

Foreign DNA detection by high-throughput sequencing to regulate genome-edited agricultural products (Supplementary Information)

Takeshi Itoh (ORCID: 0000-0002-6291-4233)^{1,2,*}, Ritsuko Onuki^{1,2,6}, Mai Tsuda³, Masao Oshima^{2,3}, Masaki Endo (ORCID: 0000-0002-9199-181X)^{2,4}, Hiroaki Sakai^{1,2}, Tsuyoshi Tanaka^{1,5}, Ryo Ohsawa³, and Yutaka Tabei^{2,4}

¹Bioinformatics Team, Advanced Analysis Center, National Agriculture and Food Research Organization, Tsukuba, Ibaraki 305-8602, Japan

²National Institute of Agrobiological Sciences, Tsukuba, Ibaraki 305-8602, Japan

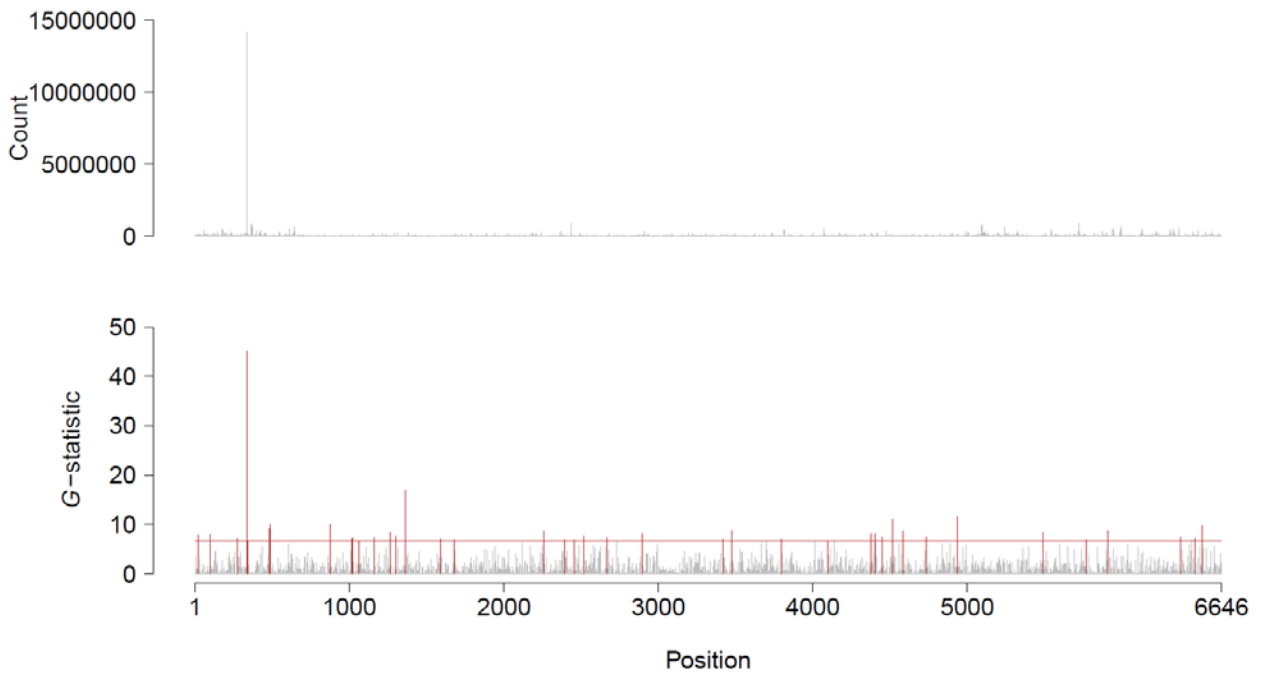
³Tsukuba Plant Innovation Research Center, University of Tsukuba, Tsukuba, Ibaraki 305-8572, Japan

⁴Institute of Agrobiological Sciences, National Agriculture and Food Research Organization, Tsukuba, Ibaraki 305-8634, Japan

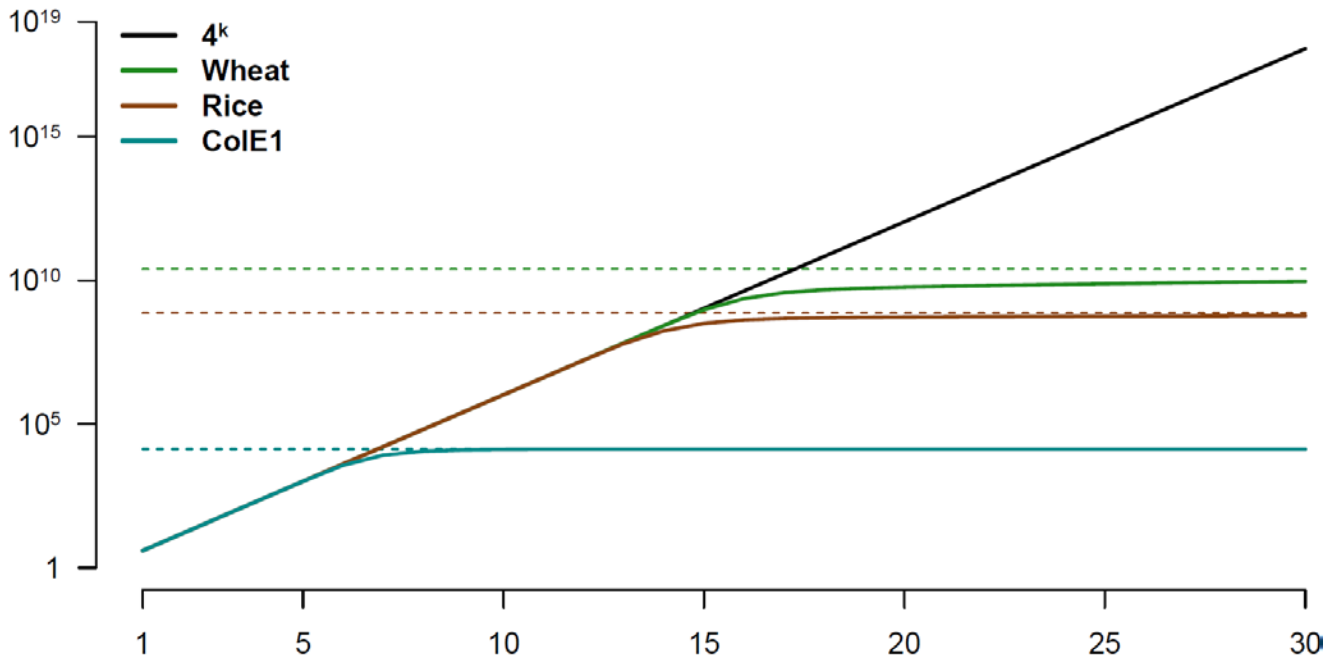
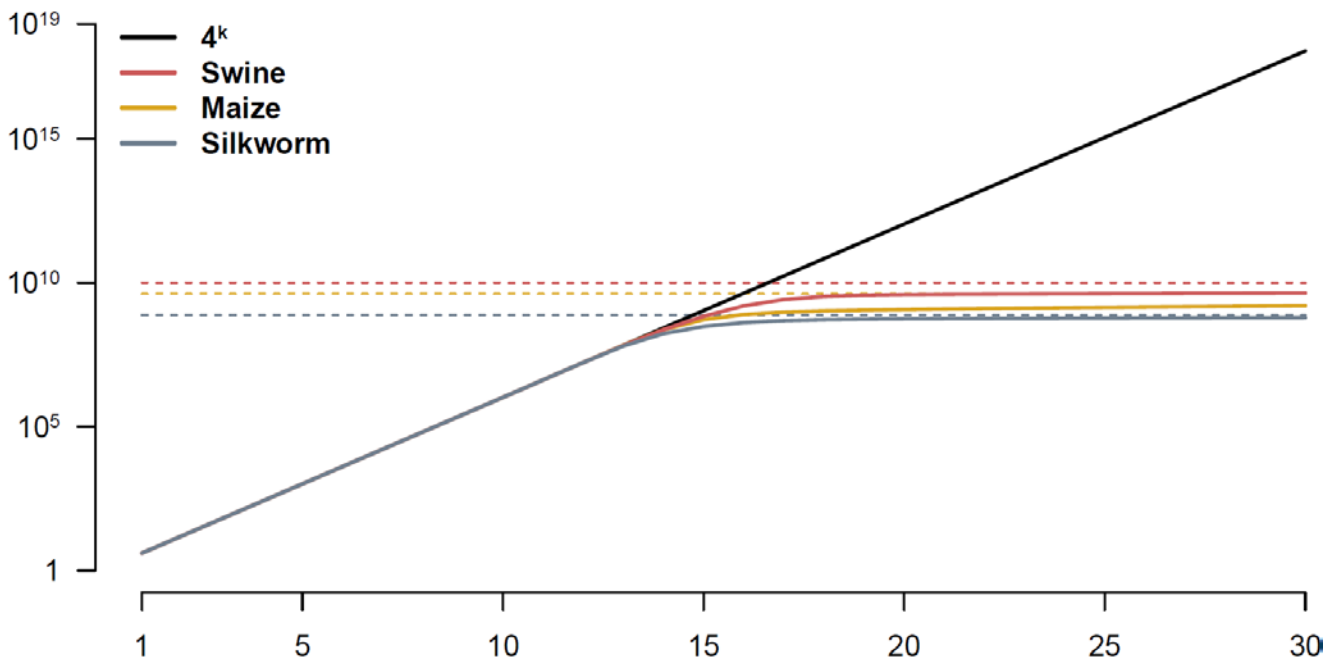
⁵Institute of Crop Science, National Agriculture and Food Research Organization, Tsukuba, Ibaraki 305-8518, Japan

