

RESEARCH NOTE

Open Access



# Identification of annotation artifacts concerning the *chalcone synthase* (CHS)

Martin Bartas<sup>1</sup> , Adriana Volna<sup>2</sup> , Jiri Cerven<sup>1</sup> and Boas Pucker<sup>3\*</sup>

## Abstract

**Objective** Chalcone synthase (CHS) catalyzes the initial step of the flavonoid biosynthesis. The CHS encoding gene is well studied in numerous plant species. Rapidly growing sequence databases contain hundreds of CHS entries that are the result of automatic annotation. In this study, we evaluated apparent multiplication of CHS domains in *CHS* gene models of four plant species.

**Main findings** *CHS* genes with an apparent triplication of the CHS domain encoding part were discovered through database searches. Such genes were found in *Macadamia integrifolia*, *Musa balbisiana*, *Musa troglodytarum*, and *Nymphaea colorata*. A manual inspection of the *CHS* gene models in these four species with massive RNA-seq data suggests that these gene models are the result of artificial fusions in the annotation process. While there are hundreds of seemingly correct CHS records in the databases, it is not clear why these annotation artifacts appeared.

**Keywords** Chalcone synthase, Flavonoid biosynthesis, Annotation error, RNA-seq mapping, Bioinformatics, Domain composition

## Introduction

Flavonoids are one of the most important groups of specialized plant metabolites. Their enormous chemical diversity results in a plethora of biological functions [1]. Most noticeable are the anthocyanins which can provide blue to red coloration to flowers. Flavonoids are distributed across the whole plant kingdom. Given the visual phenotype of several flavonoids, this pathway was established as a model system for plant metabolism and transcriptional regulation [2, 3]. Today, the flavonoid

biosynthesis and its regulation are among the best studied processes in plants. Numerous studies are published every year which investigate the flavonoid biosynthesis and the corresponding regulators in different plant species.

One of the best studied enzymes in the flavonoid biosynthesis pathway is the chalcone synthase (CHS). This enzyme catalyzes the first committed step of the flavonoid biosynthesis, particularly the reaction leading to the formation of the naringenin chalcone from one p-coumaroyl-CoA and three malonyl-CoA [4]. CHS belongs to the type III polyketide synthases, a larger protein family that also harbors several closely related enzymes like the stilbene synthase (STS) [5]. CHS and STS differ by only two functionally important residues which are Q166 and Q167 in the *Arabidopsis thaliana* CHS sequence [6]. The gene structure of *CHS* comprises usually two coding exons [7].

\*Correspondence:

Boas Pucker

b.pucker@tu-braunschweig.de

<sup>1</sup>Department of Biology and Ecology, Faculty of Science, University of Ostrava, Ostrava, Czech Republic

<sup>2</sup>Department of Physics, Faculty of Science, University of Ostrava, Ostrava, Czech Republic

<sup>3</sup>Institute of Plant Biology & BRICS, TU Braunschweig, Braunschweig, Germany



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Numerous plant genomes are sequenced every year [8] and the rapid development of long read sequencing technologies enables the generation of high quality genome sequences [9]. Since an automatic annotation of all genes in the flavonoid biosynthesis and many regulators is possible [10, 11], it is feasible to perform large scale analyses of gene families acting in this pathway. Analyses at this scale are inherently prone to errors concerning individual species and sequences that require human intervention at the final interpretation step. A systematic analysis of domains in the chalcone synthase revealed an apparent triplication in several species.

Here, we describe an investigation of several *CHS* gene models that appeared to encode a protein with a triplicated *CHS* domain. A manual inspection of multiple cases based on transcriptomic data suggests mis-annotation leading to artificially fused gene models.

