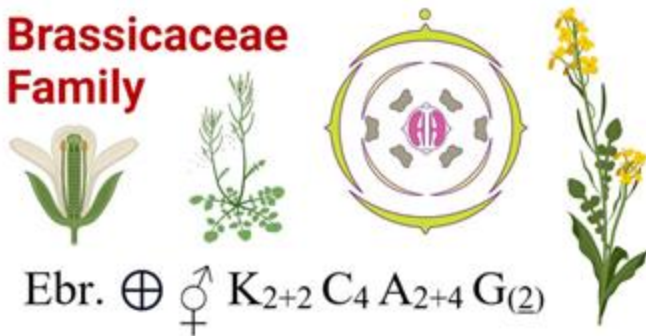


Resolution of Brassicaceae Phylogeny Using Nuclear Genes Uncovers Nested Radiations and Supports Convergent Morphological Evolution

Chien-Hsun Huang,¹ Renran Sun,¹ Yi Hu,² Liping Zeng,¹ Ning Zhang,³ Liming Cai,¹ Qiang Zhang,⁴ Marcus A. Koch,⁵ Ihsan Al-Shehbaz,⁶ Patrick P. Edger,⁷ J. Chris Pires,⁸ Dun-Yan Tan,⁹ Yang Zhong,¹ and Hong Ma^{*,1}

**Brassicaceae
Family**



Ebr. \oplus ♀ K_{2+2} C_4 A_{2+4} $\text{G}_{(2)}$

Molecular Biology and Evolution,
Volume 33, Issue 2, February 2016, Pages 394–412,
<https://doi.org/10.1093/molbev/msv226>

Published: 29 October 2015

Sirine Oueida
Valentin Goupille

M2 Bio-informatique
EMP

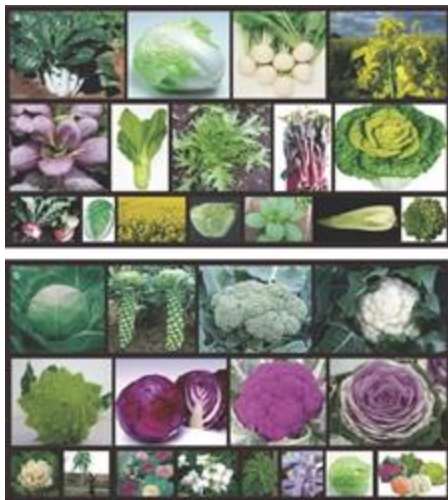


Introduction

Les Brassicacées

- L'une des familles d'angiospermes **les plus diversifiées** et **les plus rentable** sur le plan économique
 - **environ 3700 espèces** (51 tribus, 321 genres)

colza, betterave,
choux, navet,
brocoli, moutarde...



(Chen et al, 2014)

Modèle biologique végétale le plus étudié : ***Arabidopsis thaliana***



Introduction

Les Brassicacées

- **Une excellente famille modèle pour les études comparative et évolutives**
 - Diversité écologique, morphologique, génétique...
 - Histoire évolutive complexe → WGD (Whole Genome Duplication)

Pour faciliter l'utilisation des Brassicacées comme famille modèle évolutive

=> **une phylogénie avec des relations bien étayées est cruciale...**

2006 : Utilisation d'un **marqueur chloroplastique** chez 113 espèces de Brassicacées
=> ont identifié **3 lignées principales**
(avec des supports faibles à modérés)



Review

Cell
PRESS

Cabbage family affairs: the evolutionary history of Brassicaceae

Andreas Franzke¹, Martin A. Lysak², Ihsan A. Al-Shehbaz³, Marcus A. Koch⁴ and Klaus Mummenhoff⁵

¹Heidelberg Botanic Garden, Centre for Organismal Studies Heidelberg, Heidelberg University, D-69120 Heidelberg, Germany

²Department of Functional Genomics and Proteomics, Faculty of Science, Masaryk University, and CEITEC, CZ-602 00 Brno, Czech Republic

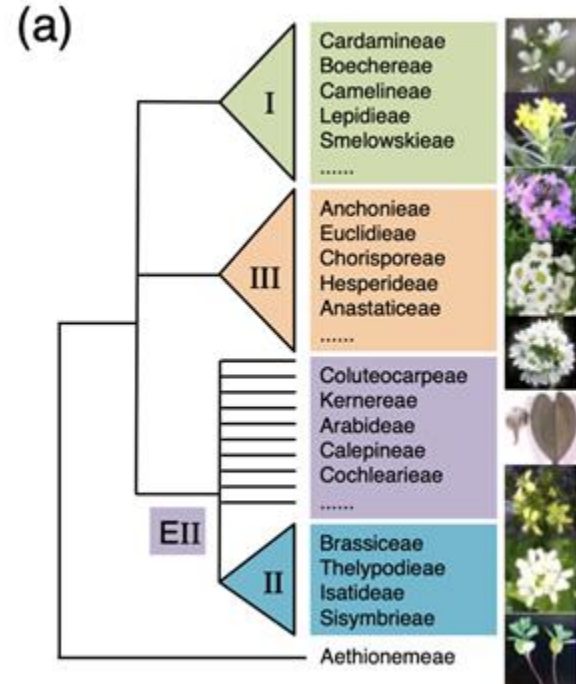
³Missouri Botanical Garden, St. Louis, MO 63166-0299, USA

⁴Biodiversity and Plant Systematics, Centre for Organismal Studies Heidelberg, Heidelberg University, D-69120 Heidelberg, Germany

⁵Biology Department, Botany, Osnabrück University, D-49069 Osnabrück, Germany

2011 : Le groupe de genres supplémentaires associés à LII avec un faible support a été appelé lignée étendue II (EII)

Figure : Phylogénie des Brassicacées
(Franzke et al, 2011)



Introduction

Les relations entre les principales lignées de la famille des Brassicacées restent floues...