⁶Present address: Research Institute, National Cancer Center Japan, Chuo-ku, Tokyo 104-0045, Japan.

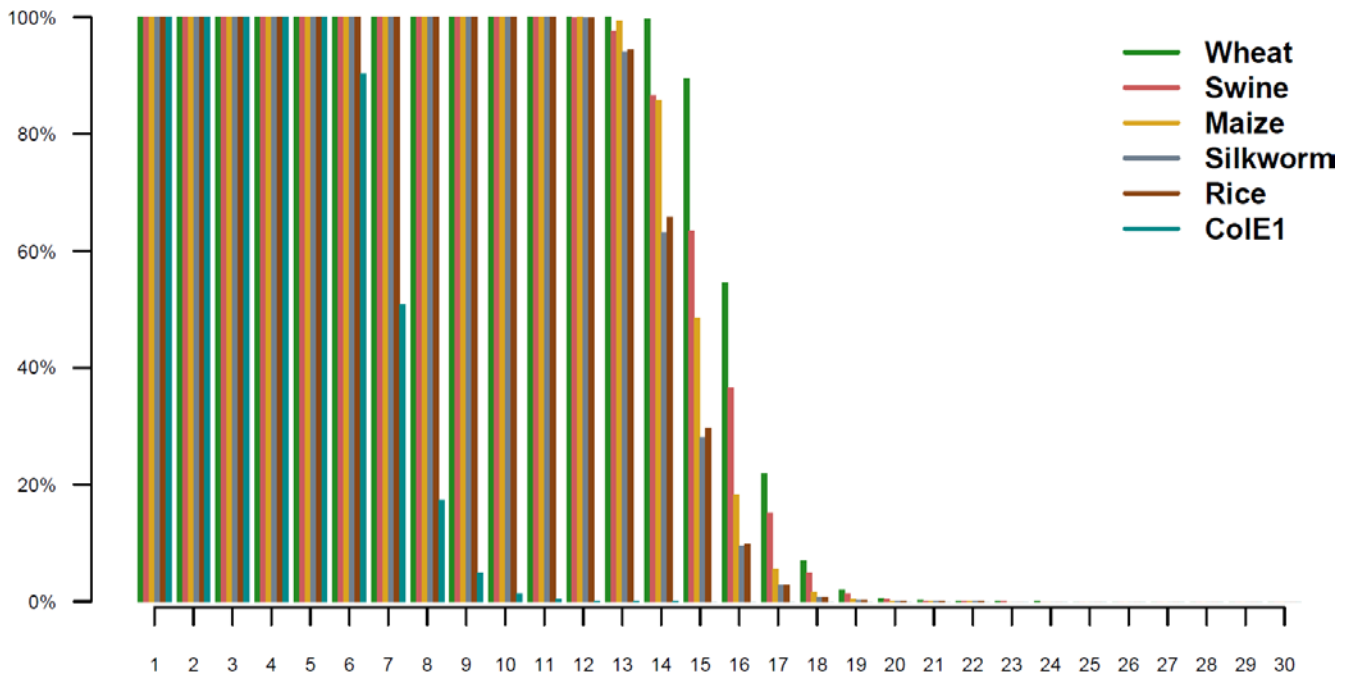
*To whom correspondence may be addressed. Email: taitoh@affrc.go.jp.



