



Bioinformatics evaluation of GMO reference genes and their amplification systems

Colaiacovo, M.*, Bonfini, L., Mazzara, M., Savini, C.

* Seidor srl

2024

This document is a publication by the Joint Research Centre (JRC), the European Commission's science and knowledge service. It aims to provide evidence-based scientific support to the European policymaking process. The contents of this publication do not necessarily reflect the position or opinion of the European Commission. Neither the European Commission nor any person acting on behalf of the Commission is responsible for the use that might be made of this publication. For information on the methodology and quality underlying the data used in this publication for which the source is neither Eurostat nor other Commission services, users should contact the referenced source. The designations employed and the presentation of material on the maps do not imply the expression of any opinion whatsoever on the part of the European Union concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries.

EU Science Hub

<https://joint-research-centre.ec.europa.eu>

JRC136933

Geel: European Commission, 2024

© European Union, 2024



The reuse policy of the European Commission documents is implemented by the Commission Decision 2011/833/EU of 12 December 2011 on the reuse of Commission documents (OJ L 330, 14.12.2011, p. 39). Unless otherwise noted, the reuse of this document is authorised under the Creative Commons Attribution 4.0 International (CC BY 4.0) licence (<https://creativecommons.org/licenses/by/4.0/>). This means that reuse is allowed provided appropriate credit is given and any changes are indicated.

For any use or reproduction of photos or other material that is not owned by the European Union permission must be sought directly from the copyright holders.

How to cite this report: European Commission, Joint Research Centre, Colaiacovo, M., Bonfini, L., Mazzara, M. and Savini, C., *Bioinformatics evaluation of GMO reference genes and their amplification systems*, European Commission, Geel, 2024, JRC136933.

The work was carried out for the JRC in its operation of the EURL-GMFF at the JRC Ispra site..

Contents

Abstract.....	2
1 Introduction.....	3
2 Maize.....	6
2.1 Analysis of the <i>adh1</i> reference gene.....	6
2.2 Analysis of the <i>hmg</i> reference gene.....	6
3 Soybean.....	8
3.1 Analysis of the <i>le1</i> reference gene.....	8
4 Cotton.....	9
4.1 Analysis of the <i>sah7</i> reference gene.....	9
4.2 Analysis of the <i>adhC</i> reference gene.....	9
5 Rapeseed.....	11
5.1 Analysis of the <i>fatA(A)</i> reference gene.....	11
6 Conclusions.....	13
References.....	14
List of abbreviations and definitions.....	15
List of figures.....	16
List of tables.....	17
Annexes.....	18
Annex 1. Genome assemblies used for the PCR simulations and the similarity searches.....	18
Annex 2. Primer alignments in the reference genomes.....	20

Abstract

The present report summarises the bioinformatics analyses performed by the European Union Reference Laboratory for Genetically Modified Food and Feed (EURL GMFF) on the reference genes commonly used for detecting GM events in maize, soybean, cotton and rapeseed. The analysis had two objectives: first, to estimate the copy number of these genes in the available genomes of these species to ensure their low and constant copy number in the varieties possibly present on the market; second, to predict expected and unexpected amplifications using *in silico* PCR simulations to verify the specificity of the related taxon-specific methods. In maize, the *adh1* gene was found in three copies in one genome and in a single copy in all others. All of the copies are predicted to be amplified by the QT-TAX-ZM-001 method. In the same species, one copy of the *hmg* gene was found in all genomes, and is predicted to be amplified by the QT-TAX-ZM-002 method. In soybean, one copy of the *le1* gene was found in all genomes, and is predicted to be amplified by the QT-TAX-GM-002 and QT-TAX-GM-003 methods. In cotton, two copies of the *sah7* gene were found in all genomes; both copies are predicted to be amplified by the QT-TAX-GH-021 method. In the same species, two copies of the *adhC* gene were found in all genomes, however only one gene copy is predicted to be amplified by the QT-TAX-GH-019 method. In rapeseed, two copies of the *fatA(A)* gene were found in 5 out of 6 genomes, but only one gene copy is predicted to be amplified by the QT-TAX-BN-001 method. In the remaining genome, four copies of the *fatA(A)* gene were found, but only two of them are predicted to be amplified by the QT-TAX-BN-001 method.

