even 80%, of large genomes represented by relatively few
- abundant retrotransposon families that comprise highly similar sequences. Analyses of the length
- required for true unique representation have been made (Figure 5); for single-occurrence
- frequencies, sequences must be ~400 bp even in compact crop genomes such as rice, although 100
- bp (not 20 bp!) is sufficient for the small genome of Arabidopsis.



Figure 5 The k-mer uniqueness ratio for some assembled plant genomes as a function of k. The uniqueness ratio is the ratio of *k*-mers occurring exactly once relative to all *k*-mers in the set. It is computed for every *k* between 10 and 500. Extrapolating beyond the tested *k*-mer interval, it appears as though poplar, rice, and grape approach unity at a much slower rate than Arabidopsis. Source: Kurtz *et al.* (2008).

229 Moreover, at least currently, most targets for gene editing are protein-coding sequences, which are 230 highly non-random. Eukaryotic proteins are encoded by only 64 specific triplets (codons), which form

the genetic code for the amino acids, with some of these being much preferred for particular amino

acids (De Amicis and Marchetti, 2000). Furthermore, some amino acids are over-represented in the

encoded set of cellular proteins, the proteome. Hence the frequencies of 20mers, each representing

234 ~7 amino acids and among the likely gene editing targets, are much higher than expected for random

sequences of that length. Certainly, if we consider the contiguous protein-coding segments of a gene

236 (exons), random insertions and deletions would equally likely destroy the protein's function unless,

at least, they were precisely phased with the reading frame of the gene. Hence, 20 bp, the insertion

limit under NGT1 is insufficient to specify a unique, non-random insertion in typical plant genomes;

sequences at least 20-fold longer are still within the range of statistically random variation in plants.

240 Alternative approaches to uniqueness for insertions under NGT1

241 From the considerations above, random mutations found in nature provide no obvious limit to

insertion size based on a naturalness criterion. If the uniqueness argument is used, the data (Figure

5) would indicate a limit of at least 400 bp would be needed under the NGT1 standard. An alternative

approach would be to choose as the limit the largest insertion that would not contain the coding

sequence of a full-length protein, in order to maintain a distinction between NGT1 and a gene

insertion, i.e., achievable with transgenesis or cisgenesis.

247 Of the conventional cellular proteins, one-finger (Dof) proteins are plant-specific zinc finger proteins

and typically contain 200 to 400 amino acids (Waschburger *et al.*, 2024), equivalent to 600—1200 bp.

Plant haemoglobins are still smaller, ~150 amino acids, encoded by 450 bp (Becana *et al.*, 2020).

250 "Miniproteins," recently discovered, are the smallest proteins to be found in plants (Gruber et al.,

251 2008). They are generally only 50 to 60 amino acids long, hence equivalent to 150 bp, but despite

their small size, they can play important regulatory functions (Molesini *et al.*, 2012). For example, the

253 cyclotides, a special class of miniproteins found in the family *Violaceae*, have antimicrobial and

antifungal properties (Kim *et al.*, 2023; Lian *et al.*, 2024). Among mammalian proteins, insulin is

exceptionally small, the mature form comprising two chains of 21 and 30 amino acids respectively;

another example, the bioactive thymosin alpha 1 peptide, is 28 amino acids long (Tao *et al.*, 2023).

257 Given that proteins of less than 50 or 60 amino acids in length are however unlikely to fold into an

active form (Linsky *et al.*, 2022), 150 bp (given 3 bp per amino acid) is a reasonable insertion size for

distinguishing protein-coding sequences. This could serve as the maximum insertion size qualifyingas NGT1.

261 It is worth noting that a limit of 50 amino acids or 150 bp generally distinguishes functional proteins, 262 but that short peptides, even less than ten amino acids, if expressed, may have functionality, e.g. 263 through their binding to enzymes in a cell. Coding sequences for short peptides may be generated 264 naturally through point mutations or indels, such as those resulting from double-strand break repair 265 processes, which are discussed above, and of course through proteolytic digestion. However, unless 266 they are near an active promotor, are contained within an mRNA that will be translated, are 267 produced in significant quantity, and have biological function, they are of no consequence. Possible 268 formation of such peptides in GMO events, for example, is routinely checked against toxin databases. 269 Insertion and deletion numbers in plant genomes vs the 20-insertion limit

The idea in the proposed legislation is that what is achieved by NGT1 should be equivalent to, and no

271 more than, what can be reached through conventional breeding (which includes radiation and

chemical mutagenesis). The 14204/23 report states that, in the literature, "the total number of

273 genetic modifications in individual viable plants ranged from thirty to one hundred. The mutation

274 frequency after using random mutagenesis was higher compared to natural mutation rates. It

275 remained nevertheless below the total number of accumulated single nucleotide polymorphisms

276 naturally occurring between different cultivars."