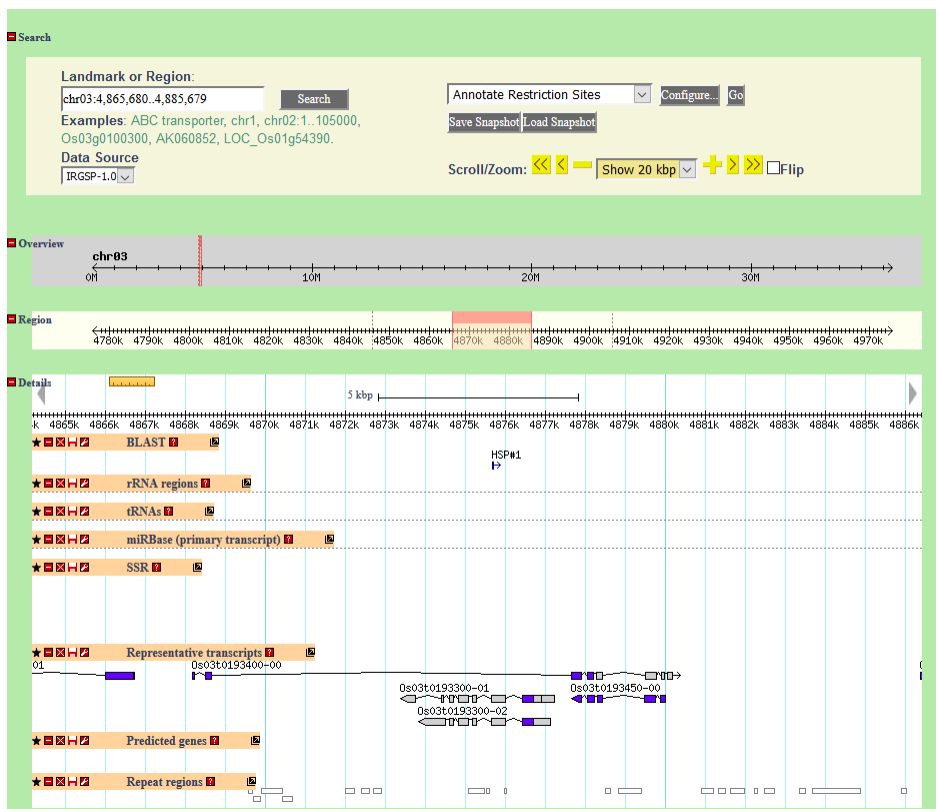
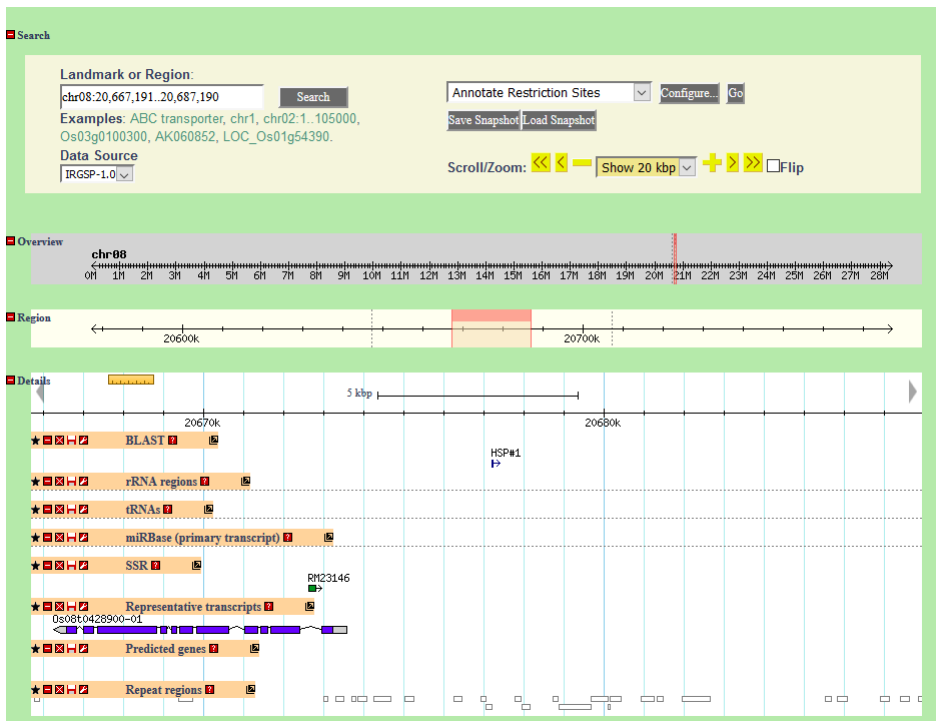
Supplementary Figure 1. A complete version of Fig. 2b.

a**b**

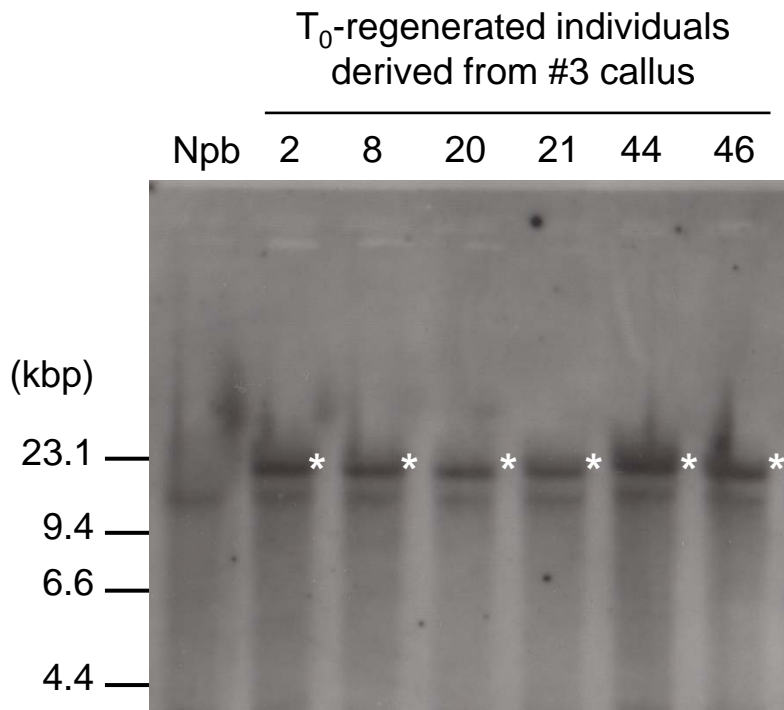
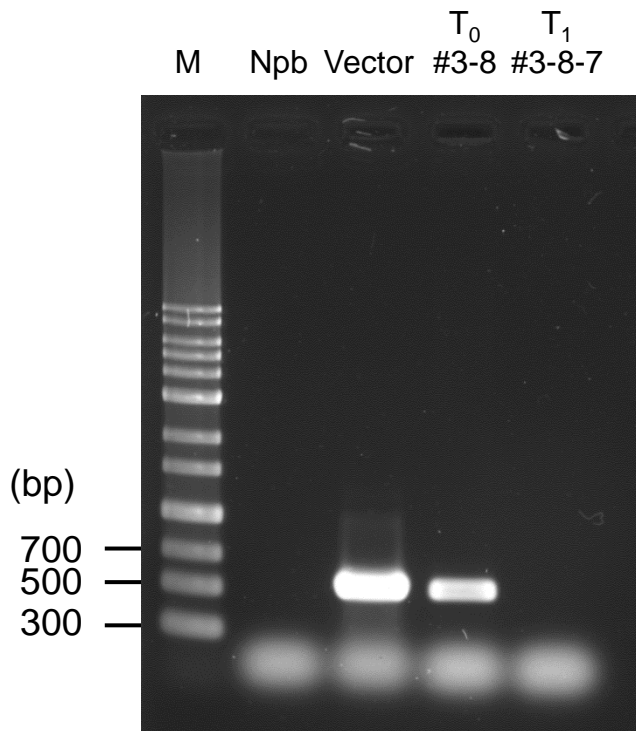
Supplementary Figure 2. The total number of k -mer patterns in the genomes. (a) Wheat, rice and ColE1. (b) Swine, maize and silkworm. The x-axis indicates the k values. The y-axis is log-scaled. The maximum numbers (4^k) are represented by a black line. The dashed lines are the theoretical upper limit of the k -mer numbers. The genome sequences were downloaded from the following sites: wheat (TGACv1), <http://plants.ensembl.org/>; swine (Sscrofa11.1), <http://www.ensembl.org/>; maize (Zm-B73-REFERENCE-GRAMENE-4.0), <https://www.maizegdb.org/>; and silkworm (as of December 12, 2017), <http://sgp.dna.affrc.go.jp/>. For rice and ColE1, see the Methods section.