1 Introduction

The present report summarises the bioinformatics analyses performed by the European Union Reference Laboratory for Genetically Modified Food and Feed (EURL GMFF) on the reference genes commonly used to detect GM events in maize, soybean, cotton and rapeseed. The goal of this analysis was to estimate the copy number of these genes in the available genomes of these species to ensure a low and constant copy number in the cultivated varieties possibly present on the market, and to predict expected and unexpected amplifications using *in silico* PCR simulations to verify the specificity of the related taxon-specific methods.

Both analyses were conducted on the same set of genomes downloaded from NCBI¹ on 4th December 2023 following the procedure below:

1. The *assembly_summary.txt* tabular file listing all the available plant genomes was downloaded from the following URL: https://ftp.ncbi.nlm.nih.gov/genomes/genbank/plant/assembly_summary.txt
2. The genome assemblies of the species of interest were selected
3. The genome assemblies were restricted to genome_rep equal to “Full”, and assembly_level equal to “Complete” or “Chromosome” to exclude incomplete or fragmented genomes

A total of 100 genome assemblies were selected, covering several cultivars for each species. The number of genome assemblies for each species is reported in Table 1. The full list of assemblies with additional details is shown in Annex 1.

Table 1. Number of genome assemblies selected for each species of interest.

Species (common name)	Species (scientific name)	Number of genome assemblies
Maize	<i>Zea mays</i>	66
Soybean	<i>Glycine max</i>	19
Cotton	<i>Gossypium hirsutum</i>	9
Rapeseed	<i>Brassica napus</i>	6

The PCR simulations were performed *in silico* using the e-PCR software (version 2.3.12), developed by NCBI (Rotmistrovsky et al., 2004) and the primers of the methods listed in Table 2. The software was run with the following parameters: -n 2 (maximum number of mismatches in each primer), -g 2 (maximum number of gaps in each primer), -f 3 (number of discontinuous words), -t 4 (tabular output with alignments), -m 1000 (margin), -d 20-1000 (amplicon size range). The probes were aligned to the predicted amplicons using the EMBOSS Matcher tool² (version 6.6.0; Huang et al., 1991). The latter was run with the following parameters: -auto (turn off prompts), -stdout (write first file to standard output), -aformat 3 srspair (srs pair format). The probe alignment was checked

¹ <https://www.ncbi.nlm.nih.gov/>

² <https://emboss.sourceforge.net/apps/release/6.6/emboss/apps/matcher.html>

in both orientations (parameter `-sreverse1`), and only alignments with at least 80% identity were accepted. Finally, the predictions were manually inspected.

Table 2. PCR amplification methods considered in this study.

GMOMETHODS ID ⁽¹⁾	Oligos	Target gene	Target species
QT-TAX-ZM-001	Fwd: CGTCGTTTCCCATCTCTTCCTCC	<i>adh1</i>	Maize
	Rev: CCACTCCGAGACCCTCAGTC		
	Probe: FAM-AATCAGGGCTCATTTTCTCGCTCCTCA-TAMRA		
QT-TAX-ZM-002	Fwd: TTGACTAGAAATCTCGTGCTGA	<i>hmg</i>	Maize
	Rev: GCTACATAGGGAGCCTTGCTCT		
	Probe: FAM-CAATCCACACAAACGCACGCGTA-TAMRA		
QT-TAX-GM-003	Fwd: CTTTCTCGACCAATTGACA	<i>le1</i>	Soybean
	Rev: TCAAACCTCAACAGCGACGAC		
	Probe: JOE-CCACAAACACATGCAGGTTATCTTGG-BHQ1		
QT-TAX-GM-002	Fwd: CCAGCTTCGCCGCTTCCTTC	<i>le1</i>	Soybean
	Rev: GAAGGCAAGCCCATCTGCAAGCC		
	Probe: FAM-CTTACCTTCTATGCCCTGACAC-TAMRA		
QT-TAX-GH-021	Fwd: AGTTTGTAGGTTTTGATGTTACATTGAG	<i>sah7</i>	Cotton
	Rev: GCATCTTTGAACCGCCTACTG		
	Probe: FAM-AAACATAAAATAATGGGAACAACCATGACATGT-TAMRA		
QT-TAX-GH-019	Fwd: CACATGACTTAGCCCATCTTTCG	<i>adhC</i>	Cotton
	Rev: CCCACCTTTTTTGGTTTAGC		
	Probe: VIC-TGCAGGTTTTGGTGCCACTGTGAATG-TAMRA		
QT-TAX-BN-001	Fwd: ACAGATGAAGTTCGGGACGAGTAC	<i>fatA(A)</i>	Rapeseed
	Rev: CAGGTTGAGATCCACATGCTTAAATAT		
	Probe: FAM-AAGAAGAATCATCATGCTTC-MGBNFQ		