277 This concept raises the question of what current data show for the number of indels found between 278 cultivars, landraces, and wild accessions. First, it is important to note that comparisons are based on 279 sequence assemblies representing, for a given cultivar, landrace, or wild line, either a single 280 individual or a consensus from several individuals. Heterozygosity within the accession or cultivar is 281 filtered from the published sequence. However, very recent intra-varietal long-read sequencing has 282 been made and phased into the two haplotypes of the clonally propagated "Fuji" apple (Cai et al., 283 2024), allowing the discovery of 68,965 somatic SNPs across 74 individuals, or 932 per each. Intra-284 individual mutation rates vary greatly by tissue, by propagation method (clonal vs. sexual), and by life 285 cycle (perennial vs. annual), ranging from 0.08—15.78 x 10<sup>-9</sup> per bp per year, the highest rate being 286 seen in wild strawberry (Fragaria vesca) stems (Wang et al., 2019). This rate corresponds to 6 287 changes per diploid genome in each cell per year in strawberry plants clonally propagated by 288 runners. In long-lived individuals, these changes accumulate; the same study found up to 19 289 inherited mutations (mean 11) per individual peach (Prunus persica) on one tree, which would be 290 close to the limit permitted for NGT1 insertions under the proposed legislation.

291 Coming back to consensus sequences for plant lines, just as for indel size, current long-read

- assemblies provide a perspective on indel number that was generally unavailable before 2022. Taking
- barley as an example, the recent barley pan-genome (Jayakodi *et al.*, 2024), comprising long-read
- sequence assemblies of 76 wild and domesticated genomes and short-read sequence data of 1,315
- 295 genotypes, contains a total of 155 million SNPs and 9 million indels in 315 elite cultivars, or 493,837
- 296 SNPs and 28,983 indels per accession. Moreover, the extensive mutation breeding used for barley in
- the 1960s has left a legacy of abundant inversion polymorphisms in current germplasm that confer
- various selective advantages: among 69 barley genotypes (67 domesticated and 2 wild accessions) a
- total of 42 inversions were found that ranged from 4 to 141 Mb in size (mean 23.9 Mb). An
- independent, very complete survey of the barley gene pool (Weisweiler *et al.*, 2022) shows ~100,000
- indels (lengths of 2—49 bp examined) in genic (exon + intron) regions among 23 inbred lines. Clusters
- of structural variants (SV) present per inbred ranged from less than 40,000 to more than 80,000.
- The high level of SVs and indels is not unique to barley. Regarding rice, *Oryza sativa* ssp. *javanica* is a
- large-grain landrace. A recent study (Long *et al.*, 2022) found from 164,018 to 211,135 indels and
- 305 3,313 to 4,959 longer SVs in javanica compared to the commonly cultivated japonica or indica
- 306 subspecies. In grapevine, Di Genova *et al.* (2014) identified 623,003 indels of 1 bp to 46 kb, of which
- 307 5981 were exon indels and 172,385 intron indels. In wheat, when the Chinese Spring reference
- 308 genome was compared to other bread wheat accessions, some 36,904 frameshift indels where found
- that may impact protein function (Montenegro *et al.*, 2017).
- 310 The high level of variations found by genome sequencing of crop cultivars and landraces has direct
- 311 practical implications. Conventional breeding involves crossing of elite cultivars with each other as
- 312 well as introgression of genetic material from landraces and wild relatives. Crosses will introduce the
- full complement of variations, including SNPs, indels and other SVs, present on one haploid set of
- chromosomes, amounting to 40 to 80 thousand in the case of barley. The incorporation of massive
- 315 numbers of genic and regulatory variations by crossing necessitates extensive back-crossing to the
- elite parental cultivar in most breeding programs, a process slowed by the "linkage drag" of
- unwanted variants flanking a desired introduced allele (Chitwood-Brown *et al.*, 2021; Deblieck *et al.*,
- 318 2022).
- 319 Not only conventional crossing, but also conventional random mutagenesis (not regulated under
- 320 2001/18), introduces large numbers of changes, the type and frequency depending on the
- 321 mutagenesis agent and dosage. Mutation frequencies from the commonly used chemical mutagen
- ethyl methane sulphonate (EMS) can be 1.5—4.1 x 10<sup>-6</sup>, corresponding to 7500 to 20,000 "off-target"
- mutations in the haploid barley genome of a mutagenized line (Jiang *et al.*, 2022). It is precisely the
- 324 messiness of conventional mutagenesis compared with the clean introduction of edited alleles by

- NBT that attracts breeders to gene editing (Yang *et al.*, 2023). Frequencies of off-target mutations
- induced by CRISPR-*Cas9* are very low, generally less than 5% in likely (i.e., almost identical) off-target
- 327 sites (Slaman *et al.*, 2023), which would correspond to frequencies on the order 1 x 10<sup>-9</sup> in the barley
- 328 chemical mutagenesis example above. The few off-target mutations would be segregated away
- 329 rapidly by onward breeding.
- 330 The studies described above, taken together, show that intra-plant, inter-individual, and inter-line
- indel numbers, both spontaneously occurring and obtained via conventional mutagenesis, are
- 332 generally well in excess, even by a thousand-fold, over permissible insertion numbers under the
- 333 proposed standard for NGT1. Even if we assume that half of the indels are insertions (restricted
- under NGT1) and the other half are deletions (unrestricted), targeted mutagenesis such as by
- 335 CRISPR/Cas9 or similar methods will not plausibly approach the amount of insertion-generated
- 336 variations seen in the breeders' pool.
- 337 Practical consequences of the maximum 20 permitted NGT1 insertions
- 338 While any number and size of deletions is permitted for NGT1, in cases where insertions are used to
- edit multiple members of gene families, the question of gene family size versus natural variation
- 340 within becomes relevant. Gene families in plants range from single-copy to hundreds of members.
- The many ongoing pan-genome projects in plants, in which high-quality genome assemblies for
- 342 multiple accessions can be analysed, have revealed large variations in many gene family sizes both
- 343 within and between species (Niu *et al.*, 2024). These together with structural variations, indels and
- 344 SNPs and would thereby challenge the proposed 20/20 rule because the number of targets for
- editing under NGT1, as well as their initial state, may vary from cultivar to cultivar.
- 346 The NLR genes are a good example of an important NGT target limited by the 20-insertion rule.
- 347 Plant genomes typically contain hundreds of nucleotide-binding site leucine-rich repeat (NLR) genes,
- 348 which are the largest family of plant disease resistance genes. The number of NLR genes per genome
- vary from 149 in Arabidopsis to ~3400 in bread wheat (Tong *et al.*, 2022). The NRL genes in
- Arabidopsis (Mondragon-Palomino et al., 2017), wheat (Hao et al., 2023), and soybean (Liu et al.,
- 2024), have been shown to have evolved and diversified through recombination and accumulation of
- 352 SNPs and indels, with changes displaying association with disease resistance. Resistance genes are
- often "stacked", as described below, and modified rather than knocked out. Hence, the need to edit
- 354 more than 20 by insertion approaches, especially to provide resistance against several pathogens,
- 355 can likely easily arise.
- 356 The alpha-gliadin genes as an example of the impact of limitations arising from 20-insertion rule