Fournir une phylogénie des Brassicacées en utilisant des marqueurs nucléaires

(Huang et al, 2015)

Resolution of Brassicaceae Phylogeny Using Nuclear Genes Uncovers Nested Radiations and Supports Convergent Morphological Evolution

Chien-Hsun Huang,¹ Renran Sun,¹ Yi Hu,² Liping Zeng,¹ Ning Zhang,³ Liming Cai,¹ Qiang Zhang,⁴
Marcus A. Koch,⁵ Ihsan Al-Shehbaz,⁶ Patrick P. Edger,⁷ J. Chris Pires,⁸ Dun-Yan Tan,⁹ Yang Zhong,¹ and
Hong Ma^{*1}

Objectifs :

I- Reconstruction phylogénétique des Brassicacées

II- Estimation des temps de divergences des Brassicacées

III- Reconstruction des caractères ancestraux à partir de traits morphologiques

I- Reconstruction phylogénétique des Brassicacées

Méthodologie :

1. Récupération et production des données génomiques / transcriptomiques
2. Sélection de gènes orthologues putatifs
3. Analyses Phylogénétiques
4. Évaluation de la robustesse phylogénétique et détection des biais

I- Reconstruction phylogénétique des Brassicacées

1) Récupération et production des données :

- Production de 33 transcriptomes :
 - 32 espèces de Brassicaceae
 - *Cleome serrulata*, une Cleomaceae
 - La famille sœur des Brassicaceae en tant que groupe externe
- Augmentées de 10 séquences de génomes entiers obtenues à partir de bases de données publiques et de 13 ensembles de données de transcriptomes supplémentaires



Cleome serrulata

Un total de 55 grands ensembles de données d'espèces de Brassicacées ont été utilisés pour sélectionner un ensemble de gènes marqueurs pour les reconstructions phylogénétiques => couvrant 29 des 51 tribus de Brassicacées

I- Reconstruction phylogénétique des Brassicacées

2) Sélection de gènes orthologues putatifs :

Reconstruction phylogénétique complexe chez les Brassicacées, compliquée par des événements de polyploïdisation, de pertes de gènes...

- Nécessité d'identifier des **gènes orthologues fiables** pour résoudre les relations entre lignées

=> **Trois approches ont été utilisées** pour sélectionner **113 gènes nucléaires à faible nombre de copies** comme marqueurs phylogénétiques...

I- Reconstruction phylogénétique des Brassicacées

3) Analyses phylogénétiques

A-Préparation des séquences

- **Obtention des gènes** : 113 gènes orthologues
- **Alignement des séquences** avec **MUSCLE**
- **Nettoyage des séquences** : élimination des séquences de mauvaises qualités et utilisation de **trimAl**
- **Concaténation** des séquences avec **SeaView**
=> Une **supermatrice** combinée de 113 gènes

I- Reconstruction phylogénétique des Brassicacées

B) Sélection du modèle de substitution :

- **GTRGAMMA :**

GTR (General Time Reversible) : Modèle qui permet des taux de substitution variables entre nucléotides

GAMMA : permet de prendre en compte les variations des taux de substitution entre sites

C) Utilisation de 2 méthodes de reconstruction phylogénétique :

- **Inférence Bayésienne** (BI) avec **MrBayes**
- **Maximum de Vraisemblance** (ML) avec **RAxML**

I- Reconstruction phylogénétique des Brassicacées

4) Évaluation de la robustesse phylogénétique et détection des biais

1. Évaluation de la confiance dans la sélection des arbres :

- **Test de 48 topologies** pour les relations entre les clades principaux de Brassicacées
- Analyse de robustesse avec **RAXML** et **CONSEL** (tests AU, bootstrap, PP, etc)

2. Détection des biais de séquence :

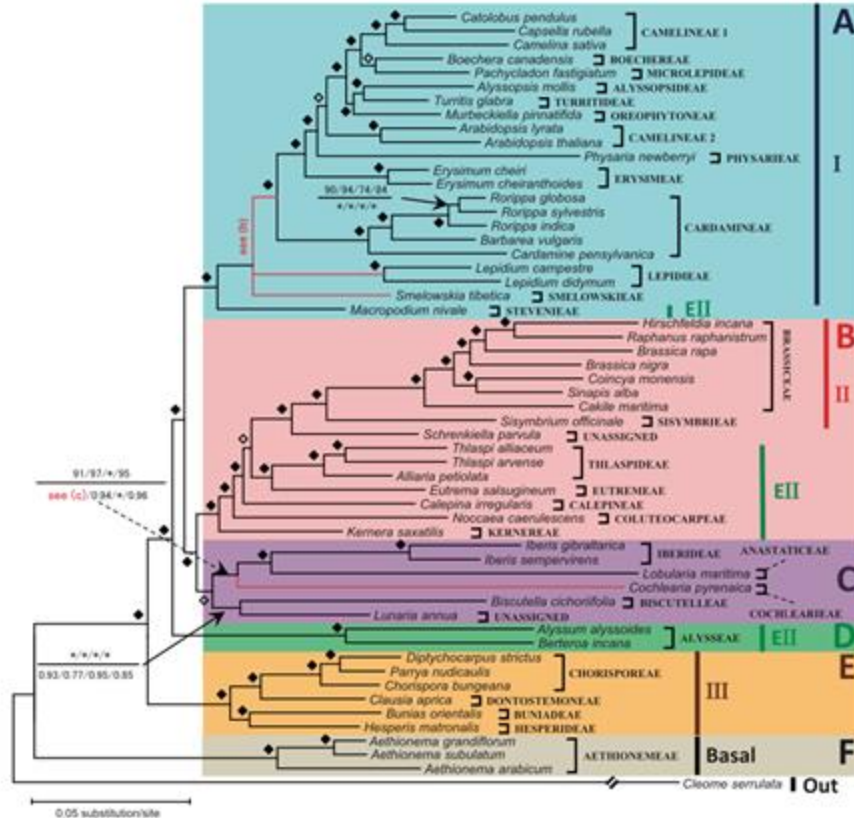
- Analyse des 113 gènes avec **TreSpEx** pour **détecter l'attraction des LB et la saturation**
- **Exclusion des gènes biaisés pour des reconstructions fiables**