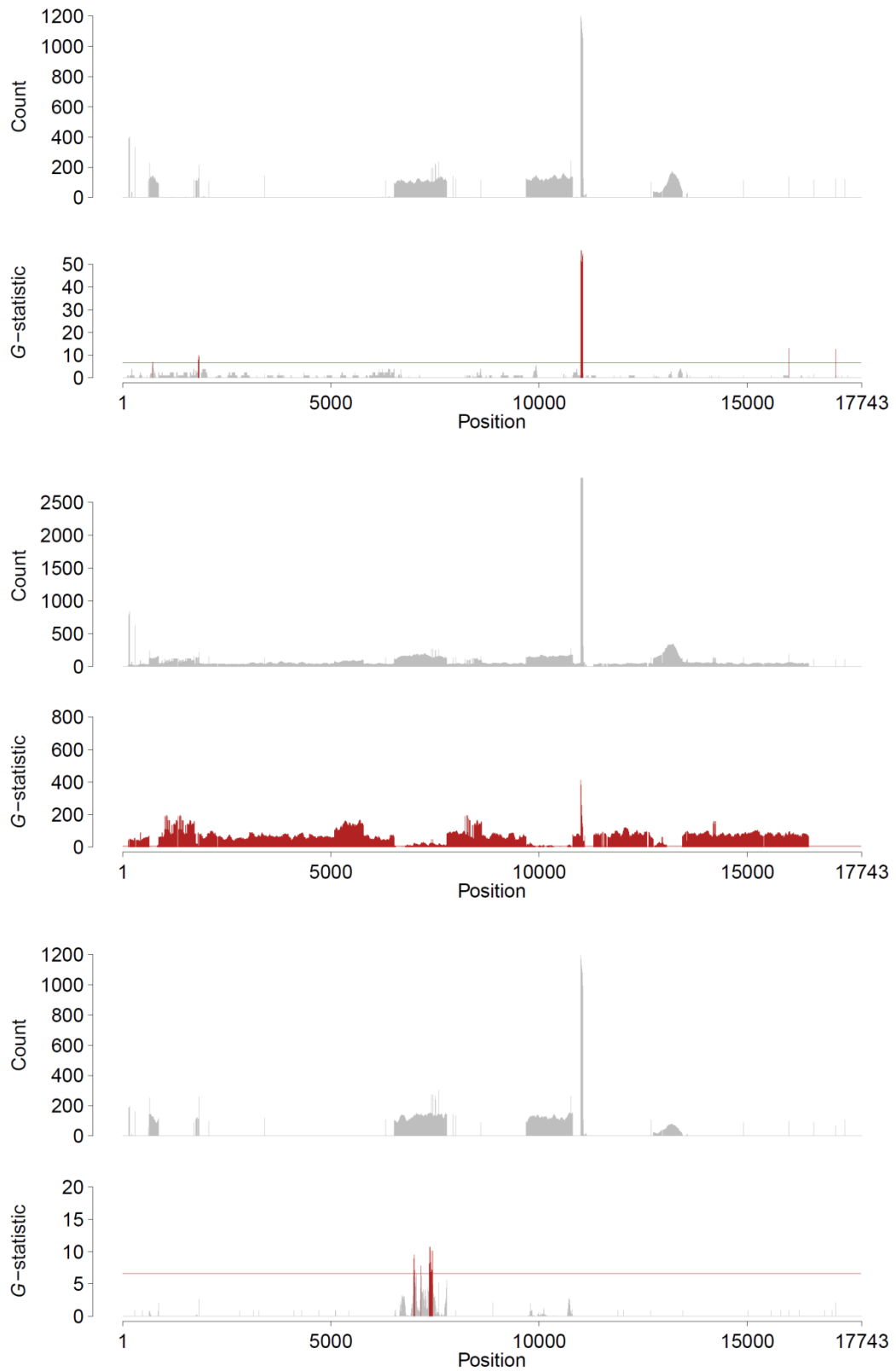
Supplementary Figure 3. The occupancy ratios of the *k*-mer patterns in genomes. The data used are the same as those in Supplementary Fig. 2.

a**b**

Supplementary Figure 4. Two 20-nt sequences that were identical between rice and ColE1. **(a)** 4,875,670-4,875,689 bp (ACAAGGAATTTCTGTTCCC) on chromosome 3 and **(b)** 20,677,181-20,677,200 bp (GCATAAATAGGTTTAATTTT) on chromosome 8. These two sequences are indicated by "HSP#1" on the BLAST track. The sequences were searched for and visualized in RAP-DB (<https://rapdb.dna.affrc.go.jp/>).

a**b**

Supplementary Figure 5. Detection of null segregants. **(a)** Southern blot analysis was conducted for the T₀-regenerated individuals derived from a callus named "#3." Each lane represents a T₀-regenerated individual derived from a single callus. Nipponbare (Npb) is used as a control. Asterisks indicate the signal specific to T₀ plants. **(b)** A PCR experiment was conducted to confirm a null segregant in T₁ plants. For T₀, "8" from #3 (#3-8) was selected and #3-8-7 was used for T₁. M: 1-kbp "DNA Ladder One" marker (Nacalai Tesque, Kyoto, Japan). Neither groupings of cropped gel images nor modifications of the images were made for **(a)** and **(b)**.



Supplementary Figure 6. A complete version of Fig. 3.