357 The genes of alpha-gliadin family of storage proteins in wheat are part of the very dynamic *Gli-2* loci. 358 The alpha gliadins are known for their importance in breadmaking as well as for their role in 359 triggering celiac disease (CD). A combination of long-read sequencing and optical mapping was used 360 to assemble the loci (Huo et al., 2018). Three loci are found in each homoeologous set of 361 chromosomes (A, B, D) of the hexaploidy bread wheat genome, in total nine loci, hence illustrating 362 the importance of using the monoploid chromosome set as the standard for the number of 363 permitted changes in plant genomes and increasing it by the ploidy level (see discussion below). Huo 364 et al. (Huo *et al.*, 2018) identified a total of 47  $\alpha$ -gliadin genes in bread wheat, with only 26 encoding 365 intact full-length protein products. Altogether 21 of the 47 were pseudogenes, 13 due to SNPs, 4 to deletions, others to rearrangements. Three contained TE insertions, premature stop codons, and 366 367 frameshift indels. However, a 20mer associated with CD epitopes is present in 2161 copies at 93— 368 100% identity in the alpha gliadin genes within the Chinese Spring genome (Schulman, unpublished). 369 Others have attempted to analyse the relative abundance of CD types (Marin-Sanz et al., 2023). An 370 in-depth analysis of transcription and protein accumulation in the bread wheat Chinese cultivar 371 Xiaoyan 81 (Wang et al., 2017) found that 52 full-length gliadin genes were transcribed, 42 of these 372 encoded proteins, 38 gliadins accumulated in mature grains, 10 did not carry any CD epitope, eight 373 had one or two epitopes in their proteins, and 20 contained more than three epitopes in their 374 proteins; of the 28 gliadins with CD epitopes, a total of 202 epitopes in the proteins were present at 375 100% match. Making the alpha-gliadins safe for CD patients by using NGT for all 28 CD-epitope-376 containing alpha-gliadin genes to alter all 202 CD epitopes would not be acceptable under the 20/20 377 rule within the current EC proposal. Removal through large deletions of the tandemly organised 378 genes (Jouanin et al., 2019), while permitted, is possible but not practical for all gliadin families if one 379 wants to maintain baking quality (Jouanin et al., 2020). 380 As a further example, receptor-like kinases (RLKs), which are critical for biotic and abiotic stress

response, and therefore likely NGT targets, are found in 100s to 1000s copies depending on the plant

species and have undergone a great degree of recombination and variation (Yan *et al.*, 2023).

383 Another example of a large gene family in plants is that of cytochrome P450 (CYP450), which

includes 100s of members in most plant genomes (Zhang et al., 2023). A subgroup of CYP450, CYP71,

385 which is connected to insect resistance, senescence, and yield-related traits, was studied in rice. In

rice, 105 *OsCYP71* genes were found, of which 36 pairs were involved in gene duplication (in essence,

387 large SVs); major indels of 20 bp affecting 20% of the varieties' promoter structures and thereby

expression patterns and trait QTLs were found. In these sorts of cases, the natural variation would

need to be confirmed in the edited and non-edited versions to confirm that the editing per se did not

390 generate more than 20 changes for NGT1 status.

#### 391 Impact on polyploid crops

392 Beyond variation in gene family number in the basic set of chromosomes, many plant species are not 393 diploid (two sets of chromosomes) but rather tetraploid (four sets), hexaploid (six), octoploid (eight), 394 or even higher. This means that gene family numbers likewise may double, triple, guadruple, or be of 395 higher multiples, as described above for CD epitopes in wheat, complicating editing within a fixed, 396 low limit of insertions under NGT1. For example, pasta (durum) wheat is tetraploid, as is potato, 397 while bread (common) wheat is hexaploid, and cultivated strawberry is octoploid, as is sugar cane. 398 Without adjustments for ploidy, the current 20/20 limits for NGT1 would therefore be far more 399 restrictive for bread wheat than for pasta wheat, and both more than for einkorn wheat, which is a 400 diploid as is barley. Cultivated roses (Rosa hybrida) can be either diploid, triploid, or tetraploid 401 (Harmon et al., 2023) but would be permitted the same maximum 20 insertions under NGT1. Clearly 402 a more rational approach is needed.