3. **Analyse de coalescence** multi-espèces

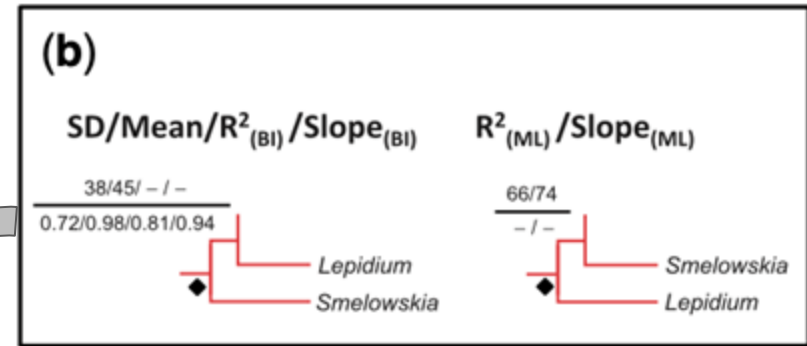
- **Astral** est utilisé pour l'**analyse multilocus** basée sur les 113 arbres génétiques
([Mirarab et al, 2014](#))

I- Reconstruction phylogénétique des Brassicacées

Concaténation



b) Relations entre *Smelowskia* et *Lepidium*



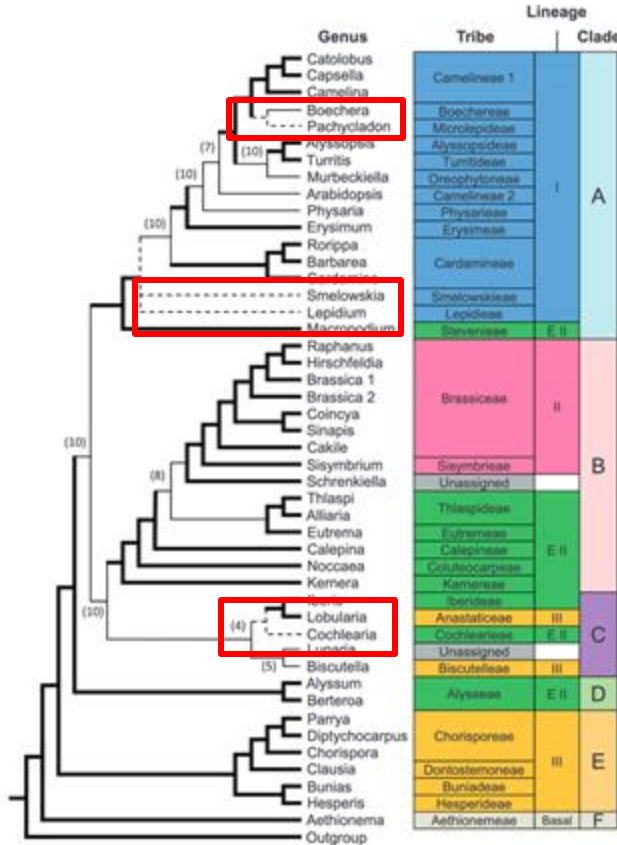
Smelowskia et *Lepidium* sont très proches phylogénétiquement dans l'arbre.

Différence entre les différentes méthodes statistiques utilisés

Figure 1 : Arbre phylogénétique des Brassicaceae basé sur les analyses ML et BI des concaténations d'ensembles de gènes après des tests de biais de séquence.

I- Reconstruction phylogénétique des Brassicacées

Consensus et Différence entre
Concaténation et Coalescence



Ce modèle est proposé selon nos résultats de **concaténation** (en utilisant 113 gènes comme marqueurs, et arbres après tests de biais et méthodes **coalescentes**

- Les **structures cohérentes en traits pleins**
 - Épais (toutes montrant un support maximal)
 - ou fins (avec des valeurs indiquant le nombre de résultats avec un support maximal).
- Les **lignes pointillées indiquent les incertitudes**

6 clades : A, B, C, D, E et F issues issues des méthodes de concaténation et de coalescence.