## Main text

### Methods

#### Identification of sequences with potential *CHS* domain duplications

A BLASTp analysis with the *Nymphae colorata* *CHS* (XP\_049936683.1) as query and nr (containing 545,546,009 sequences) as subject was performed through the NCBI webportal (<https://blast.ncbi.nlm.nih.gov/Blast.cgi?PAGE=Proteins>) on the 2nd of June 2023. The word size was set to 5 and an e-value cutoff of  $10^{-30}$  was applied. Matched sequences longer than 500 amino acid residues were manually screened for *CHS* domain duplications or triplications. Redundant sequences were removed from the list.

#### Evaluation of gene models via RNA-seq read mapping

The genome sequences and corresponding annotations of four species with apparent fusion of multiple *CHS* gene copies were selected for manual inspection: *Musa balbisiana* [12], *Musa troglodytarum* [13], *Macadamia integrifolia* [14], and *Nymphaea colorata* [15]. All available paired-end RNA-seq data sets of these species were retrieved from the Sequence Read Archive via fastq-dump [16] (Additional File 1 in [17]). Read pairs were aligned to the genome sequence with STAR v2.5.1b [18, 19] in 2-pass-mode using a minimal similarity of 95% and a minimal alignment length of 90% as previously described [20]. The resulting BAM files were processed with customized Python scripts [17]. All reads mapped to the locus of interest were extracted from each individual BAM file with samtools v1.15.1 [21]. These subset BAM files were merged to produce one BAM file per species (Additional File 2, Additional File 3, Additional File 4,

Additional File 5 in [17]). The final BAM file was visualized with Integrative Genomics Viewer (IGV) [22].

#### Analysis and visualization of RNA-seq read mapping coverage

A previously developed [23] Python script was applied to convert BAM files into coverage files that list the number of aligned reads per genomic position. These coverage files served as the basis for the generation of gene model-focused coverage plots [17]. These plots visualize the number of aligned reads per position around a given locus of interest. High coverage with RNA-seq reads indicates exon positions while introns are characterized by very low or no RNA-seq read coverage at all. Introns are not included in the final mRNAs and usually mostly mature mRNAs are extracted with standard RNA extraction protocols used for RNA-seq experiments.