⁽¹⁾ In the GMOMETHODS Database³, detection methods with identical oligos but different reporter dyes are indicated with a different ID. However, since the bioinformatics analyses described in this report do not take into account the reporter dyes, the same conclusions remain valid for all methods sharing the same primers and probe, i.e.: QT-TAX-ZM-003 and QT-TAX-ZM-014 for QT-TAX-ZM-001; QT-TAX-ZM-005 and QT-TAX-ZM-008 for QT-TAX-ZM-002; QT-TAX-GM-020 for QT-TAX-GM-003; QT-TAX-GM-005 for QT-TAX-GM-002; QT-TAX-GH-018 and QT-TAX-GH-020 for QT-TAX-GH-019.

³ <https://gmo-crl.jrc.ec.europa.eu/gmomethods/>

The similarity searches in the plant genomes were performed using the blastn⁴ software (version 2.12.0; Altschul et al., 1990) with the following parameters: -dust no (low complexity filter off) -soft_masking no (soft masking off). The results were manually inspected and filtered to keep only the alignments covering at least 80% of the query sequence. Table 3 lists the Genbank record IDs that were used as queries to perform the similarity searches for each reference gene. Since a complete record was not available for the *sah7* gene in *Gossypium hirsutum*, the homologous gene in *Gossypium barbadense* (a different cotton species accounting for about 5% of the total cotton production) was used instead (i.e. Genbank ID AY117068.1). The primer alignments for the reference genomes are reported in Annex 2.

Table 3. Genbank records used to perform similarity searches in the plant genomes.

Target gene	Target species	Genbank record ID	Sequence length (bp)
<i>adh1</i>	Maize	X04050.1	5,249
<i>hmg</i>	Maize	AJ131373.1	2,133
<i>le1</i>	Soybean	K00821.1	2,152
<i>sah7</i>	Cotton	AY117068.1	837
<i>adhC</i>	Cotton	AF036569.1	1,651
<i>fatA(A)</i>	Rapeseed	X87842.1	3,007

The results of the analysis are reported in the following sections according to species.

⁴ <https://blast.ncbi.nlm.nih.gov/>

2 Maize

The maize genome is about 2.4 Gbp in size. Two reference genes were analysed for this species: *adh1* and *hmg*, respectively targeted by PCR methods QT-TAX-ZM-001 and QT-TAX-ZM-002. A total of 66 genome assemblies were analysed, including the RefSeq representative genome GCA_902167145.1 (cultivar B73).

2.1 Analysis of the *adh1* reference gene

According to the similarity searches, the *adh1* gene was found in all genome assemblies, therefore it is likely present in most if not all maize cultivars. However, while most genomes have only one gene copy, genome assembly GCA_021307875.1 (cultivar A188) appears to have three gene copies. Considering GCA_902167145.1 as the reference assembly, the *adh1* gene is located on chromosome 1 (Genbank ID LR618874.1, coordinates 280,590,357-280,594,023 reverse strand). It can be found on the Gramene website with ID Zm00001eb056510. Interestingly, GCA_902167065.1 (cultivar Il14H) and GCA_902167135.1 (cultivar CML69) have rearrangements of the gene (partial translocation) or an error in the sequence assembly.

The PCR simulations confirmed that the *adh1* gene can be amplified by the QT-TAX-ZM-001 primers in all genome assemblies. A perfect alignment with primers and probe was obtained for all genome assemblies, including GCA_902167065.1 (cultivar Il14H) and GCA_902167135.1 (cultivar CML69), where rearrangements of the *adh1* gene were noticed. The predicted amplicon size was 135 bp in all genome assemblies.

A peculiar case is genome assembly GCA_021307875.1 (cultivar A188), where three different amplifications are predicted, which is in agreement with the fact that three copies of this gene are present in this genome. The genome of A188 was obtained by single molecule sequencing (Pacific Biosciences technology) and it can be considered of good quality. Cultivar A188 is an inbred maize line that shows great phenotypic difference from other inbred lines, such as in plant height, days to pollination, oil concentration, protein concentration or starch concentration (Ge et al., 2022). Its agronomical traits are poor, as it is shorter and has lower yield compared to other varieties, but its efficient regeneration from stem cells make it an interesting model for maize research. Moreover, line Hi-II, which is commonly used for maize transformation, was derived from the A188 cultivar⁵.