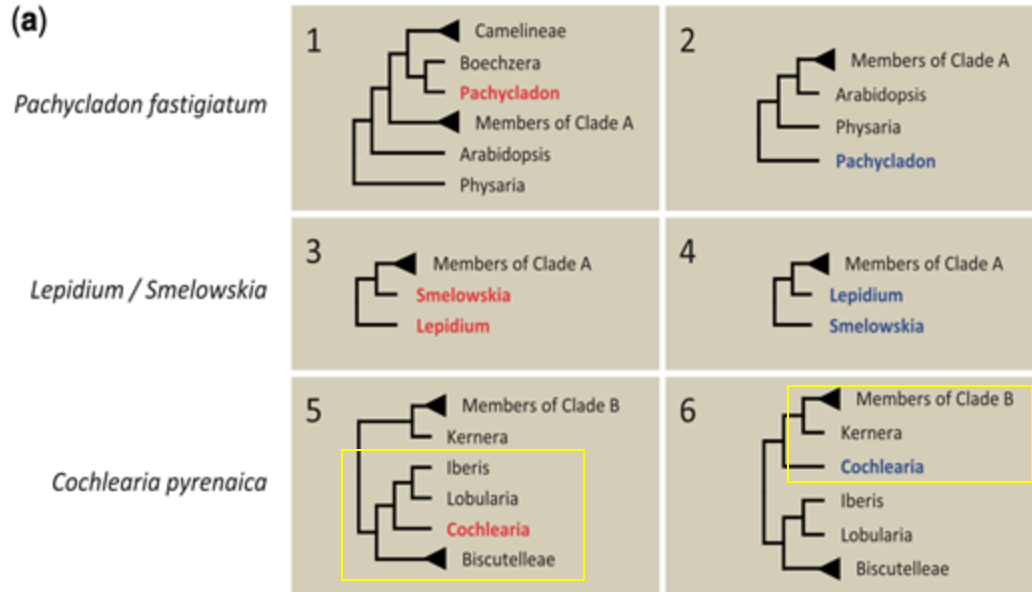
Les clades A, B et C → groupe (ABC) → groupe **monophylétique** → regroupé avec le clade D.

Clade E : Groupe séparé

Clade F (*Aethionemeae*) : lignée basale

Figure 2 : Une phylogénie résumée des Brassicacées

I- Reconstruction phylogénétique des Brassicacées



La position des genres peut varier en fonction des méthodes utilisées

Variations :

- Événements de réticulation
- Hybridation
- Polyploïdisation

Figure 3 : Topologies conflictuelles de certains genres dans cette étude.

II- Estimation des temps de divergences des Brassicacées

Méthode utilisée : Vraisemblance pénalisée (PL) implémentée dans r8s

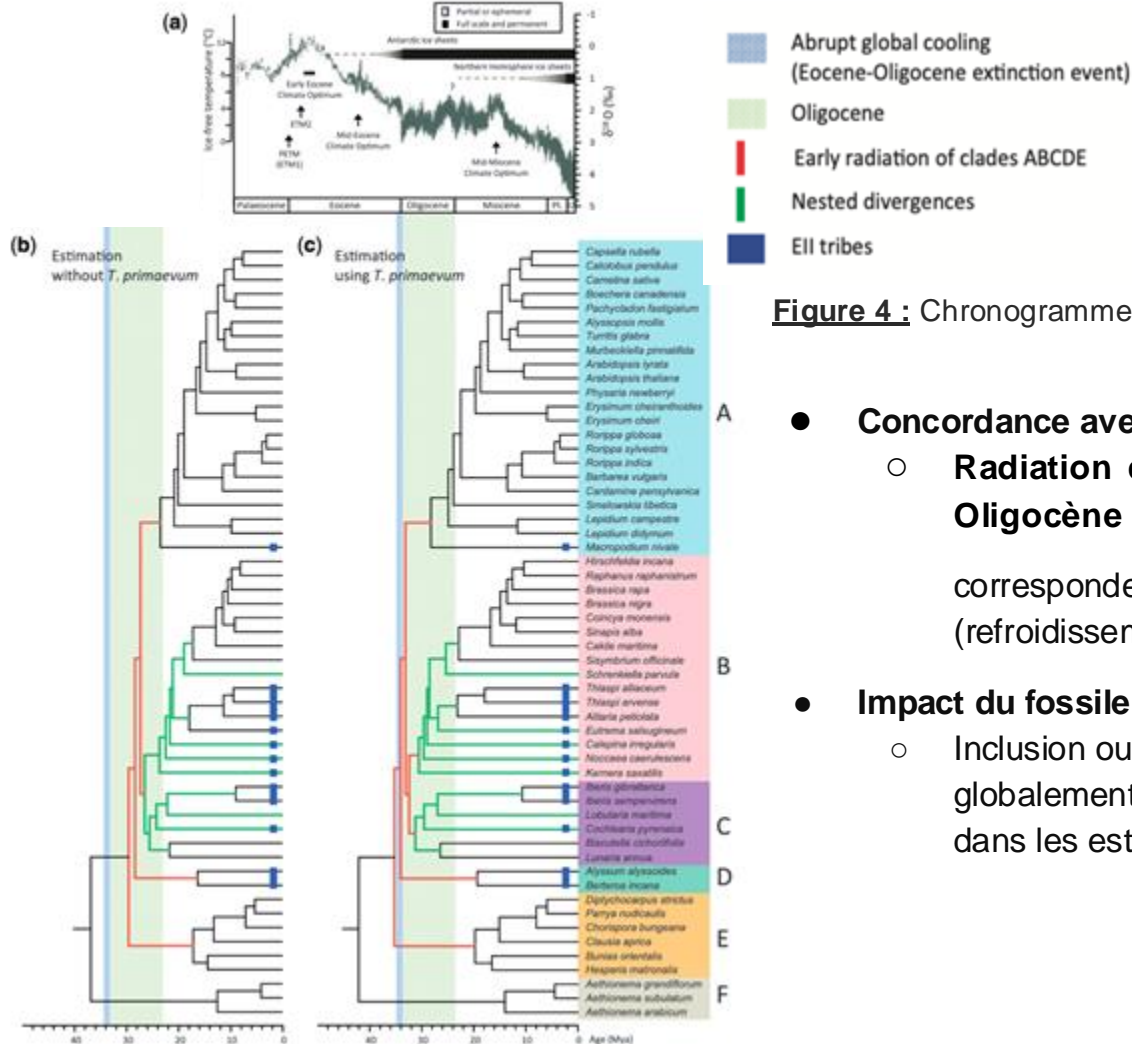
=> Hypothèse de taux variables d'évolution