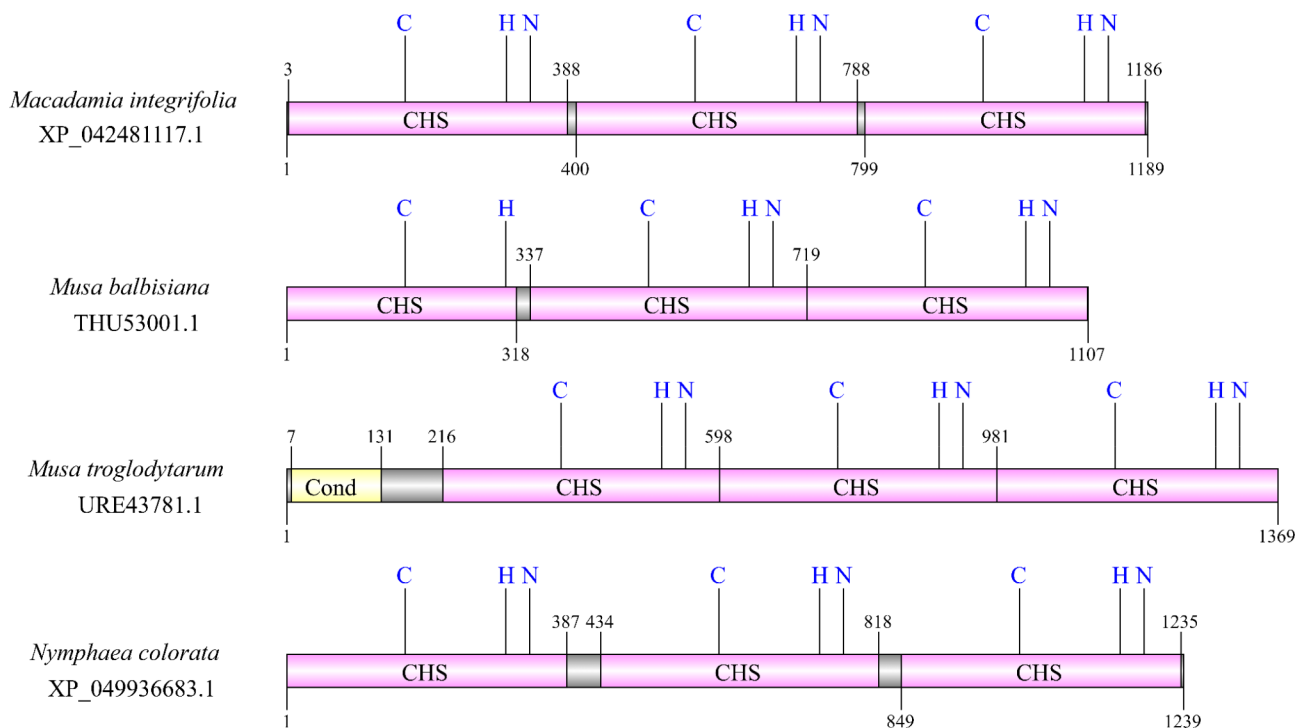
#### Comparison of sequence similarity

Sequences around the *CHS* loci of interest were retrieved from the respective genome sequence using the Python script seqex3.py v0.5 [24]. Sequence similarity between the different domains predicted in the inspected sequences was analyzed by constructing dot plots with dotplotter v0.1 [25] using a k-mer size of 31. Dotplotter compares two sequences by generating k-mers from one given sequence and identifying matches in the other sequence. The corresponding k-mer positions in both sequences are visualized in a dotplot. This allows an intuitive visualization of repeats.

### Results

An analysis of the *CHS* sequences retrieved from the NCBI revealed apparent fusion proteins harboring multiple (two or three) *CHS* domains in 28 protein sequences across different plant species (Additional File 6 in [17]). More specifically, 8 sequences showed a triplication of the *CHS* domain, while 20 sequences showed a duplication. Representative examples of predicted proteins with multiple *CHS* domains were identified in *Macadamia integrifolia*, *Musa balbisiana*, *Musa troglodytarum*, and *Nymphaea colorata* (Fig. 1). A comparison of each *CHS* locus in the four species against itself revealed the expected high similarity (Additional File 7 in [17]).

Loci with an unexpected *CHS* annotation were inspected in an RNA-seq read mapping (Fig. 2). The annotation indicates a single *CHS* gene encoding three *CHS* domains, but the results of our RNA-seq read mapping suggest that there are multiple individual *CHS* genes. The coverage continuously drops towards the ends of several exons. That is an indication of the end of a gene. In contrast, an abrupt drop in coverage to almost zero would indicate a splice site. There are almost no connecting reads between some of the annotated exons and



**Fig. 1** Domain composition of protein sequences with apparent CHS domain triplication. CHS abbreviation in figure stands for *Chalcone synthase domain*, and Cond abbreviation stands for *Chalcone and stilbene synthases, N-terminal domain*. So-called catalytic triad consisting of conserved Cysteine (C), Histidine (H), and Asparagine (N) amino acid residues is always depicted (if preserved)

individual exons show very different RNA-seq coverage. There are also substantial fractions of annotated exons without RNA-seq support. It is not expected that exons at the 5'- and 3'-end of a gene have high coverage, while the enclosed exons have low coverage. The underlying RNA-seq read mapping files of *Macadamia integrifolia* (Additional File 2 in [17]), *Musa balbisiana* (Additional File 3 in [17]), *Musa troglodytarum* (Additional File 4 in [17]), and *Nymphaea colorata* (Additional File 5 in [17]) are available for in depth inspection. Detailed instructions are provided to enable researchers to perform a similar investigation of other cases of potential annotation artifacts [17].