2.2 Analysis of the *hmg* reference gene

According to the similarity searches, the *hmg* gene was found in all genome assemblies, therefore it is likely present in most if not all maize cultivars and always in a single copy. Considering GCA_902167145.1 as the reference assembly, the *hmg* gene is located on chromosome 5 (Genbank ID LR618878.1, coordinates 65,694,795-65,696,975 reverse strand). It can be found on the Gramene⁶ website with ID Zm00001d014858.

The PCR simulations confirmed that the *hmg* gene can be amplified by the QT-TAX-ZM-002 primers in all genome assemblies. In 50 assemblies out of 66, the alignment with the primers was perfect;

⁵ <https://botany.one/2021/10/new-maize-reference-genome-may-help-grow-genetically-modified-lines-of-corn-more-easily/>

⁶ <https://www.gramene.org/>

in 15 assemblies an alignment gap was found in one of the primers, whereas in one assembly a mismatch was found in one of the primers. In all genome assemblies, the probe was predicted to anneal to the amplicon, although one mismatch was found in GCA_902167105.1 (cultivar KI3). The amplicon size was 79 bp.

3 Soybean

The soybean genome is about 1.1 Gbp in size. The *le1* reference gene (lectin) was analysed for this species, targeted by two different PCR methods (QT-TAX-GM-002 and QT-TAX-GM-003). A total of 19 genome assemblies were analysed, including the RefSeq representative genome GCA_000004515.5 (cultivar Williams 82).

3.1 Analysis of the *le1* reference gene

According to the similarity searches, the *le1* gene was found in all genome assemblies, suggesting that it may be present in most if not all soybean cultivars. Moreover, only one gene copy is present in each genome. Considering GCA_000004515.5 as the reference assembly, the *le1* gene is located on chromosome 2 (Genbank ID CM000835.4, coordinates 1,123,507-1,125,658 forward strand). It can be found on the Gramene website with ID GLYMA_02G012600.

The PCR simulations confirmed that the *le1* gene can be amplified by the QT-TAX-GM-002 and QT-TAX-GM-003 methods in all genome assemblies. The alignment with primers and probe was perfect in all genome assemblies, with an amplicon size of 74 bp and 102 bp respectively for the QT-TAX-GM-002 and QT-TAX-GM-003 methods.

4 Cotton

The cotton genome is about 2.3 Gbp in size. Two reference genes were analysed for this species: *sah7* and *adhC*, respectively targeted by PCR methods QT-TAX-GH-021 and QT-TAX-GH-019. A total of 9 genome assemblies were analysed, including the RefSeq representative genome GCA_007990345.1 (unknown cultivar).

4.1 Analysis of the *sah7* reference gene

According to the similarity searches, the *sah7* gene was found in all genome assemblies, therefore it is likely that most cotton cultivars contain this gene. Moreover, two gene copies are present in each genome. Considering GCA_007990345.1 as the reference assembly, the *sah7* gene was found in two genome locations:

- Chromosome A2 (Genbank ID CM017663.1, coordinates 15,958,024-15,959,338 forward strand). Found on Phytozome⁷ with ID Gohir.A02G075700
- Chromosome D2 (Genbank ID CM017676.1, coordinates 14,221,520-14,222,762 forward strand). Found on Phytozome with ID Gohir.D02G082800.

The PCR simulations confirmed that the *sah7* gene can be amplified by the QT-TAX-GH-021 primers in all genome assemblies, with a perfect alignment of the primers. However, since the gene is duplicated in two different genomic locations, two PCR products are predicted (amplicon size of 115 and 123 bp). In both cases, the probe aligns perfectly to the amplicon.

4.2 Analysis of the *adhC* reference gene

According to the similarity searches, the *adhC* gene was found in all genome assemblies, therefore it is likely that most cotton cultivars contain this gene. Moreover, two gene copies are present in each genome. Considering GCA_007990345.1 as the reference assembly, the *adhC* gene was found in the following two genome locations:

- Chromosome D4 (Genbank ID CM017678.1, coordinates 28,702,584-28,706,281 forward strand). Found on Phytozome with ID Gohir.D04G092650.
- Chromosome A4 (Genbank ID CM017665.1, coordinates 34,413,696-34,416,940 reverse strand). Found on Phytozome with ID Gohir.A04G057600.