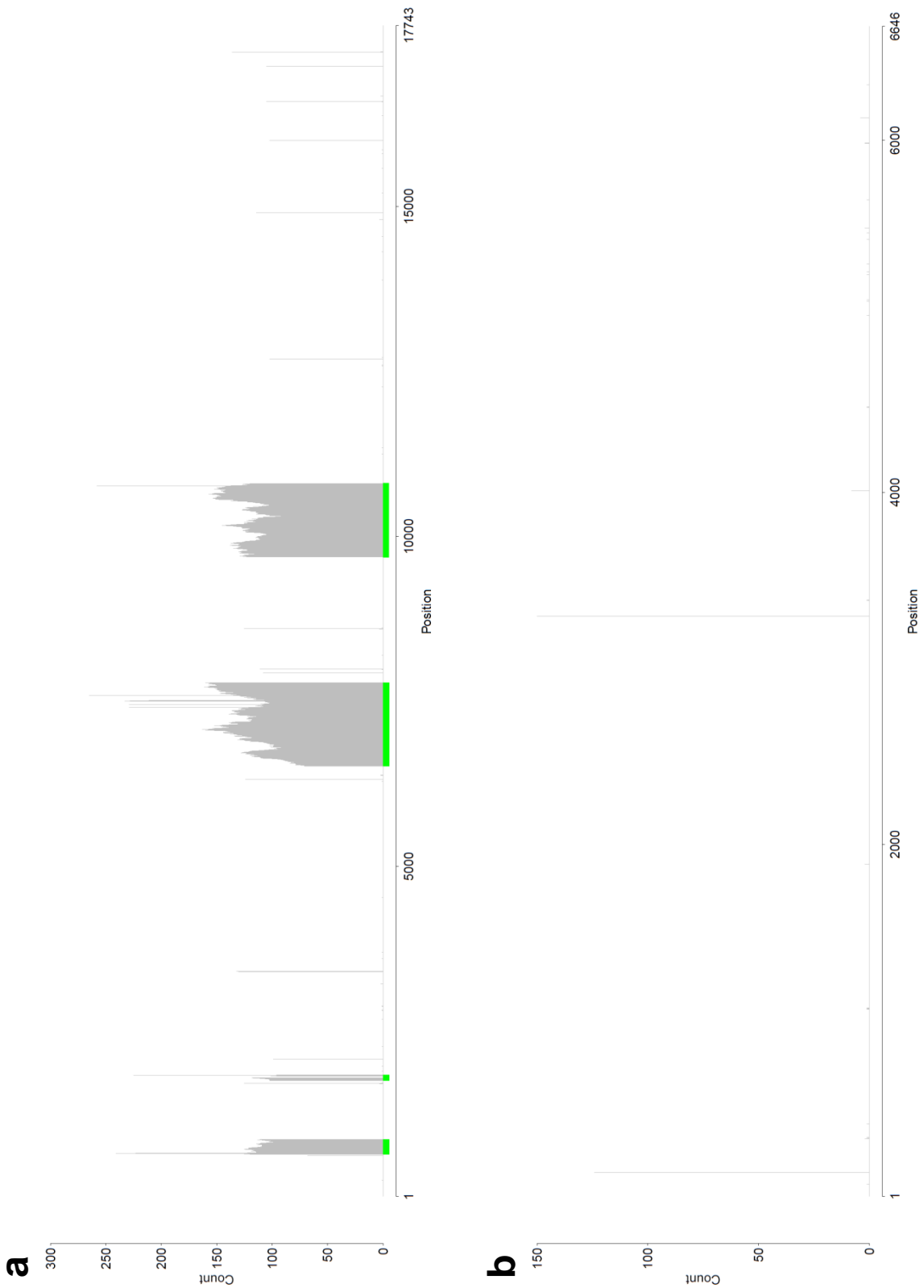
CLUSTAL 2.1 multiple sequence alignment

```

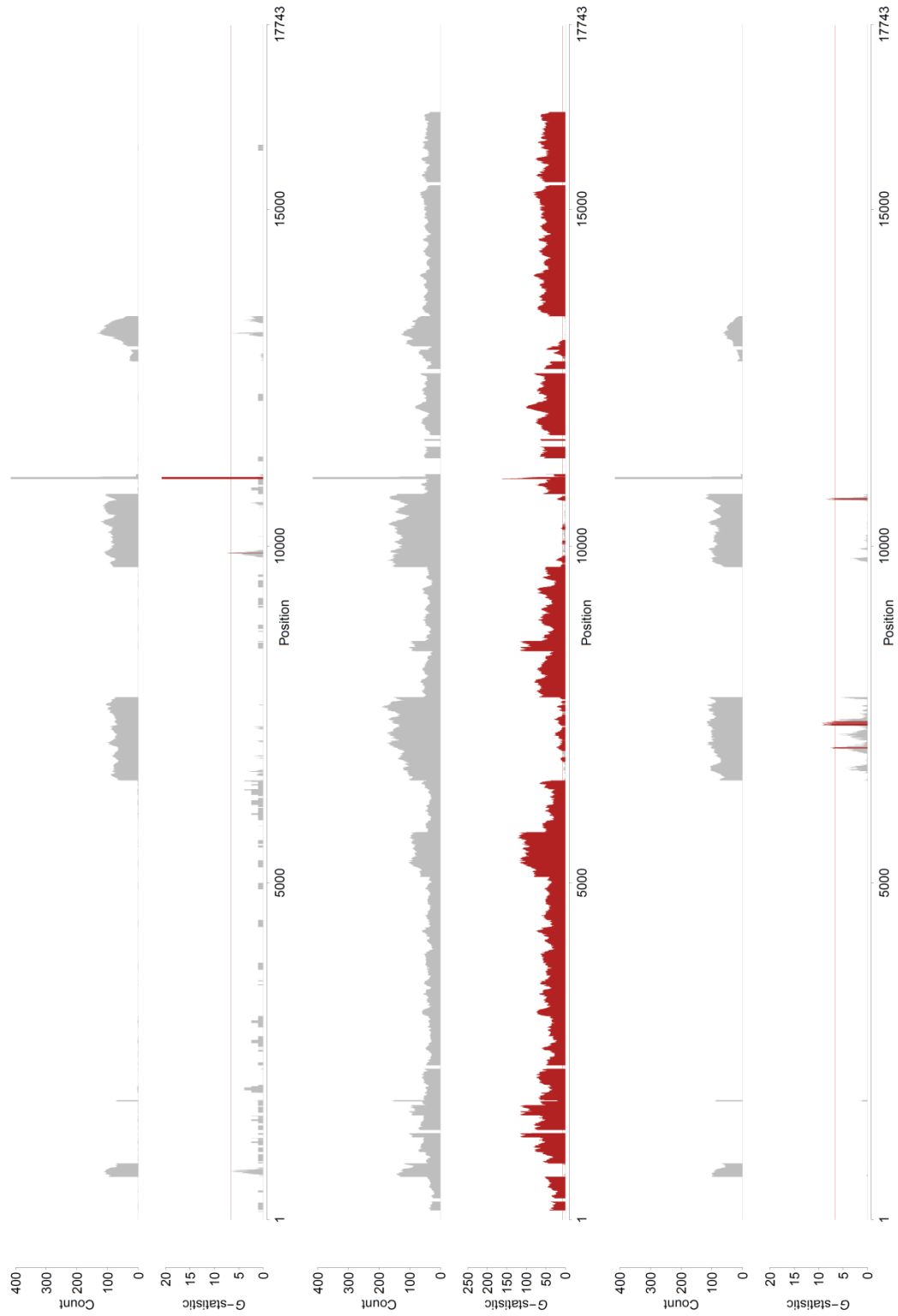
wild_type      ---GTCGTATTACGTGCGCTCACTGGCCGTCGTTTTACAACTGCTGACTGGGAAAC      wild_type      AATTGCGCCGAAAGATGTGGCGCGGATTGGCCATGGAACGCTCGCCGCGCAGTATCA
T1             GTGAGTCGTATTACGTGCGCTCACTGGCCGTCGTTTTACAACTGCTGACTGGGAAAC      T1             AATTGCGCCGAAAGATGTGGCGCGGATTGGCCATGGAACGCTCGCCGCGCAGTATCA
                *****
wild_type      CCTGGCGTTACCCAACCTTAATCGCCTTGCAGCAGTCCCCCTTCGCGAGCTGGCGTAAT      wild_type      GCCGGTGATTCTGGAAGCGCTCAGGCGTATCTGGCCAGGAAGAAGTCTGTCGCCAG
T1             CCTGGCGTTACCCAACCTTAATCGCCTTGCAGCAGTCCCCCTTCGCGAGCTGGCGTAAT      T1             GCCGGTGATTCTGGAAGCGCTCAGGCGTATCTGGCCAGGAAGAAGTCTGTCGCCAG
                *****
wild_type      AGCGAAGAGGCCCGCACCGAAACGCCCTTCCCAACAGTTCGCGAGCCTGAATGGCGAATG      wild_type      CCGTGCGGATCAGCTGGAAGAATTTGTGCATCTACGTGAAAGCGAAATACCAAAAGTGT
T1             AGCGAAGAGGCCCGCACCGAAACGCCCTTCCCAACAGTTCGCGAGCCTGAATGGCGAATG      T1             CCGTGCGGATCAGCTGGAAGAATTTGTGCATCTACGTGAAAGCGAAATACCAAAAGTGT
                *****
wild_type      GGAGCGCCCTGTAGCGGCCACTCAACCTATCTCGGCTATTCTTTGATTATAAGGGA      wild_type      GGTAATAATAACTGTCAGACCAAGTTTACTCATATATACTTTAGATTGATTTAAACT
T1             GGAGCGCCCTGTAGCGGCCACTCAACCTATCTCGGCTATTCTTTGATTATAAGGGA      T1             GGTAATAATAACTGTCAGACCAAGTTTACTCATATATACTTTAGATTGATTTAAACT
                *****
wild_type      TTTTGGCGATTTCGGCCTATTGGTTAAAAATGAGCTGATTTAACAAAAATTAACGCGA      wild_type      TCATTTTTAATTTAAAAGGATCTAGGTGAAGTCTTTTTGATAATCTCATGACCAAAAT
T1             TTTTGGCGATTTCGGCCTATTGGTTAAAAATGAGCTGATTTAACAAAAATTAACGCGA      T1             TCATTTTTAATTTAAAAGGATCTAGGTGAAGTCTTTTTGATAATCTCATGACCAAAAT
                *****
wild_type      ATTTTAAACAAATATTAACGCTTACAATTTAGGTGGCACTTTTCGGGAAATGTGCGCGG      wild_type      CCCTTAACGTGAGTTTTCTGTCACCTGAGCGTCAAGCCCGTAGAAAAGATCAAAAGTATC
T1             ATTTTAAACAAATATTAACGCTTACAATTTAGGTGGCACTTTTCGGGAAATGTGCGCGG      T1             CCCTTAACGTGAGTTTTCTGTCACCTGAGCGTCAAGCCCGTAGAAAAGATCAAAAGTATC
                *****
wild_type      AACCCCTATTGTTTATTTTCTAAATACATTCAAATATGATCCGCTCATGAGACAATA      wild_type      TTCTTGAGATCCTTTTTCTGCGCGTAATCTGCTGTGCAAAACAAAAAACCCAGCT
T1             AACCCCTATTGTTTATTTTCTAAATACATTCAAATATGATCCGCTCATGAGACAATA      T1             TTCTTGAGATCCTTTTTCTGCGCGTAATCTGCTGTGCAAAACAAAAAACCCAGCT
                *****
wild_type      ACCCTGATAAATGCTTCAATAATTTGAAAAGGAAGAGTATGCTAGCCGTAATGGAG      wild_type      ACCAGCGTGTGTTGTTTGGCGGATCAAGAGCTACCAACTCTTTTTCCGAAGTAACTGG
T1             ACCCTGATAAATGCTTCAATAATTTGAAAAGGAAGAGTATGCTAGCCGTAATGGAG      T1             ACCAGCGTGTGTTGTTTGGCGGATCAAGAGCTACCAACTCTTTTTCCGAAGTAACTGG
                *****
wild_type      CCGTACCCTGACCGAAGTACGGCGGTAATGGTGGCGTGGCGGTTTATGGCGTGCTA      wild_type      CTTGAGCAGAGCGCAGATACCAAACTGTTCTTCTAGTGTAGCCGTAGTGGCCACCA
T1             CCGTACCCTGACCGAAGTACGGCGGTAATGGTGGCGTGGCGGTTTATGGCGTGCTA      T1             CTTGAGCAGAGCGCAGATACCAAACTGTTCTTCTAGTGTAGCCGTAGTGGCCACCA
                *****
wild_type      TGATTGCTTTTTTGGCGTGCAGAGCATGCGCGTGCAGCAACAGCAGGCGGTTATGC      wild_type      CTTCAAGAAGTCTGTAGCAGCCCTACATACCTGCTCTGCTAATCTGTTACCAGTGGC
T1             TGATTGCTTTTTTGGCGTGCAGAGCATGCGCGTGCAGCAACAGCAGGCGGTTATGC      T1             CTTCAAGAAGTCTGTAGCAGCCCTACATACCTGCTCTGCTAATCTGTTACCAGTGGC
                *****
wild_type      GGTGGCCGTTGCTGATGCTGTGGAGCAGCAACGATGTGACCCAGCAGGCGAGCCGTC      wild_type      TGCTGCCAGTGGCGATAAGTGTGTTTACCGGGTTGGACTCAAGACGATAGTTACCGBA
T1             GGTGGCCGTTGCTGATGCTGTGGAGCAGCAACGATGTGACCCAGCAGGCGAGCCGTC      T1             TGCTGCCAGTGGCGATAAGTGTGTTTACCGGGTTGGACTCAAGACGATAGTTACCGBA
                *****
wild_type      GAAAACCAACTGAACATTATGCGTGAAGCGGATGTCGGAAGTGAACCCAGCTGCTG      wild_type      TAAGGCGCAGCGTGGCGTGAACGGGGGTTCTGTGCACACAGCCAGCTTGGAGCAAC