403 Gene stacking versus NGT insertion number

404 An important goal in plant breeding is to "stack" or combine multiple beneficial traits into the same 405 plant line, e.g. to improve an already commercially successful variety. This is to meet the widespread 406 goal and urgent need for several classes of phenotypes: simultaneous resistance to multiple plant 407 diseases; robust resistance to individual pathogens through combined use of different independently 408 acting genes; both abiotic (e.g. drought) stress tolerance and disease resistance in a crop plant; both 409 healthy crop plants and a harvest with human-health-promoting qualities (e.g. CD-safety). In some 410 cases, even a single trait, for example CD-epitope-free gluten protein in wheat, requires the stepwise 411 stacking of alleles. This is possible but slow by conventional breeding, requiring support by marker-412 assisted selection (MAS) and epitope immuno-assays. At each stage, the properties of the gliadins for 413 baking quality would need to be preserved and tested. In fact, even use of GMO approaches to 414 introduce multiple genes in parallel is technically highly challenging (Halpin, 2005). In practice, also 415 transgenes therefore have been stacked through conventional crosses (Li et al., 2023), with at most 416 seven genes currently stacked (https://www.isaaa.org/gmapprovaldatabase/eventslist/), which is a 417 maize line providing herbicide tolerance, multiple insect resistance, a modified alpha amylase, and 418 altered mannose metabolism. The practical limitations to gene stacking raise several important 419 questions: First, should the current technical limits of older conventional approaches serve as the 420 basis for limiting NGT target numbers? Second, if so, can this limit be justified by some risk specific to 421 gene stacking and not merely the sum of the individual risks? Conventional gene stacking is a moving 422 target, as both biochemical phenotyping and marker-assisted selection improves. Moreover, no 423 restrictions are imposed on stacked-gene conventional cultivars; rather, they command a premium 424 price and are welcomed in the marketplace.

#### 425 Conclusions and Future Prospects

426 We find that current genomic data indicate that natural variation in the germplasm used by breeders

- 427 is much greater that earlier understood and that both conventional breeding and mutagenesis can
- 428 introduce genomic changes that are more extensive in size and more frequent than the Category 1
- 429 (NGT1) "20 insertions of maximum 20 bp" rule would allow. Regarding genome size and polyploidy,
- 430 the 20/20 rule for NGT1 does not take into account varying plant gene family sizes, the dynamic
- 431 variation of gene family number and genome size in evolution, the effect of the limitation on
- improvement of the many polyploid crops in agriculture. Neither does it address the need for gene
- 433 stacking to combine the traits needed for future-ready crops, which may lead to the limit being easily
- 434 exceeded for many practical breeding goals. We conclude, moreover, that the criteria of
- 435 "naturalness" and "uniqueness" which form the standards for the proposed rule are not met by the
- 436 proposed NGT1 limits.

437 An approach based on the current state of knowledge, which would imply broadening the 20/20

rule, would better support the development of NGT1 plants while still ensuring they are equivalent

- 439 to conventional plants. Such an approach would facilitate breeding for climate resilience, farming
- sustainability, and nutritional security. In March 2025, proposed amendments to Annex I introduced
- by the Polish Presidency of the Council of the EU, which serves through June 2025, appeared to
- 442 achieve a qualified majority of the Council (Permanent Representatives Committee) to proceed to
- the trialogue (negotiations on the terms of the legislation and Annex I between the Council, the
- 444 Commission, and the European Parliament). An important proposed amendment is that NGT1 limits
- would apply per monoploid genome. The 20/20 limit for NGT1, however, would remain in place per
- 446 monoploid genome. The standards reached will greatly influence both the use of NGTs in Europe for
- research and applications and the introduction of NGT products into the marketplace.
- 448

### 449 Acknowledgements

- This work was financially supported by Grant 210021 from the Jane and Aatos Erkko Foundation toA.H.S.
- 452 Author Contributions
- 453 All authors contributed to the review article. A.H.S. conceived of the study and drafted the
- 454 manuscript. F.H., M.J.M.S., J.F.S., R.W., O.-A.R., and K.M. contributed data, text, and interpretation
- and made revisions. A.H.S. prepared the figures. All authors approved the final version of the
- 456 manuscript.

16

- 457 Data availability statement
- 458 Only publicly available data was used and in this study; no new data were generated.
- 459 Conflict of interest disclosure
- 460 The authors declare no conflicts of interest.
- 461 Permission to reproduce material from other sources
- 462 Figure 1 is created by the authors. Figures 2, 3, 4 and 5 is an open access article distributed under the
- terms of the Creative Commons CC Attribution licenses. These permits unrestricted use, distribution,
- 464 and reproduction in any medium, provided the original work is properly cited.
- 465