The PCR simulations confirmed that the *adhC* gene can be amplified by the QT-TAX-GH-019 primers in all genome assemblies, with a perfect alignment of primers and probe. Interestingly, although the gene appears to be duplicated in two different chromosomes, only the first copy of the gene in chromosome D4 could be amplified, with an amplicon size of 73 bp. As shown in Figure 1, the second copy of the gene has three mismatches and one insertion of two base pairs in the region that should be annealed by the forward primer, and should therefore not be amplified.

⁷ <https://phytozome-next.jgi.doe.gov/>

5 Rapeseed

The rapeseed genome is about 1.0 Gbp in size. The *fatA(A)* reference gene was analysed for this species, targeted by PCR method QT-TAX-BN-001. A total of 6 genome assemblies were analysed, including the RefSeq representative genome GCA_020379485.1 (cultivar Da-Ae).

5.1 Analysis of the *fatA(A)* reference gene

According to the similarity searches, the *fatA(A)* gene was found in all genome assemblies, therefore it is likely that most rapeseed cultivars contain this gene. However, while most genomes have only one gene copy, genome assembly GCA_000686985.2 (cultivar ZS11) appears to have two gene copies. Considering GCA_020379485.1 as the reference assembly, the *fatA(A)* gene is located on chromosome C4 (Genbank ID CM035459.1, coordinates 40,807,478-40,809,230 reverse strand).

The PCR simulations confirmed that the *fatA(A)* gene can be amplified by the QT-TAX-BN-001 primers, in all genome assemblies. Differently from what was expected, two PCR products were predicted for each genome assembly, with the exception of GCA_000686985.2 (cultivar ZS11), where four amplifications were predicted. We further investigated the results of the Blast alignment by manual inspection and found another copy of gene *fatA(A)* on chromosome A4 (Genbank ID CM035449.1, coordinates 6,279,586-6,281,411 reverse strand). This second copy was initially missed due to a divergent segment in the central part of the gene, compared to the query sequence; because of this difference, the BLAST software detected two distinct alignments for two gene segments, instead of one contiguous alignment for the entire gene.

As shown in Figure 2, the copy located on chromosome A4 is the one providing the best amplification (perfect alignment with primers and probe, amplicon size 126 bp), while the copy located on chromosome C4 has two mismatches in the region bound by the reverse primer, as well as four mismatches and an insertion of 2 bp in the region bound by the probe, and is therefore unlikely to be detected with the QT-TAX-BN-001 method. This result is in agreement with the data reported in Henderson et al., 2016, where the method QT-TAX-BN-001 was first described. The authors developed this method to target specifically canola-quality *Brassica* species, i.e. the diploid *B. rapa* (turnip rape, AA) and the allotetraploids *B. napus* (rape, AACC) and *B. juncea* (brown mustard, AABB). These species share the sub-genome A, which is not present in the other *Brassica* species, i.e. the allotetraploid *B. carinata* (BBCC) and the diploids *B. nigra* (black mustard, BB) and *B. oleracea* (black mustard, CC). As our results show, QT-TAX-BN-001 can only amplify the copy of *fatA* gene from the A sub-genome, thus confirming the predictions of Henderson and colleagues. The same conclusion is reported in Jacchia et al., 2018, where the authors measured the λ values with ddPCR for several reference gene assays.

Figure 2. Alignment of the predicted PCR amplicons on chromosome A4 (Genbank ID CM035449.1) and on chromosome C4 (Genbank ID CM035459.1). The regions bound by the forward primer, the probe and the reverse primer of QT-TAX-BN-001 are highlighted in red, blue and green, respectively.

```

QT-TAX-BN-001 CAGGTTGAGATCCACATGCTTAAATAT
CM035449.1 (chrA4) : 6280137 CAGGTTGAGATCCACATGCTTAAATATTAAGTGTAACTAGCAATTATAAGG--- 6280189
|||
CM035459.1 (chrC4) : 40808019 CATGTTGAGATCCACATGCTTAAATA--ACTGTTAACTAGCAATTATAATTATT 40808071

QT-TAX-BN-001 GAAGC--ATGATGATTCTTCTT GTACTCGTCCCGAACTTCATCTGT
CM035449.1 (chrA4) : 6280190 GAAGC--ATGATGATTCTTCTTCACCTGGGTTCTCGAGGACAGAAAACCAA GTACTCGTCCCGAACTTCATCTGT 6280262
|||||
CM035459.1 (chrC4) : 40808072 GAAGGGAAGCATGATTCTTTTCACCTGGGTTCTCGAGGACAGAAAACCAA GTACTCGTCCCGAACTTCATCTGT 40808146

```

6 Conclusions

The similarity searches and PCR simulations performed on the genomes of the species of interest allowed us to draw the following conclusions.