T1             GAAAACCAACTGAACATTATGCGTGAAGCGGATGTCGGAAGTGAACCCAGCTGCTG      T1             TAAGGCGCAGCGTGGCGTGAACGGGGGTTCTGTGCACACAGCCAGCTTGGAGCAAC
                *****
wild_type      TGAAGTGGTGGCGGATGAAACGTCATCTGGAACCCAGCCCTGCTGGCCGTGATCTGTA      wild_type      GACCTACACCGAAGTGAATACCTACAGCTGAGCTATGAGAAAGCCACAGCTTCCCGA
T1             TGAAGTGGTGGCGGATGAAACGTCATCTGGAACCCAGCCCTGCTGGCCGTGATCTGTA      T1             GACCTACACCGAAGTGAATACCTACAGCTGAGCTATGAGAAAGCCACAGCTTCCCGA
                *****
wild_type      TGGCAGCGCGTGGATGGCGCCTGAAACCCGATAGCGATATTGATCTGCTGGTGNMNNN      wild_type      AGGGAGAAAGCGGACAGGTATCCGGTAAGCGCAGGTCGGAAACAGGAGGCGCAGG
T1             TGGCAGCGCGTGGATGGCGCCTGAAACCCGATAGCGATATTGATCTGCTGGTGNMNNN      T1             AGGGAGAAAGCGGACAGGTATCCGGTAAGCGCAGGTCGGAAACAGGAGGCGCAGG
                *****
wild_type      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN      wild_type      GGAGCTTCCAGGGGAAACGCTGGTATCTTTANTAGTCTGTCGGGTTTCGCCACCTCT
T1             NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN      T1             GGAGCTTCCAGGGGAAACGCTGGTATCTTTANTAGTCTGTCGGGTTTCGCCACCTCT
                *****
wild_type      GCCCGGTTGAAAGCGAAATCTGCGTGGCGTGAAGTGAACCTGCTGGTGCATGATGACA      wild_type      GACTTGAGCGTGAATTTTGTGATGCTGCTCAGGGGGGCGGAGCTATGAAAAACGCCA
T1             GCCCGGTTGAAAGCGAAATCTGCGTGGCGTGAAGTGAACCTGCTGGTGCATGATGACA      T1             GACTTGAGCGTGAATTTTGTGATGCTGCTCAGGGGGGCGGAGCTATGAAAAACGCCA
                *****
wild_type      TTATCCCGTGGCCTTATCCGCGAAACGTCAGCTGCACT-----                wild_type      GCAACGCGCCTTTTTACGGTTCTGGCCTTTTGTGCTTTTGTGCTATTAGGACCCC
T1             TTATCCCGTGGCCTTATCCGCGAAACGTCAGCTGCACT-----                T1             GCAACGCGCCTTTTTACGGTTCTGGCCTTTTGTGCTTTTGTGCTATTAGGACCCC
                *****
wild_type      -----                wild_type      AGGCTTTACCGAACGACGCGCAGCGAGTCACTGAGCGAGGAGCGGAGAGCGCCCA
T1             -----                T1             AGGCTTTACCGAACGACGCGCAGCGAGTCACTGAGCGAGGAGCGGAGAGCGCCCA
                *****
wild_type      -----                wild_type      ATACGCAAGGAAACAGCTATGACCATGTTAATGACAGCTGGCACGAGTTCGCCACTG
T1             -----                T1             ATACGCAAGGAAACAGCTATGACCATGTTAATGACAGCTGGCACGAGTTCGCCACTG
                *****
wild_type      -----TTGGCAATGGCAGCTAACGATATTCTGGCCGCAT                wild_type      GAAAGCGGCGAGTGAAGTGAACGCCATGAGGCCAGNNNNNNNNNNNNNNNNNNNNNN
T1             -----TTGGCAATGGCAGCTAACGATATTCTGGCCGCAT                T1             GAAAGCGGCGAGTGAAGTGAACGCCATGAGGCCAGNNNNNNNNNNNNNNNNNNNNNN
                *****
wild_type      TTTTGAACCGGACCACTGATATCGATCTGGCCATCTGCTGACCAAGCGCGTGAACA      wild_type      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
T1             TTTTGAACCGGACCACTGATATCGATCTGGCCATCTGCTGACCAAGCGCGTGAACA      T1             NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
                *****
wild_type      TAGCGTGGCGTGGTGGTCCGCGCAGCGAAGAAGTGTGATCCGCTGCCGGAACAGGA      wild_type      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
T1             TAGCGTGGCGTGGTGGTCCGCGCAGCGAAGAAGTGTGATCCGCTGCCGGAACAGGA      T1             NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
                *****
wild_type      CCTGTTTGAAGCGCTGAACGAAACCTGACCTGTGGAACAGCCCGCGGATTGGCGGG      wild_type      NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
T1             CCTGTTTGAAGCGCTGAACGAAACCTGACCTGTGGAACAGCCCGCGGATTGGCGGG      T1             NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
                *****
wild_type      TGATGAACGCAACGTTGCTGACCTGAGCCGATTTGGTATAGCGCGGTGACCGCAA      wild_type      GAAGTTTTAGATTGAGTTCTACTGTCAGCGCGCGGATATCTGCAGA
T1             TGATGAACGCAACGTTGCTGACCTGAGCCGATTTGGTATAGCGCGGTGACCGCAA      T1             GAAGTTTTAGATTGAGTTCTACTGTCAGCGCGCGGATATCTGCAGA
                *****