L'arbre ML reconstruit à l'aide de 113 gènes avec une longueur de branche générée par RAxML a été utilisé comme arbre d'entrée

- **Calibration basée sur des fossiles**

- ***D. bicarpelata* (89,3 Mya)** : Fossile utilisé pour calibrer l'âge du nœud de tige des Brassicales (incluant les Malvaceae comme groupe sœur).
- ***T. primaevum* (23,03 Mya)** : Fossile supposé de Brassicaceae, utilisé pour contraindre l'âge minimum de la divergence entre *Thlaspi* et *Alliaria*.
- **16 autres étalonnages** ont été utilisés en dehors du groupe des Brassicaceae pour estimer l'âge des groupes voisins

- **Estimations réalisées avec et sans fossile incertain (*T. primaevum*)**



II- Estimation des temps de divergences des Brassicacées

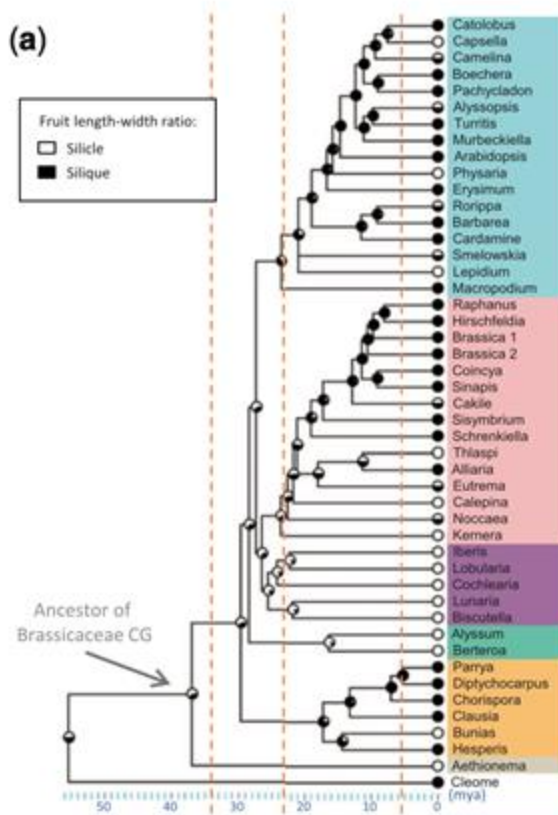
Figure 4 : Chronogramme des Brassicacées déduit à l'aide de r8s

- **Concordance avec le paléoclimat :**
 - **Radiation des Brassicacées → période Eocène-Oligocène**

correspondent aux changements climatiques du passé (refroidissement)
- **Impact du fossile :**
 - Inclusion ou non de *T. primaevum* : Résultats globalement cohérents : Différences minimales (~5 Ma) dans les estimations des temps de divergence

III- Reconstruction des caractères ancestraux à partir de traits morphologiques

Analyse des caractères morphologiques :



Reconstruction des états ancestraux pour plusieurs caractères morphologiques à l'aide de **Mesquite**

Silicule : **Silique** dont la longueur n'excède pas 3 fois la largeur

La reconstruction des caractères ancestraux révèle de **nombreux événements d'évolution convergente**

Figure 5 : Schémas évolutifs du rapport longueur-largeur des fruits

Conclusion

I- Reconstruction phylogénétique des Brassicacées

- Utilisation des marqueurs nucléaires
- Utilisation des méthodes de concaténation et de coalescence
- Une **phylogénie des Brassicacées largement confirmée**

→ Reste des incertitudes → **Incompréhension partielle** de la phylogénie globale

II- Estimation des temps de divergences des Brassicacées

- Changement climatique a joué un rôle dans la diversification des Brassicacées

III- Reconstruction des caractères ancestraux à partir de traits morphologiques

- La reconstruction des caractères ancestraux révèle de nombreux événements **d'évolution convergents**

Limites et perspectives

Limites :

- **Échantillonnage incomplet** : 19 tribus non représentées.
- **Incompréhension partielle** de la phylogénie globale.

Perspectives :

- **Espèces comme *Pachycladon* et *Cochlearia*** nécessitent plus de données pour clarifier leur **histoire évolutive**.
- **Augmenter l'échantillonnage** des espèces et des gènes pour **améliorer la compréhension** des relations évolutives.

Resolving the backbone of the Brassicaceae phylogeny for investigating trait diversity

Lachezar A. Nikolov¹ , Philip Shushkov², Bruno Nevado³ , Xiangchao Gan¹, Ihsan A. Al-Shehbaz⁴ , Dmitry Filatov³ , C. Donovan Bailey⁵ and Miltos Tsiantis¹

¹Department of Comparative Development and Genetics, Max Planck Institute for Plant Breeding Research, Cologne 50829, Germany; ²Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA 91125, USA; ³Department of Plant Sciences, University of Oxford, Oxford, OX1 3RB, UK; ⁴Missouri Botanical Garden, 4344 Shaw Boulevard, St Louis, MO 63110, USA; ⁵Department of Biology, New Mexico State University, Las Cruces, NM 88003, USA

Summary

Author for correspondence:
Miltos Tsiantis
Tel: +49 22 15062105
Email: tsiantis@mpipz.mpg.de

Received: 9 August 2018
Accepted: 10 January 2019

New Phytologist (2019) **222**: 1638–1651
doi: 10.1111/nph.15732

Key words: anchored phylogenomics, comparative transcriptomics, crucifers, leaf evolution, targeted sequence capture.