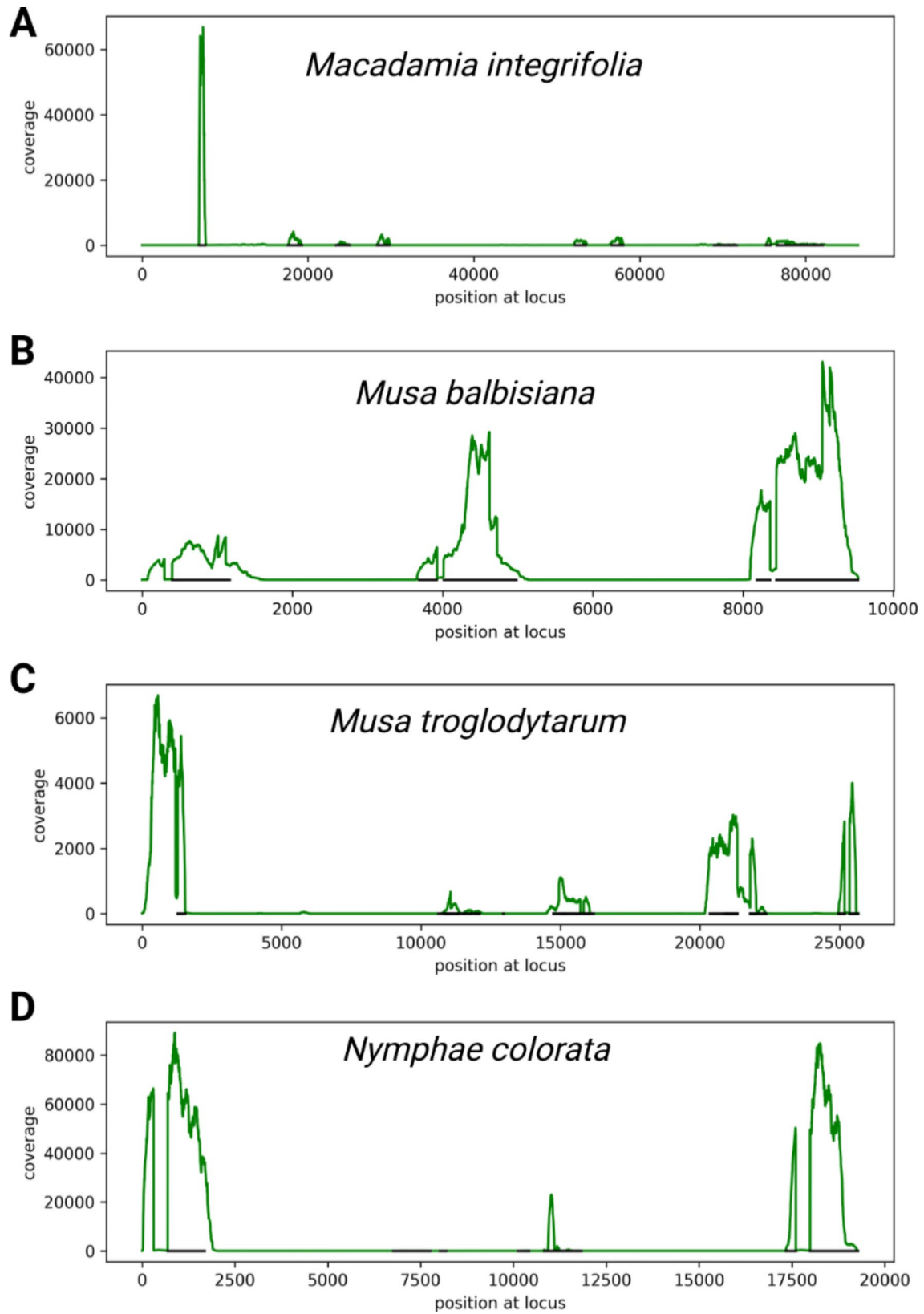
## Discussion

We outlined how surprising results of a quick database search and in-depth inspection of these entries resulted in the identification of most likely annotation artifacts. We investigated the assembly and annotation process underlying the annotations of the four analyzed species. The *Macadamia integrifolia* assembly was produced with MaSuRCA v3.2.6 and ALLMAPS v.Jul-2019 [14, 26, 27]. The *Musa balbisiana* and *Musa troglodytarum* genome sequences were assembled with wtdbg v1.2.8 and NextDenovo v2.4.0, respectively [12, 13, 28, 29]. The *Nymphaea colorata* assembly was produced by Canu v1.3 [15, 30]. All these assemblers are well established tools that