```

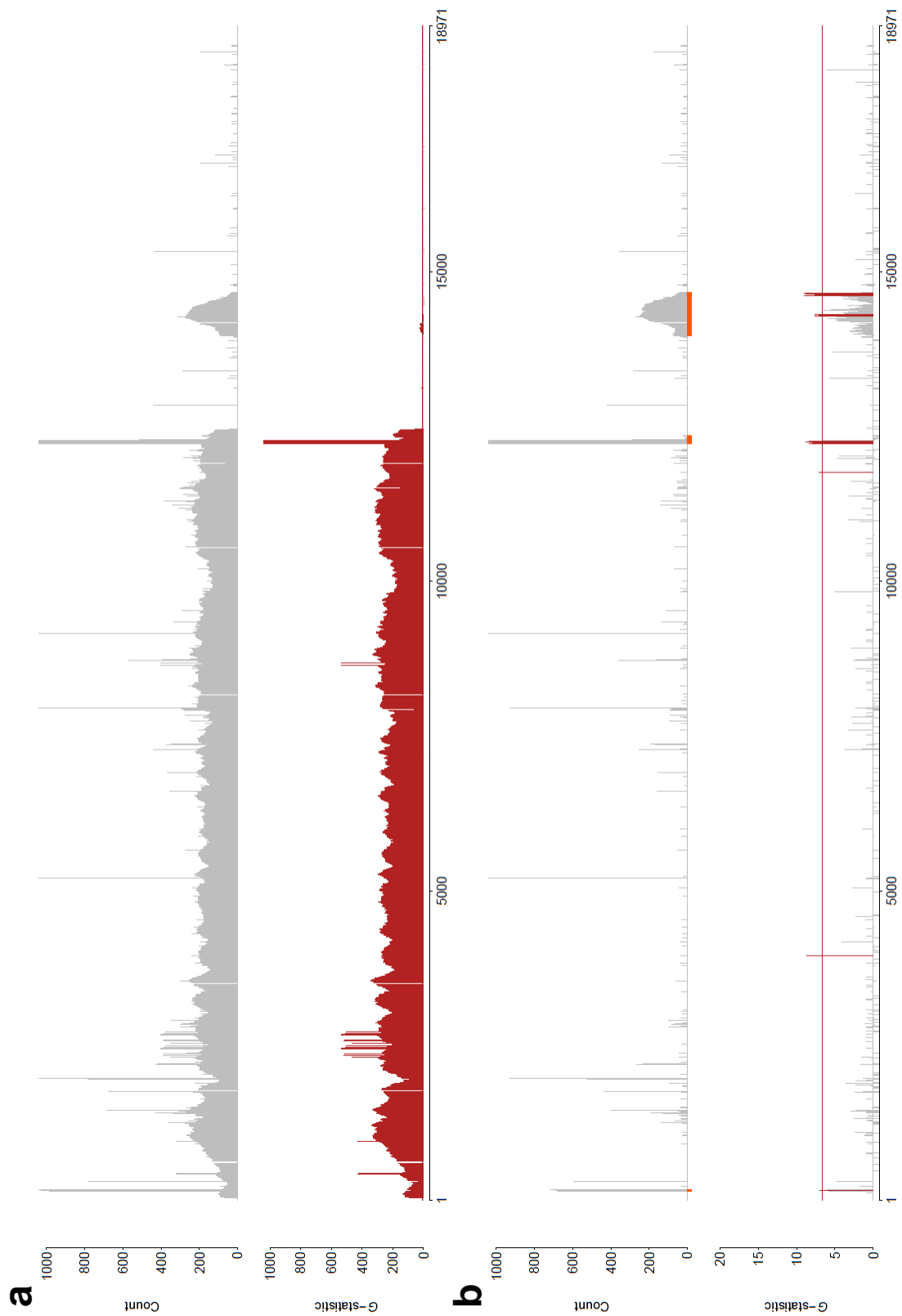
Supplementary Figure 7. The assembled vector-like sequences were obtained from the wild type and T₁ samples in which no vector sequences were expected. The alignment was performed by Clustal X (ver. 2.1) downloaded from <http://www.clustal.org/clustal2/>.