- 466 References
- Ajawatanawong, P. and Baldauf, S.L. (2013) Evolution of protein indels in plants, animals and fungi.
   *BMC Evol Biol* 13, 140.
- Becana, M., Yruela, I., Sarath, G., Catalan, P. and Hargrove, M.S. (2020) Plant hemoglobins: a journey
   from unicellular green algae to vascular plants. *New Phytol.* 227, 1618-1635.
- 471 Ben-Tov, D., Mafessoni, F., Cucuy, A., Honig, A., Melamed-Bessudo, C. and Levy, A.A. (2024)
  472 Uncovering the dynamics of precise repair at CRISPR/Cas9-induced double-strand breaks.
  473 Nature Comm. 15, 5096.
- Broothaerts, W., Jacchia, S., Angers, A., Petrillo, M., Querci, M., Savini, C., Van Den Eede, G. and
  Emons, H. (2021) *New Genomic Techniques: State-of-the-Art Review, EUR 30430 EN*.
  Luxembourg:Publications Office of the European Union.
- Cai, Y., Gao, X., Mao, J., Liu, Y., Tong, L., Chen, X., Liu, Y., Kou, W., Chang, C., Foster, T., Yao, J., Cornille,
  A., Tahir, M.M., Liu, Z., Yan, Z., Lin, S., Ma, F., Ma, J., Xing, L., An, N., Zuo, X., Lv, Y., Zhao, Z., Li,
  W., Li, Q., Zhao, C., Hu, Y., Liu, H., Wang, C., Shi, X., Ma, D., Fei, Z., Jiang, Y. and Zhang, D.
  (2024) Genome sequencing of 'Fuji' apple clonal varieties reveals genetic mechanism of the
  spur-type morphology. *Nature Comm.* 15, 10082.
- Chitwood-Brown, J., Vallad, G.E., Lee, T.G. and Hutton, S.F. (2021) Characterization and elimination of
   linkage-drag associated with Fusarium wilt race 3 resistance genes. *Theor. Appl. Genet.* 134,
   2129-2140.
- 485 Commission Services (2023) Regulation on new genomic techniques (NGT) Technical paper on the
   486 rationale for the equivalence criteria in Annex I.
- https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CONSIL:ST\_14204\_2023\_INIT
   De Amicis, F. and Marchetti, S. (2000) Intercodon dinucleotides affect codon choice in plant genes.
- 489 *Nucleic Acids Res.* 28, 3339-3345.
  490 Deblieck, M., Szilagyi, G., Andrii, F., Saranga, Y., Lauterberg, M., Neumann, K., Krugman, T., Perovic,
  491 D., Pillen, K. and Ordon, F. (2022) Dissection of a grain yield QTL from wild emmer wheat
- 491 D., Pillen, K. and Ordon, F. (2022) Dissection of a grain yield QTL from wild emmer wheat
  492 reveals sub-intervals associated with culm length and kernel number. *Front. Genet.* 13,
  493 955295.
- 494 Di Genova, A., Almeida, A.M., Munoz-Espinoza, C., Vizoso, P., Travisany, D., Moraga, C., Pinto, M.,
  495 Hinrichsen, P., Orellana, A. and Maass, A. (2014) Whole genome comparison between table
  496 and wine grapes reveals a comprehensive catalog of structural variants. *BMC Plant Biol.* 14,
  497 7.
- 498 EFSA GMO Panel (EFSA Panel on Genetically Modified Organisms). (2012) Scientific opinion
   499 addressing the safety assessment of plants developed using Zinc Finger Nuclease 3 and other
   500 Site-Directed Nucleases with similar function. *EFSA J.* 10, 2943.
- EFSA GMO Panel (EFSA Panel on Genetically Modified Organisms), Mullins E, Bresson J-L, Dalmay T,
   Dewhurst IC, Epstein MM, Firbank LG, Guerche P, Hejatko J, Moreno FJ, Naegeli H, Nogué F,
   Rostoks N, Sánchez Serrano JJ, Savoini G, Veromann E, Veronesi F, Fernandez A, Gennaro A,
   Demadanaulau N, Deffaello T and Schepping D, Mulling F, Bresson LL, Dalmay T,
- 504 Papadopoulou N, Raffaello T and Schoonjans R, Mullins, E., Bresson, J.L., Dalmay, T., 505 Dowburst J.C. Epstein M.M. Firbank J.C. Cuerche P. Heiatke J. Moreno E.L. Nagao
- 505Dewhurst, I.C., Epstein, M.M., Firbank, L.G., Guerche, P., Hejatko, J., Moreno, F.J., Naegeli, H.,506Nogue, F., Rostoks, N., Sanchez Serrano, J.J., Savoini, G., Veromann, E., Veronesi, F.,
- Fernandez, A., Gennaro, A., Papadopoulou, N., Raffaello, T. and Schoonjans, R. (2022) Criteria
   for risk assessment of plants produced by targeted mutagenesis, cisgenesis and intragenesis.
   *EFSA J.* 20, e07618.
- EFSA GMO Panel (EFSA Panel on Genetically Modified Organisms), Naegeli, H., Bresson, J.L., Dalmay,
  T., Dewhurst, I.C., Epstein, M.M., Firbank, L.G., Guerche, P., Hejatko, J., Moreno, F.J., Mullins,
  E., Nogue, F., Sanchez Serrano, J.J., Savoini, G., Veromann, E., Veronesi, F., Casacuberta, J.,
  Gennaro, A., Paraskevopoulos, K., Raffaello, T. and Rostoks, N. (2020) Applicability of the
  EFSA Opinion on site-directed nucleases type 3 for the safety assessment of plants
  developed using site-directed nucleases type 1 and 2 and oligonucleotide-directed
  mutagenesis. *EFSA J.* 18, e06299.