- The Brassicaceae family comprises *c.* 4000 species including economically important crops and the model plant *Arabidopsis thaliana*. Despite their importance, the relationships among major lineages in the family remain unresolved, hampering comparative research.
- Here, we inferred a Brassicaceae phylogeny using newly generated targeted enrichment sequence data of 1827 exons (> 940 000 bases) representing 63 species, as well as sequenced genome data of 16 species, together representing 50 of the 52 currently recognized Brassicaceae tribes. A third of the samples were derived from herbarium material, facilitating broad taxonomic coverage of the family.
- Six major clades formed successive sister groups to the rest of Brassicaceae. We also recovered strong support for novel relationships among tribes, and resolved the position of 16 taxa previously not assigned to a tribe. The broad utility of these phylogenetic results is illustrated through a comparative investigation of genome-wide expression signatures that distinguish simple from complex leaves in Brassicaceae.
- Our study provides an easily extendable dataset for further advances in Brassicaceae systematics and a timely higher-level phylogenetic framework for a wide range of comparative studies of multiple traits in an intensively investigated group of plants.

Article

Global Brassicaceae phylogeny based on filtering of 1,000-gene dataset

Kasper P. Hendriks,^{1,2,38,39,*} Christiane Kiefer,³ Ihsan A. Al-Shehbaz,⁴ C. Donovan Bailey,⁵ Alex Hooft van Huysduynen,^{2,6} Lachezar A. Nikolov,⁷ Lars Nauheimer,⁸ Alexandre R. Zuntini,⁹ Dmitry A. German,¹⁰ Andreas Franzke,¹¹ Marcus A. Koch,³ Martin A. Lysak,¹² Óscar Toro-Núñez,¹³ Banş Özüdoğru,¹⁴ Vanessa R. Invernón,¹⁵ Nora Walden,³ Olivier Maurin,⁹ Nikolai M. Hay,¹⁶ Philip Shushkov,¹⁷ Terezie Mandáková,¹² M. Eric Schranz,¹⁸ Mats Thulin,¹⁹ Michael D. Windham,¹⁶ Ivana Rešetnik,²⁰ Stanislav Španiel,²¹ Elfy Ly,^{2,22,23} J. Chris Pires,²⁴ Alex Harkess,²⁵ Barbara Neuffer,¹ Robert Vogt,²⁶

SUMMARY

The mustard family (Brassicaceae) is a scientifically and economically important family, containing the model plant *Arabidopsis thaliana* and numerous crop species that feed billions worldwide. Despite its relevance, most phylogenetic trees of the family are incompletely sampled and often contain poorly supported branches. Here, we present the most complete Brassicaceae genus-level family phylogenies to date (Brassicaceae Tree of Life or BrassiToL) based on nuclear (1,081 genes, 319 of the 349 genera; 57 of the 58 tribes) and plastome (60 genes, 265 genera; all tribes) data. We found cytonuclear discordance between the two, which is likely a result of rampant hybridization among closely and more distantly related lineages. To evaluate the impact of such hybridization on the nuclear phylogeny reconstruction, we performed five different gene sampling routines, which increasingly removed putatively paralog genes. Our cleaned subset of 297 genes revealed high support for the tribes, whereas support for the main lineages (supertribes) was moderate. Calibration based on the 20 most clock-like nuclear genes suggests a late Eocene to late Oligocene origin of the family. Finally, our results strongly support a recently published new family classification, dividing the family into two subfamilies (one with five supertribes), together representing 58 tribes. This includes five recently described or re-established tribes, including Arabidopsidae, a monogeneric tribe accommodating *Arabidopsis* without any close relatives. With a worldwide community of thousands of researchers working on Brassicaceae and its diverse members, our new genus-level family phylogeny will be an indispensable tool for studies on biodiversity and plant biology.