have been deployed in numerous plant genome sequencing projects. The annotation was produced by Gnomon [31] which is the default annotation pipeline operated at the NCBI. There is no indication how any of these steps could have caused the mis-annotation. Gnomon was also applied on many other plant genome sequences and usually generated CHS annotations comprising only a single CHS domain [32–34]. The repetitive regions with multiple CHS copies in an array might contribute to the mis-annotation as suggested for repeats before [35].

To the best of our knowledge, there are no reports that validated a CHS protein with multiple repeated domains. Despite all efforts, it might not be possible to automatically maintain a perfect database free of any mis-annotations. Despite technological advances, manual inspection and correction might be required in some cases [36]. Swiss-Prot is an initiative to establish a collection of sequences that were curated by experts [37, 38]. However, such a manual inspection is not a scalable approach given the rapid growth of sequence collections due to improved sequencing technologies [8, 9]. As most records in the database are likely correct, users should carefully inspect any substantially deviating records. Especially sequences, which appear as exciting discoveries, should be considered as suspicious and thorough inspection is required.

Unexpected sequences should be carefully checked, because annotation artifacts are more likely than striking



**Fig. 2** RNA-seq read coverage of *CHS* loci in *Macadamia integrifolia*(A), *Musa balbisiana*(B), *Musa troglodytarum*(C), and *Nymphae colorata*(D). The annotation suggests that a single gene is spanning the entire displayed locus, but the continuously dropping coverage towards the end of exons (black lines) in our RNA-seq read mapping suggests that there are multiple individual genes

differences between closely related species. We describe our strategy for the gene model investigation in detail to enable repetition by others [17].

- (1) A comparison with orthologs in many other species is often a way to identify unexpected sequence properties. Sequences should only differ systematically if there are particularly striking events during evolution. If such events are not known, any major sequence differences could point to artifacts and should be considered as such.
- (2) Generally, the alignment of RNA-seq reads or other types of transcriptomic evidence should be used to gain insights into the structure of genes. It is important to use a proper split read aligner like STAR [18, 19] or HISAT2 [39] for this step. The coverage should be consistent or should drop continuously towards the ends of a gene. Splice sites within a gene are characterized by abrupt changes of the coverage to almost zero. In addition, introns should be spanned by a number of reads that is almost equivalent to the coverage at the border of the flanking exons.
- (3) If no suitable datasets are available, it is also possible to validate the gene structure and sequence through amplification via RT-PCR followed by Sanger sequencing [40]. As this approach is more time consuming and requires more financial resources, the aforementioned data-reuse approaches should be applied first.

### Limitations

Our investigation was restricted to four examples of apparent CHS domain triplication events within a single *CHS* gene in four different plant species. The gene models in question were contradicted by the analyzed RNA-seq datasets. However, we cannot rule out the possibility that such modified CHS versions exist elsewhere.

### Abbreviations

CHS	Chalcone synthase
STS	Stilbene synthase
ER	Endoplasmic reticulum
IGV	Integrative Genomics Viewer

### Acknowledgements

This work was supported by the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure (de.NBI) (031A532B, 031A533A, 031A533B, 031A534A, 031A535A, 031A537A, 031A537B, 031A537C, 031A537D, 031A538A). Many thanks to the Bioinformatics Resource Facility (BRF) at the Center for Biotechnology (CeBiTec) at Bielefeld University for providing an environment to perform the computational analyses.