- Gruber, C.W., Elliott, A.G., Ireland, D.C., Delprete, P.G., Dessein, S., Goransson, U., Trabi, M., Wang,
  C.K., Kinghorn, A.B., Robbrecht, E. and Craik, D.J. (2008) Distribution and evolution of circular
  miniproteins in flowering plants. *Plant Cell* 20, 2471-2483.
- Guo, X., Ren, J., Zhou, X., Zhang, M., Lei, C., Chai, R., Zhang, L. and Lu, D. (2024) Strategies to improve
   the efficiency and quality of mutant breeding using heavy-ion beam irradiation. *Crit. Rev. Biotechnol.* 44, 735-752.
- Halpin, C. (2005) Gene stacking in transgenic plants--the challenge for 21st century plant
  biotechnology. *Plant Biotechnol. J.* 3, 141-155.
- Hao, Y., Pan, Y., Chen, W., Rashid, M.A.R., Li, M., Che, N., Duan, X. and Zhao, Y. (2023) Contribution of
   duplicated nucleotide-binding leucine-rich repeat (NLR) genes to wheat disease resistance.
   *Plants* 12, 2794.
- Harmon, D.D., Chen, H., Byrne, D., Liu, W. and Ranney, T.G. (2023) Cytogenetics, ploidy, and genome
   sizes of rose (Rosa spp.) cultivars and breeding lines. *Ornam. Plant Res.* 3, 10.
- Himmelbach A, Walde I, Mascher M, Stein N. (2018) Tethered chromosome conformation capture
   sequencing in Triticeae: A valuable tool for genome assembly. *Bio Protoc.* 8, e2955.
- Hon, T., Mars, K., Young, G., Tsai, Y. C., Karalius, J. W., Landolin, J. M., *et al.* (2020) Highly accurate
   long-read HiFi sequencing data for five complex genomes. *Sci Data* 7, 399.
- Huo, N., Zhu, T., Altenbach, S., Dong, L., Wang, Y., Mohr, T., Liu, Z., Dvorak, J., Luo, M.C. and Gu, Y.Q.
  (2018) Dynamic evolution of alpha-gliadin prolamin gene family in homeologous genomes of hexaploid wheat. *Sci. Rep.* 8, 5181.
- Jayakodi, M., Lu, Q., Pidon, H., Rabanus-Wallace, M.T., Bayer, M., Lux, T., Guo, Y., Jaegle, B., Badea, A.,
  Bekele, W., Brar, G.S., Braune, K., Bunk, B., Chalmers, K.J., Chapman, B., Jorgensen, M.E.,
  Feng, J.W., Feser, M., Fiebig, A., Gundlach, H., Guo, W., Haberer, G., Hansson, M.,
- 540 Himmelbach, A., Hoffie, I., Hoffie, R.E., Hu, H., Isobe, S., Konig, P., Kale, S.M., Kamal, N.,
- 541 Keeble-Gagnere, G., Keller, B., Knauft, M., Koppolu, R., Krattinger, S.G., Kumlehn, J.,
- 542 Langridge, P., Li, C., Marone, M.P., Maurer, A., Mayer, K.F.X., Melzer, M., Muehlbauer, G.J.,
- 543 Murozuka, E., Padmarasu, S., Perovic, D., Pillen, K., Pin, P.A., Pozniak, C.J., Ramsay, L., Pedas,
- P.R., Rutten, T., Sakuma, S., Sato, K., Schuler, D., Schmutzer, T., Scholz, U., Schreiber, M.,
  Shirasawa, K., Simpson, C., Skadhauge, B., Spannagl, M., Steffenson, B.J., Thomsen, H.C.,
- 545 Shirasawa, K., Simpson, C., Skadhauge, B., Spannagl, M., Steffenson, B.J., Thomsen, H.C 546 Tibbits, J.F., Nielsen, M.T.S., Trautewig, C., Veguaud, D., Voss, C., Wang, P., Waugh, R.,
- Westcott, S., Rasmussen, M.W., Zhang, R., Zhang, X.Q., Wicker, T., Dockter, C., Mascher, M.
  and Stein, N. (2024) Structural variation in the pangenome of wild and domesticated barley. *Nature* 636, 654-662.
- Jayakodi, M., Padmarasu, S., Haberer, G., Bonthala, V.S., Gundlach, H., Monat, C., Lux, T., Kamal, N.,
  Lang, D., Himmelbach, A., Ens, J., Zhang, X.Q., Angessa, T.T., Zhou, G., Tan, C., Hill, C., Wang,
  P., Schreiber, M., Boston, L.B., Plott, C., Jenkins, J., Guo, Y., Fiebig, A., Budak, H., Xu, D., Zhang,
  J., Wang, C., Grimwood, J., Schmutz, J., Guo, G., Zhang, G., Mochida, K., Hirayama, T., Sato, K.,
  Chalmers, K.J., Langridge, P., Waugh, R., Pozniak, C.J., Scholz, U., Mayer, K.F.X., Spannagl, M.,
  Li, C., Mascher, M. and Stein, N. (2020) The barley pan-genome reveals the hidden legacy of
  mutation breeding. *Nature* 588, 284-289.
- Jiang, C., Lei, M., Guo, Y., Gao, G., Shi, L., Jin, Y., Cai, Y., Himmelbach, A., Zhou, S., He, Q., Yao, X., Kan,
  J., Haberer, G., Duan, F., Li, L., Liu, J., Zhang, J., Spannagl, M., Liu, C., Stein, N., Feng, Z.,
  Mascher, M. and Yang, P. (2022) A reference-guided TILLING by amplicon-sequencing
  platform supports forward and reverse genetics in barley. *Plant Commun.* 3, 100317.
- Jouanin, A., Borm, T., Boyd, L., Cockram, L., Leigh, F., Santos, B., Visser, R. and Smulders, M. (2019)
   Development of the GlutEnSeq capture system for sequencing gluten gene families in
   hexaploid bread wheat with deletions or mutations induced by γ-irradiation or CRISPR/Cas9.
   *J Cereal Sci* 88, 157-166.
- Jouanin, A., Gilissen, L., Schaart, J.G., Leigh, F.J., Cockram, J., Wallington, E.J., Boyd, L.A., van den
   Broeck, H.C., van der Meer, I.M., America, A.H.P., Visser, R.G.F. and Smulders, M.J.M. (2020)