Maize:

- One copy of the *adh1* gene was found in 65 out of 66 genome assemblies. In one genome assembly (cultivar A188), three copies of the gene were found. Method QT-TAX-ZM-001 should detect this gene in all the assemblies with a perfect alignment of primers and probe. Three PCR amplifications were predicted for the genome assembly of cultivar A188.
- One copy of the *hmg* gene was found in all genome assemblies. Method QT-TAX-ZM-002 should detect this gene in all the assemblies. In 50 out of 66 assemblies, the alignment of the primers is perfect, while in the remaining assemblies one gap or one mismatch was found for one of the primers, although this should not prevent the PCR amplification. In one assembly, a mismatch was found in the alignment of the probe, while in the remaining assemblies the alignment was perfect.

Soybean:

- One copy of the *le1* gene was found in all genome assemblies. Methods QT-TAX-GM-002 and QT-TAX-GM-003 should detect this gene in all the assemblies. In all cases, the alignment of primers and probe is perfect.

Cotton:

- Two copies of the *sah7* gene were found in all genome assemblies. Both gene copies can be amplified by method QT-TAX-GH-021, with a perfect alignment of primers and probe.
- Two copies of the *adhC* gene were found in all genome assemblies. However, only one copy of this gene can be amplified by method QT-TAX-GH-019, with a perfect alignment of primers and probe.

Rapeseed:

- Two copies of the *fatA(A)* gene were found in 5 out of 6 genome assemblies. In one genome assembly (cultivar ZS11), four gene copies were found. However, only one copy of this gene can be amplified by method QT-TAX-BN-001 (two copies for cultivar ZS11), with a perfect alignment of primers and probe.

References

Altschul, S.F. et al., 'Basic local alignment search tool', *Journal of Molecular Biology*, Vol. 215, Issue 3, 1990, pp. 403-410.

Ge, F. et al., 'Genome assembly of the maize inbred line A188 provides a new reference genome for functional genomics', *The Crop Journal*, Vol. 10, Issue 1, 2022, pp. 47-55.

Henderson, N. et al., 'PCR-based detection and quantification of a transgenic glyphosate-tolerant canola using a novel reference gene system', *Food Analytical Methods*, Vol. 9, 2015, pp. 353-361.

Huang, X. and Miller, W., 'A time-efficient, linear-space local similarity algorithm', *Advances in Applied Mathematics*, Vol. 12, Issue 13, 1991, pp. 337-357.

Jacchia, S. et al., 'Identification of single target taxon-specific reference assays for the most commonly genetically transformed crops using digital droplet PCR', *Food Control*, Vol. 93, 2018, pp. 191-200.

Rotmistrovsky, K. et al., 'A web server for performing electronic PCR', *Nucleic Acids Research*, Vol. 32, Issue suppl_2, 2004, pp. W108-W112.

List of abbreviations and definitions

Abbreviations	Definitions
Gbp	Giga base pairs (1 Gbp = 1,000,000,000 bp)
EURL GMFF	European Union Reference Laboratory for Genetically Modified Food and Feed
NCBI	National Center for Biotechnology Information
PCR	Polymerase Chain Reaction

List of figures

Figure 1. Alignment of the predicted PCR amplicons on chromosome D4 (Genbank ID CM017678.1) and on chromosome A4 (Genbank ID CM017665.1). The regions bound by the forward primer, the probe and the reverse primer of QT-TAX-GH-019 are highlighted in red, blue and green, respectively.
.....10

Figure 2. Alignment of the predicted PCR amplicons on chromosome A4 (Genbank ID CM035449.1) and on chromosome C4 (Genbank ID CM035459.1). The regions bound by the forward primer, the probe and the reverse primer of QT-TAX-BN-001 are highlighted in red, blue and green, respectively.
.....12

List of tables

Table 1. Number of genome assemblies selected for each species of interest	3
Table 2. PCR amplification methods considered in this study.....	4
Table 3. Genbank records used to perform similarity searches in the plant genomes.....	5