### Authors' contributions

MB and BP planned the study. BP wrote the software and performed the bioinformatic analysis. MB, AV, JC and BP interpreted the results, and wrote the manuscript. All authors have read the final version of the manuscript and approved its submission.

### Funding

We acknowledge support by the Open Access Publication Funds of Technische Universität Braunschweig. Open Access funding enabled and organized by Projekt DEAL.

### Data Availability

The developed scripts, a detailed description for the inspection of suspicious gene models, and additional sequence collections are available via GitHub: <https://github.com/bpucker/CHS>. All analyzed data sets are publicly available (Additional File 1 in [17]). Genome sequences and corresponding annotations were retrieved from the BananaGenomeHub (<https://banana-genome-hub.southgreen.fr/>) and the NCBI (<https://www.ncbi.nlm.nih.gov/>), respectively.

### Declarations

#### Competing interests

The authors declare no competing interests.

#### Ethics approval and consent to participate

Not applicable.

#### Consent for publication

Not applicable.

Received: 27 March 2023 / Accepted: 13 June 2023

Published online: 20 June 2023

### References

1. Winkel-Shirley B. Flavonoid biosynthesis. A colorful model for genetics, biochemistry, cell biology, and biotechnology. *Plant Physiol.* 2001;126:485–93.
2. Dubos C, Le Gourrierec J, Baudry A, Huet G, Lanet E, Debeaujon I. MYB2 is a new regulator of flavonoid biosynthesis in *Arabidopsis thaliana*. *Plant J.* 2008;55:940–53.
3. Ramsay NA, Glover BJ. MYB-bHLH-WD40 protein complex and the evolution of cellular diversity. *Trends Plant Sci.* 2005;10:63–70.
4. Dao TTH, Linthorst HJM, Verpoorte R. Chalcone synthase and its functions in plant resistance. *Phytochem Rev.* 2011;10:397–412.
5. Flores-Sanchez IJ, Verpoorte R. Plant Polyketide Synthases: a fascinating group of enzymes. *Plant Physiol Biochem.* 2009;47:167–74.
6. Schröder G, Schröder J. A single change of histidine to glutamine alters the substrate preference of a stilbene synthase. *J Biol Chem.* 1992;267:20558–60.
7. Yang J, Gu H. Duplication and divergent evolution of the CHS and CHS-like genes in the chalcone synthase (CHS) superfamily. *Chin Sci Bull.* 2006;51:505–9.
8. Marks RA, Hotaling S, Frandsen PB, VanBuren R. Representation and participation across 20 years of plant genome sequencing. *Nat Plants.* 2021;7:1571–8.
9. Pucker B, Irisarri I, de Vries J, Xu B. Plant genome sequence assembly in the era of long reads: Progress, challenges and future directions. *Quant Plant Biology.* 2022;3:e5.
10. Rempel A, Pucker B. KIPES3: Automatic annotation of biosynthesis pathways. 2022;2022.06.30.498365.
11. Pucker B. Automatic identification and annotation of MYB gene family members in plants. *BMC Genomics.* 2022;23:220.
12. Wang Z, Miao H, Liu J, Xu B, Yao X, Xu C. *Musa balbisiana* genome reveals subgenome evolution and functional divergence. *Nat plants.* 2019;5:810–21.
13. Li Z, Wang J, Fu Y, Jing Y, Huang B, Chen Y. The *Musa troglodytarum* L. genome provides insights into the mechanism of non-climacteric behaviour and enrichment of carotenoids. *BMC Biol.* 2022;20:186.
14. Nock CJ, Baten A, Mauleon R, Langdon KS, Topp B, Hardner C. Chromosome-scale assembly and annotation of the macadamia genome (*Macadamia integrifolia* HAES 741). *G3: genes, genomes. Genetics.* 2020;10:3497–504.
15. Zhang L, Chen F, Zhang X, Li Z, Zhao Y, Lohaus R. The water lily genome and the early evolution of flowering plants. *Nature.* 2020;577:79–84.
16. NCBI. sra-tools. 2020. <https://github.com/ncbi/sra-tools>.
17. Pucker B. Manual inspection of *CHS* gene models. 2023. <https://github.com/bpucker/CHS>.
18. Dobin A, Gingeras TR. Mapping RNA-seq reads with STAR. *Curr protocols Bioinf.* 2015;51:11–4.

19. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.
20. Haak M, Vinke S, Keller W, Droste J, Rückert C, Kalinowski J. High quality de novo transcriptome assembly of *Croton tiglium*. *Front Mol Biosci*. 2018;5:62.
21. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25:2078–9.
22. Thorvaldsdóttir H, Robinson JT, Mesirov JP. Integrative Genomics viewer (IGV): high-performance genomics data visualization and exploration. *Brief Bioinform*. 2013;14:178–92.
23. Pucker B, Brockington SF. Genome-wide analyses supported by RNA-Seq reveal non-canonical splice sites in plant genomes. *BMC Genomics*. 2018;19:1–13.
24. Pucker B, Schilbert HM, Schumacher SF. Integrating Molecular Biology and Bioinformatics Education. *J Integr Bioinform*. 2019;16.
25. Pucker B. PBBtools v0.1. 2023. <https://github.com/bpucker/PBBtools>.
26. Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL, Yorke JA. The MaSuRCA genome assembler. *Bioinformatics*. 2013;29:2669–77.
27. Tang H, Zhang X, Miao C, Zhang J, Ming R, Schnable JC. ALLMAPS: robust scaffold ordering based on multiple maps. *Genome Biol*. 2015;16:3.
28. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods*. 2020;17:155–8.
29. GrandOmic. NextDenovo. 2023. <https://github.com/Nextomics/NextDenovo>.
30. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*. 2017;27:722–36.
31. Souvorov A, Kapustin Y, Kiryutin B, Chetvernin V, Tatusova T, Lipman D. Gnomon-theNCBIeukaryoticgenepredictiontool.2018.[https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/gnomon/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/gnomon/). Accessed 13Nov2018.
32. Liu Y-J, Wang X-R, Zeng Q-Y. De novo assembly of white poplar genome and genetic diversity of white poplar population in Irtysh River basin in China. *Sci China Life Sci*. 2019;62:609–18.
33. Daccord N, Celton J-M, Linsmith G, Becker C, Choisne N, Schijlen E. High-quality de novo assembly of the apple genome and methylome dynamics of early fruit development. *Nat Genet*. 2017;49:1099–106.
34. Guo L, Winzer T, Yang X, Li Y, Ning Z, He Z. The opium poppy genome and morphinan production. *Science*. 2018;362:343–7.
35. Tørresen OK, Star B, Mier P, Andrade-Navarro MA, Bateman A, Jarnot P. Tandem repeats lead to sequence assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res*. 2019;47:10994–1006.
36. Vaattovaara A, Leppälä J, Salojärvi J, Wrzaczek M. High-throughput sequencing data and the impact of plant gene annotation quality. *J Exp Bot*. 2019;70:1069–76.
37. Boeckmann B, Bairoch A, Apweiler R, Blatter M-C, Estreicher A, Gasteiger E. The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003. *Nucleic Acids Res*. 2003;31:365–70.
38. Schneider M, Tognolli M, Bairoch A. The swiss-prot protein knowledgebase and ExPASy: providing the plant community with high quality proteomic data and tools. *Plant Physiol Biochem*. 2004;42:1013–21.
39. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol*. 2019;37:907–15.
40. Pucker B, Holtgräwe D, Weisshaar B. Consideration of non-canonical splice sites improves gene prediction on the *Arabidopsis thaliana* Niederzenz-1 genome sequence. *BMC Res Notes*. 2017;10:1–6.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.