567 CRISPR/Cas9 Gene Editing of Gluten in Wheat to Reduce Gluten Content and Exposure-568 Reviewing Methods to Screen for Coeliac Safety. Front. Nutr. 7, 51. 569 Kim, D.E., Jensen, D.R., Feldman, D., Tischer, D., Saleem, A., Chow, C.M., Li, X., Carter, L., Milles, L., 570 Nguyen, H., Kang, A., Bera, A.K., Peterson, F.C., Volkman, B.F., Ovchinnikov, S. and Baker, D. 571 (2023) De novo design of small beta barrel proteins. Proc. Natl. Acad. Sci. USA 120, 572 e2207974120. 573 Kitamura, S., Satoh, K., Hase, Y., Yoshihara, R., Oono, Y. and Shikazono, N. (2024) Differential 574 contributions of double-strand break repair pathways to DNA rearrangements following the 575 irradiation of Arabidopsis seeds and seedlings with ion beams. Plant J. 120, 445-458. 576 Kurtz, S., Narechania, A., Stein, J.C. and Ware, D. (2008) A new method to compute K-mer 577 frequencies and its application to annotate large repetitive plant genomes. BMC Genomics 9, 578 517. 579 Li, B., Chen, Z., Chen, H., Wang, C., Song, L., Sun, Y., Cai, Y., Zhou, D., Ouyang, L., Zhu, C., He, H. and 580 Peng, X. (2023) Stacking multiple genes improves resistance to Chilo suppressalis, 581 Magnaporthe oryzae, and Nilaparvata lugens in transgenic rice. Genes (Basel) 14. 582 Lian, Y., Tang, X., Hu, G., Miao, C., Cui, Y., Zhangsun, D., Wu, Y. and Luo, S. (2024) Characterization and 583 evaluation of cytotoxic and antimicrobial activities of cyclotides from Viola japonica. Sci. Rep. 584 14, 9733. 585 Linsky, T.W., Noble, K., Tobin, A.R., Crow, R., Carter, L., Urbauer, J.L., Baker, D. and Strauch, E.M. (2022) 586 Sampling of structure and sequence space of small protein folds. Nature Comm. 13, 7151. 587 Liu, G., Fang, Y., Liu, X., Jiang, J., Ding, G., Wang, Y., Zhao, X., Xu, X., Liu, M., Wang, Y. and Yang, C. 588 (2024) Genome-wide association study and haplotype analysis reveal novel candidate genes 589 for resistance to powdery mildew in soybean. Front. Plant Sci. 15, 1369650. 590 Long, W., Luo, L., Luo, L., Xu, W., Li, Y., Cai, Y. and Xie, H. (2022) Whole genome resequencing of 20 591 accessions of rice landraces reveals Javanica genomic structure variation and allelic 592 genotypes of a grain weight gene TGW2. Front. Plant Sci. 13, 857435. 593 Marin-Sanz, M., Barro, F. and Sanchez-Leon, S. (2023) Unraveling the celiac disease-related 594 immunogenic complexes in a set of wheat and tritordeum genotypes: implications for low-595 aluten precision breeding in cereal crops. Front. Plant Sci. 14, 1171882. 596 Molesini, B., Pandolfini, T., Pii, Y., Korte, A. and Spena, A. (2012) Arabidopsis thaliana AUCSIA-1 597 regulates auxin biology and physically interacts with a kinesin-related protein. PLoS One 7, 598 e41327. 599 Mondragon-Palomino, M., Stam, R., John-Arputharaj, A. and Dresselhaus, T. (2017) Diversification of 600 defensins and NLRs in Arabidopsis species by different evolutionary mechanisms. BMC Evol. 601 Biol. 17, 255. 602 Montenegro, J.D., Golicz, A.A., Bayer, P.E., Hurgobin, B., Lee, H., Chan, C.K., Visendi, P., Lai, K., Dolezel, 603 J., Batley, J. and Edwards, D. (2017) The pangenome of hexaploid bread wheat. Plant J. 90, 604 1007-1013. 605 Niu, J., Wang, W., Wang, Z., Chen, Z., Zhang, X., Qin, Z., Miao, L., Yang, Z., Xie, C., Xin, M., Peng, H., 606 Yao, Y., Liu, J., Ni, Z., Sun, Q. and Guo, W. (2024) Tagging large CNV blocks in wheat boosts 607 digitalization of germplasm resources by ultra-low-coverage sequencing. Genome Biol. 25, 608 171. 609 Paajanen P, Kettleborough G, López-Girona E, Giolai M, Heavens D, Baker D, Lister A, Cugliandolo F, 610 Wilde G, Hein I, Macaulay I, Bryan GJ, Clark MD. (2019) A critical comparison of technologies 611 for a plant genome sequencing project. *Gigascience* 8, ngiy163. 612 Shivaprasad, K.M., Aski, M., Mishra, G.P., Sinha, S.K., Gupta, S., Mishra, D.C., Singh, A.K., Singh, A., Tripathi, K., Kumar, R.R., Kumar, A., Kumar, S. and Dikshit, H.K. (2024) Genome-wide 613 discovery of InDels and validation of PCR-Based InDel markers for earliness in a RIL 614 615 population and genotypes of lentil (Lens culinaris Medik.). PLoS One 19, e0302870.