Annexes

Annex 1. Genome assemblies used for the PCR simulations and the similarity searches

assembly accession	species	cultivar	assembly level
GCA_000987745.1	<i>Gossypium hirsutum</i>	TM-1	Chromosome
GCA_006980745.1	<i>Gossypium hirsutum</i>	TM1	Chromosome
GCA_006980775.1	<i>Gossypium hirsutum</i>	ZM24	Chromosome
GCA_007990345.1 *	<i>Gossypium hirsutum</i>	not available	Chromosome
GCA_018997965.1	<i>Gossypium hirsutum</i>	NDM8	Chromosome
GCA_021461695.1	<i>Gossypium hirsutum</i>	not available	Chromosome
GCA_021461685.1	<i>Gossypium hirsutum</i>	BARBREN-713	Chromosome
GCA_024704785.1	<i>Gossypium hirsutum</i>	JBM	Chromosome
GCA_024600755.1	<i>Gossypium hirsutum</i>	Zhongzhimian No.2	Chromosome
GCA_000686985.2	<i>Brassica napus</i>	ZS11	Chromosome
GCA_020379485.1 *	<i>Brassica napus</i>	Da-Ae	Chromosome
GCA_026770265.1	<i>Brassica napus</i>	Zy821	Chromosome
GCA_026770255.1	<i>Brassica napus</i>	GH06	Chromosome
GCA_905183035.1	<i>Brassica napus</i>	not available	Chromosome
GCA_033556945.1	<i>Brassica napus</i>	NTS57	Chromosome
GCA_002905335.2	<i>Glycine max</i>	Lee	Chromosome
GCA_003349995.2	<i>Glycine max</i>	Zhonghuang 13	Chromosome
GCA_014282185.1	<i>Glycine max</i>	Williams 82	Chromosome
GCA_014282065.1	<i>Glycine max</i>	Hefeng 25	Chromosome
GCA_014282095.1	<i>Glycine max</i>	Jinyuan	Chromosome
GCA_014282075.1	<i>Glycine max</i>	Zhonghuang 35	Chromosome
GCA_014282145.1	<i>Glycine max</i>	Huaxia 3	Chromosome
GCA_014282035.1	<i>Glycine max</i>	Wenfeng 7	Chromosome
GCA_014282085.1	<i>Glycine max</i>	Zhonghuang 13	Chromosome
GCA_015227745.1	<i>Glycine max</i>	Tianlong1	Chromosome
GCA_000004515.5 *	<i>Glycine max</i>	Williams 82	Chromosome
GCA_019321705.1	<i>Glycine max</i>	PI594527	Chromosome
GCA_012273815.2	<i>Glycine max</i>	EMBRAPA BRS 537	Chromosome
GCA_020497155.1	<i>Glycine max</i>	Hwangkeum	Chromosome
GCA_021733175.1	<i>Glycine max</i>	Jidou 17	Chromosome
GCA_022114995.1	<i>Glycine max</i>	ZN6	Chromosome
GCA_029216145.1	<i>Glycine max</i>	ENREI	Chromosome
GCA_030864155.1	<i>Glycine max</i>	Williams 82	Complete Genome
GCA_033623075.1	<i>Glycine max</i>	Jack	Complete Genome
GCA_000005005.6	<i>Zea mays</i>	B73	Chromosome
GCA_003185045.1	<i>Zea mays</i>	inbred line Mo17	Chromosome
GCA_003709335.1	<i>Zea mays</i>	DK105	Chromosome
GCA_003704525.1	<i>Zea mays</i>	PE0075	Chromosome
GCA_009176585.1	<i>Zea mays</i>	K0326Y	Chromosome
GCA_902166995.1	<i>Zea mays</i>	M162W	Chromosome