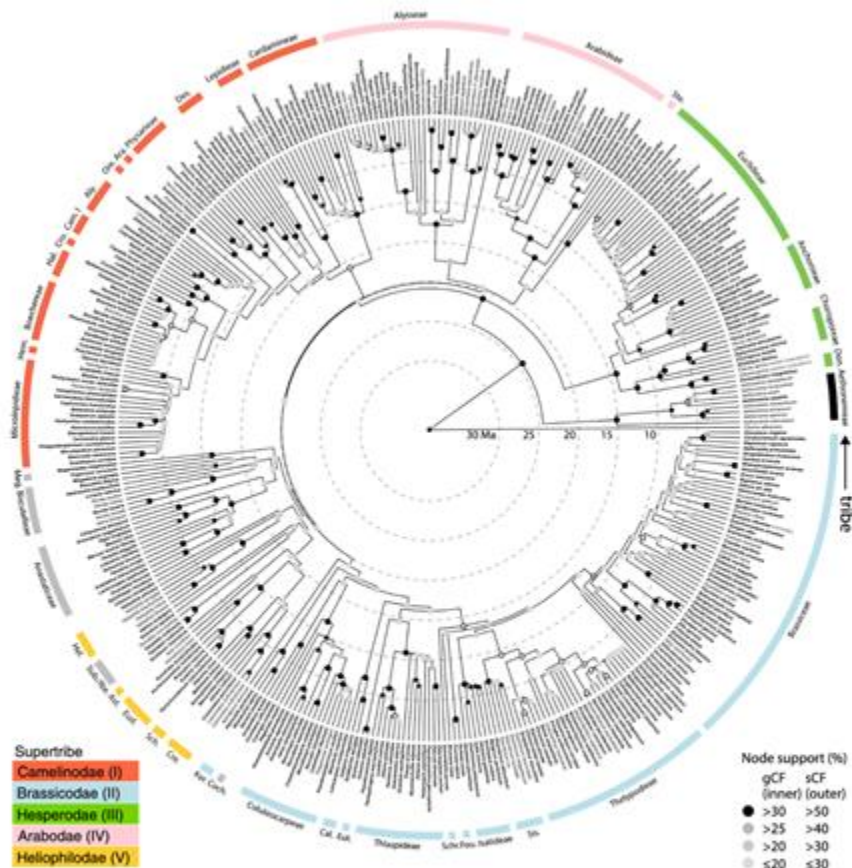


Figure 2. Time-calibrated genus-level Brassicaceae Tree of Life (BrassiToL) from a maximum likelihood analysis of a 297 nuclear genes supermatrix (superstrict routine)

Genus type species are highlighted in bold. All tribes with more than a single representative are listed. Abbreviations of tribes are as follows: Aly., Alysiaceae; Ara., Arabidoideae trib. nov.; Ast., Asteae; Cal., Calepineae; Cam. I, Camelinoideae I; Coch., Cochleariaceae; Ore., Cruciferales; Des., Descurainiaceae; Don., Drostostemonaceae; Eud., Eudemeae; Eut., Eutemeae; Fou., Fourcraeae; Hal., Halimolobae; Hel., Heliophila; Hem., Hemiphragmaceae; Ibe., Iberideae; Kar., Kerneraceae; Mag., Megacarpaceae; Ore., Orophytoneae; Sch., Schizopetalaceae; Schr., Schrenkiaeae trib. nov.; Sis., Sisymbrieae; Ste., Steuriaceae; Sub., Subulaceae. See also [Data S1A](#) for a fully annotated version of this phylogeny (including bootstrap, gCF, sCF and node age 95% HPD intervals) with outgroups representing all families within the order Brassicales and calibration nodes. See also [Data S1A](#).

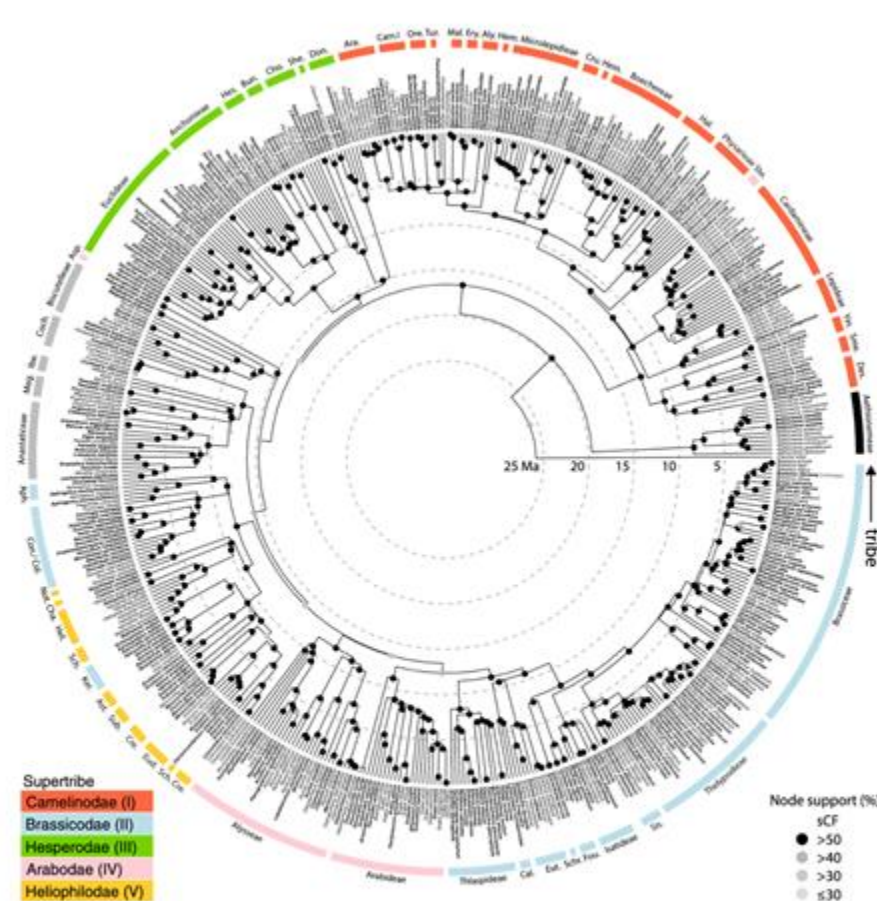


Figure 3. Time-calibrated genus-level Brassicaceae Tree of Life from a maximum likelihood analysis of a 60 plastome genes supermatrix. Genus type species are highlighted in bold. All tribes with more than a single representative are listed. Abbreviations of tribes follow those of Figure 2, with additionally: Aph., Aphragmaceae; Asp., Asperuginaceae trib. nov.; Bun., Bunodiaceae; Cha., Chamiraceae; Cho., Choripontaceae; Col., Coluteocarpaceae; Con., Conringiaceae; Don., Drostostemonaceae; Ery., Erysimeae; Hes., Hesperideae; Mal., Malcoideae; Not., Notothlaspiaceae; Ore., Orophytoneae; She., Shebaziaceae; Sme., Smelowskaceae; Tur., Turridaceae; Yin., Yinshaniaceae. See also [Data S1B](#) for a fully annotated version of this phylogeny (including bootstrap, sCF and node age 95% HPD intervals) with outgroups representing all families within the order Brassicales and calibration nodes. See also [Data S1B](#).