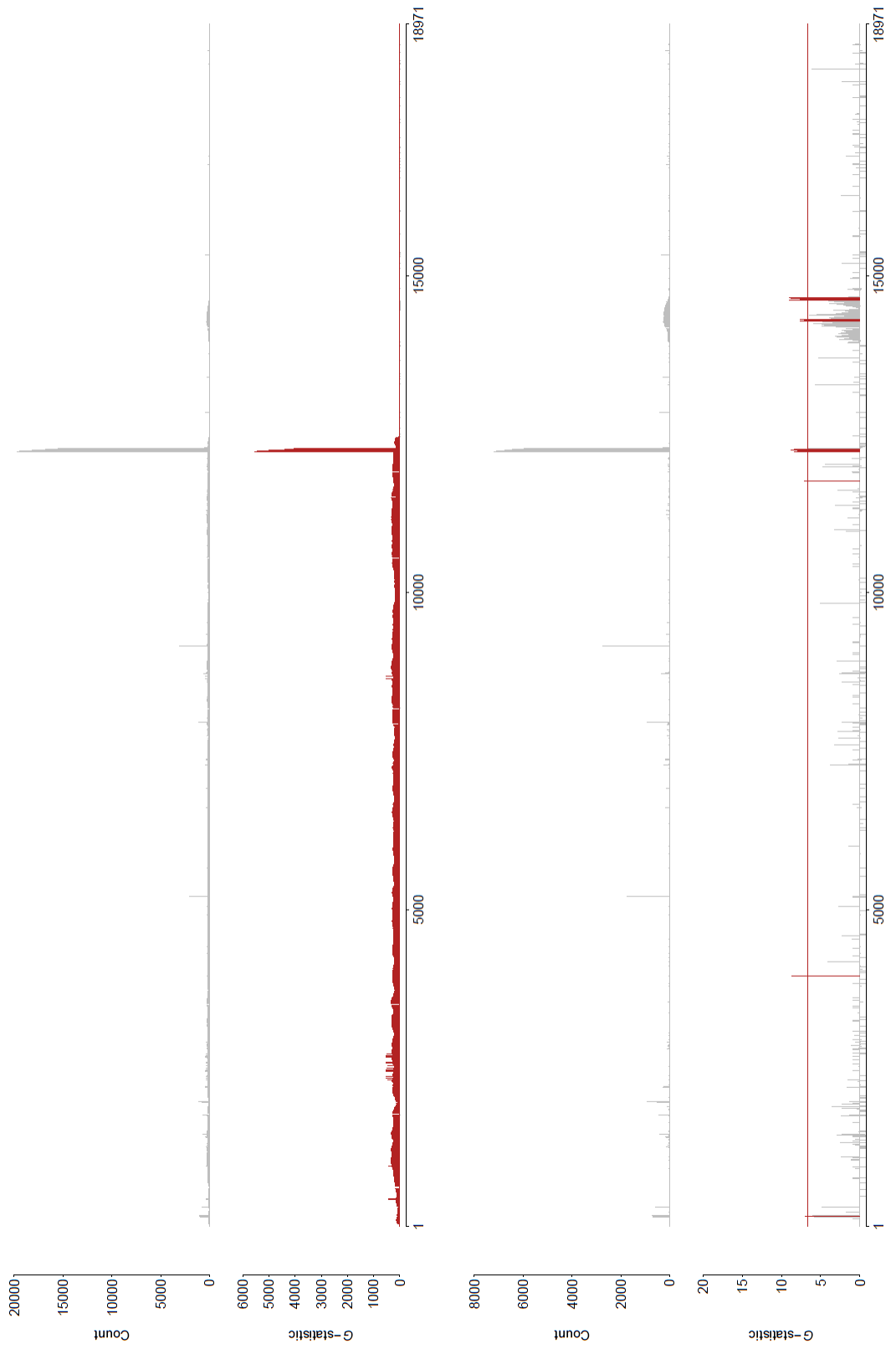
Supplementary Figure 8. Detection of identical 20-mers between the real genome and vector sequences. KAPA's library preparation kit was used. The reads were aligned to (a) the vector used for rice and (b) ColE1. Although some cisgenic regions (green boxes) had a number of hits in (a), the previously detected contamination (orange boxes in Fig. 3a) was not observed.



Supplementary Figure 9. Detection of identical 50-mers between the real genome and vector sequences. For details, see the legend of Fig. 3.

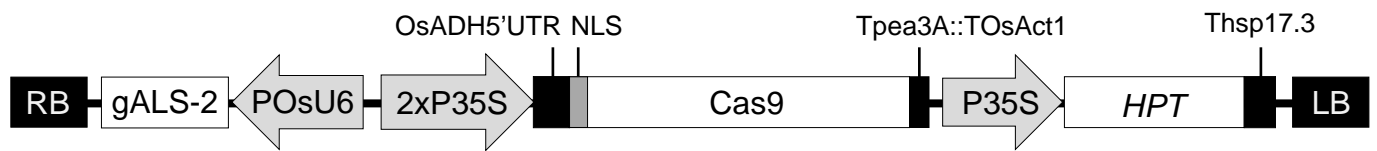


Supplementary Figure 10. Detection of identical 20-mers between the wheat genome and vector sequences. **(a)** T_0 and **(b)** T_1 were examined. The orange boxes in **b** are the regions that were previously reported as a cloning vector-like sequences (Abe, F. *et al.*, Cell Reports, 2019), which are probably due to contamination in the library preparation kit used.



Supplementary Figure 11. A complete version of Supplementary Fig. 10.

pZH_gOsALS-2_Cas9



Supplementary Figure 12. Expression vector construct for sgRNA and SpCas9 in rice.

Supplementary Table 1. Number of properly mapped read pairs to the rice genome.

	Number of pairs	Mapped ratio
Total Reads (50x Coverage)	93,282,216	-
BWA-MEM	93,281,848	99.99961%
NovoAlign	88,697,864	95.08550%
SOAPaligner	92,616,102	99.28592%

Supplementary Table 2. Detection of a foreign DNA insert by *k*-mer search.

Insert length (nt)	<i>k</i>	# of positions with statistical significance*	# of erroneously detected positions
100	50	100	0
10	10	398	398

*Note: true and false hits are not distinguished.

Supplementary Table 4. Average number and standard deviation of false positive hits in the detection

Coverage				
10x	20x	30x	40x	50x
0.55±0.88	0.67±0.94	0.67±0.95	0.68±0.94	0.69±0.96

Supplementary Table 5. Total number of unmapped nucleotides.

	Wild type #1	Wild type #2
BWA-MEM	59,411	58,207
NovoAlign	4,251,829	4,094,610
SOAPaligner	131,133	124,402

Supplementary Table 6. Primers used in this study for PCR of *HPT*.

Primer name	Sequence (5'-3')
HPT-F01	CAAAGATCGTTATGTTTATCGGCACTTTG
HPT-R01	GAAGAAGATGTTGGCGACCTCGTATTG

Supplementary Note 1

Throughout the real data analyses, the reads were all preprocessed by Trimmomatic so that low-quality regions would be removed prior to the k -mer analysis. Theoretically, this process may be omitted because as long as identical k -mers are our concern, all the erroneous sequence segments will simply be disregarded. We compared the preprocessed data with non-preprocessed data from the wild type sample and found that the total number of false positive sites increased from 79 to 102 in the non-preprocessed data. This is probably because sequencing errors led to a specific number of spuriously identical k -mers. Even though the probability that sequencing errors create these false identical hits is extremely low, the enormous amounts of reads generated by modern sequencing technologies could result in tens or hundreds of false positives. Therefore, to reduce false hits, the trimming of the low-quality regions by an appropriate program is recommended.

Supplementary Note 2

The unexpected contamination found in the wild type sequences (orange boxes in Fig. 3a) was examined as follows. The entire read sets from the wild type and T₁ samples, where no vector sequences were expected, were respectively subjected to *de novo* assembly by SOAPdenovo2. These independent assembly attempts obtained essentially the same vector-like sequence (Supplementary Fig. 7), which was partially identical to but largely different from our vector sequence. This sequence was sent to Illumina, Inc. and it was confirmed that this unexpected DNA was related to the use of DNA-binding proteins in the manufacturing of the Illumina reagent kit components (Rooz Golshani, personal communication). While this type of subtle DNA contamination would not hamper an ordinary analysis of high-throughput sequencing data, the assessment results of a generally-used vector sequence, which is a major concern for GMO and NBT regulations, should be interpreted with caution.

Since this contamination was caused by the TruSeq DNA PCR-Free Library Preparation Kit, we additionally used the KAPA Hyper Prep Kit/PCR Free and DNA sequencing was conducted by using Illumina HiSeq X to obtain 151-nt paired-end reads. After trimming by Trimmomatic, a total of 55,310,502,778 nucleotides were subjected to our *k*-mer analysis. As a result, although some false hits were observed in limited small regions, no obvious contamination signals that span several hundreds or thousands of nucleotides were observed (Supplementary Fig. 8). This indicates that if an appropriate library preparation kit was used, the contamination problem does not occur.