GCA_902167005.1	<i>Zea mays</i>	NC350	Chromosome
GCA_902167045.1	<i>Zea mays</i>	Ky21	Chromosome
GCA_902167075.1	<i>Zea mays</i>	B97	Chromosome
GCA_902167135.1	<i>Zea mays</i>	CML69	Chromosome
GCA_902167145.1 *	<i>Zea mays</i>	B73	Chromosome
GCA_902167095.1	<i>Zea mays</i>	CML322	Chromosome
GCA_902167185.1	<i>Zea mays</i>	TZ18	Chromosome
GCA_902166985.1	<i>Zea mays</i>	CML277	Chromosome
GCA_902166975.1	<i>Zea mays</i>	CML247	Chromosome
GCA_902167035.1	<i>Zea mays</i>	NC358	Chromosome
GCA_902167105.1	<i>Zea mays</i>	KI3	Chromosome
GCA_902167155.1	<i>Zea mays</i>	CML228	Chromosome
GCA_902167375.1	<i>Zea mays</i>	MS71	Chromosome
GCA_902166965.1	<i>Zea mays</i>	P39	Chromosome
GCA_902167015.1	<i>Zea mays</i>	OH43	Chromosome
GCA_902167065.1	<i>Zea mays</i>	Il14H	Chromosome
GCA_902167085.1	<i>Zea mays</i>	Hp301	Chromosome
GCA_902167175.1	<i>Zea mays</i>	M37W	Chromosome
GCA_902167205.1	<i>Zea mays</i>	Tx303	Chromosome
GCA_902166955.1	<i>Zea mays</i>	OH7B	Chromosome
GCA_902167025.1	<i>Zea mays</i>	CML333	Chromosome
GCA_902167055.1	<i>Zea mays</i>	CML103	Chromosome
GCA_902167115.1	<i>Zea mays</i>	KI11	Chromosome
GCA_902167165.1	<i>Zea mays</i>	Mo18W	Chromosome
GCA_902373975.1	<i>Zea mays</i>	CML52	Chromosome
GCA_902714155.1	<i>Zea mays</i>	B73 Ab10	Chromosome
GCA_014529475.1	<i>Zea mays</i>	not available	Chromosome
GCA_905067065.1	<i>Zea mays</i>	not available	Chromosome
GCA_016432965.1	<i>Zea mays</i>	la453-sh2	Chromosome
GCA_019096015.1	<i>Zea mays</i>	PHB47	Chromosome
GCA_019095995.1	<i>Zea mays</i>	PHJ40	Chromosome
GCA_019095975.1	<i>Zea mays</i>	LH145	Chromosome
GCA_019096025.1	<i>Zea mays</i>	NKH8431	Chromosome
GCA_019095955.1	<i>Zea mays</i>	B84	Chromosome
GCA_910593975.1	<i>Zea mays</i>	B104	Chromosome
GCA_021307875.1	<i>Zea mays</i>	A188	Chromosome
GCA_022117705.1	<i>Zea mays</i>	Mo17-2021	Complete Genome
GCA_024505845.1	<i>Zea mays</i>	inbred line Dan340	Chromosome
GCA_024718635.1	<i>Zea mays</i>	not available	Chromosome
GCA_026226675.1	<i>Zea mays</i>	S37	Chromosome
GCA_026226635.1	<i>Zea mays</i>	PH207	Chromosome
GCA_026228945.1	<i>Zea mays</i>	huangzaosi	Chromosome
GCA_026226085.1	<i>Zea mays</i>	Zheng58	Chromosome
GCA_026226705.1	<i>Zea mays</i>	Jing724	Chromosome
GCA_026229665.1	<i>Zea mays</i>	Chang7-2	Chromosome
GCA_026226095.1	<i>Zea mays</i>	Xu178	Chromosome

GETTING IN TOUCH WITH THE EU

In person

All over the European Union there are hundreds of Europe Direct centres. You can find the address of the centre nearest you online (european-union.europa.eu/contact-eu/meet-us_en).

On the phone or in writing

Europe Direct is a service that answers your questions about the European Union. You can contact this service:

- by freephone: 00 800 6 7 8 9 10 11 (certain operators may charge for these calls),
- at the following standard number: +32 22999696,
- via the following form: european-union.europa.eu/contact-eu/write-us_en.

FINDING INFORMATION ABOUT THE EU

Online

Information about the European Union in all the official languages of the EU is available on the Europa website (european-union.europa.eu).

EU publications

You can view or order EU publications at op.europa.eu/en/publications. Multiple copies of free publications can be obtained by contacting Europe Direct or your local documentation centre (european-union.europa.eu/contact-eu/meet-us_en).

EU law and related documents

For access to legal information from the EU, including all EU law since 1951 in all the official language versions, go to EUR-Lex (eur-lex.europa.eu).

Open data from the EU

The portal data.europa.eu provides access to open datasets from the EU institutions, bodies and agencies. These can be downloaded and reused for free, for both commercial and non-commercial purposes. The portal also provides access to a wealth of datasets from European countries.

Science for policy

The Joint Research Centre (JRC) provides independent, evidence-based knowledge and science, supporting EU policies to positively impact society



EU Science Hub

joint-research-centre.ec.europa.eu