- Slaman, E., Lammers, M., Angenent, G.C. and de Maagd, R.A. (2023) High-throughput sgRNA testing
  reveals rules for Cas9 specificity and DNA repair in tomato cells. *Front. Genome Ed.* 5,
  1196763.
- Sun, G., Pourkheirandish, M. and Komatsuda, T. (2009) Molecular evolution and phylogeny of the
   RPB2 gene in the genus Hordeum. *Ann. Bot.* 103, 975-983.
- Sun H, Tusso S, Dent CI, Goel M, Wijfjes RY, Baus LC, Dong X, Campoy JA, Kurdadze A, Walkemeier B,
  Sänger C, Huettel B, Hutten RCB, van Eck HJ, Dehmer KJ, Schneeberger K. (2025) The phased
  pan-genome of tetraploid European potato. *Nature* doi: 10.1038/s41586-025-08843-0.
- Tao, N., Xu, X., Ying, Y., Hu, S., Sun, Q., Lv, G. and Gao, J. (2023) Thymosin alpha1 and Its Role in Viral
   Infectious Diseases: The Mechanism and Clinical Application. *Molecules* 28.
- Tong, C., Zhang, Y. and Shi, F. (2022) Genome-wide identification and analysis of the NLR gene family
   in Medicago ruthenica. *Front. Genet.* 13, 1088763.
- Wang, D.W., Li, D., Wang, J., Zhao, Y., Wang, Z., Yue, G., Liu, X., Qin, H., Zhang, K., Dong, L. and Wang,
  D. (2017) Genome-wide analysis of complex wheat gliadins, the dominant carriers of celiac
  disease epitopes. *Sci. Rep.* 7, 44609.
- Wang, L., Ji, Y., Hu, Y., Hu, H., Jia, X., Jiang, M., Zhang, X., Zhao, L., Zhang, Y., Jia, Y., Qin, C., Yu, L.,
  Huang, J., Yang, S., Hurst, L.D. and Tian, D. (2019) The architecture of intra-organism
  mutation rate variation in plants. *PLoS Biol.* 17, e3000191.
- Waschburger, E.L., Filgueiras, J.P.C. and Turchetto-Zolet, A.C. (2024) DOF gene family expansion and
   diversification. *Genet. Mol. Biol.* 46, e20230109.
- Weisweiler, M., Arlt, C., Wu, P.Y., Van Inghelandt, D., Hartwig, T. and Stich, B. (2022) Structural
  variants in the barley gene pool: precision and sensitivity to detect them using short-read
  sequencing and their association with gene expression and phenotypic variation. *Theor. Appl. Genet.* 135, 3511-3529.
- Yadav, C.B., Bhareti, P., Muthamilarasan, M., Mukherjee, M., Khan, Y., Rathi, P. and Prasad, M. (2015)
  Genome-wide SNP identification and characterization in two soybean cultivars with
  contrasting Mungbean Yellow Mosaic India Virus disease resistance traits. *PLoS One* 10,
  e0123897.
- Yan, J., Su, P., Meng, X. and Liu, P. (2023) Phylogeny of the plant receptor-like kinase (RLK) gene family
   and expression analysis of wheat RLK genes in response to biotic and abiotic stresses. *BMC Genomics* 24, 224.
- Yang, T., Ali, M., Lin, L., Li, P., He, H., Zhu, Q., Sun, C., Wu, N., Zhang, X., Huang, T., Li, C.B., Li, C. and
  Deng, L. (2023) Recoloring tomato fruit by CRISPR/Cas9-mediated multiplex gene editing. *Hortic. Res.* 10, uhac214.
- Zeng, Y.X., Wen, Z.H., Ma, L.Y., Ji, Z.J., Li, X.M. and Yang, C.D. (2013) Development of 1047 insertion deletion markers for rice genetic studies and breeding. *Genet. Mol. Res.* 12, 5226-5235.
- Zhang, N., Roberts, H.M., Van Eck, J. and Martin, G.B. (2020) Generation and Molecular
  Characterization of CRISPR/Cas9-Induced Mutations in 63 Immunity-Associated Genes in
  Tomato Reveals Specificity and a Range of Gene Modifications. *Front. Plant Sci.* 11, 10.

Zhang, W., Li, H., Li, Q., Wang, Z., Zeng, W., Yin, H., Qi, K., Zou, Y., Hu, J., Huang, B., Gu, P., Qiao, X. and
Zhang, S. (2023) Genome-wide identification, comparative analysis and functional roles in
flavonoid biosynthesis of cytochrome P450 superfamily in pear (Pyrus spp.). *BMC Genom. Data* 24, 58.