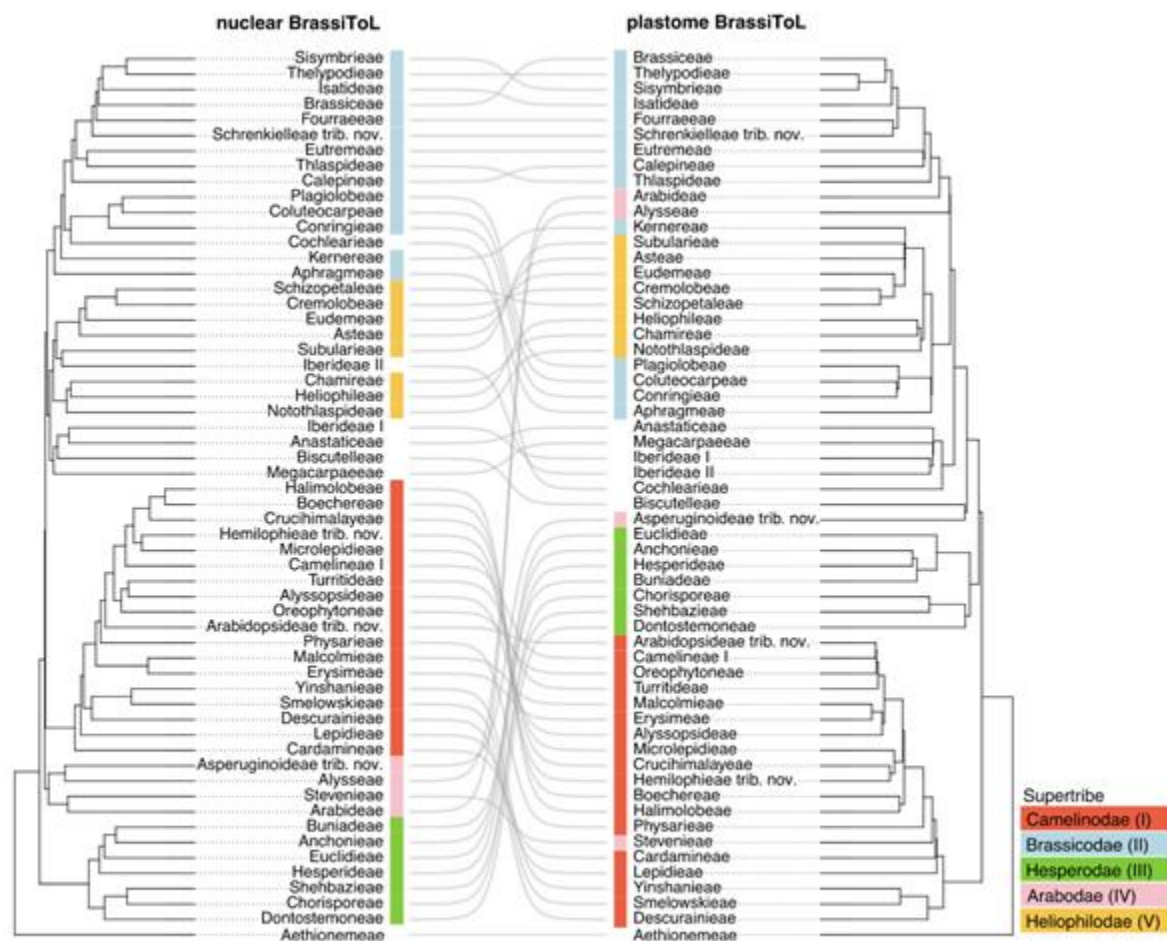


Figure 4. Cytonuclear discordance at tribe level in newly derived nuclear and plastome Brassicaceae Trees of Life

Curved lines between tip labels from the two phylogenies link the same tribes. Tribes are represented by a randomly chosen sample from each tribe. Rogue tribes have not yet been assigned to a supertribe due to their changing topological position from different sampling routines and phylogenetic approaches.

Table 2. Brassicaceae divergence time estimates

Studies	Crown Brassicaceae (Ma)	Crown core Brassicaceae ^a (Ma)	Method	Dataset	Calibration
Koch et al. ⁷⁹	–	25.9–23.1	synonymous substitution rate	<i>Adh</i> and <i>Chs</i>	synonymous substitution rate
Franzke et al. ¹⁸	35.0–15.0–1.0	28.0–11.0–1.0	BEAST	<i>nad4</i>	one secondary calibration
Beilstein et al. ¹⁹	64.2–54.3–45.2	54.3–46.9–39.4	BEAST	<i>ndhF</i> and <i>PHYA</i>	four fossils
Couvreur et al. ²⁰	49.4–37.6–24.2	43.8–32.3–20.9	BEAST	eight genes from nuclei, chloroplast, and mitochondria	one fossil
Kagale et al. ⁸⁰	–	26.6	synonymous substitution rate	213 nuclear orthologs	synonymous substitution rate
Edger et al. ¹¹	45.9–31.8–16.8	–	BEAST	1,115 single-copy nuclear genes	two fossils
Hohmann et al. ⁹	38.6–32.4–27.1	27.3–23.4–19.9	BEAST	plastomes	four fossils
Huang et al. ²²	37.8–37.1–36.3	30.3–29.7–29.1	r8s	113 low-copy nuclear orthologs	18 fossils ^b
Cardinal-McTeague et al. ⁵⁴	44.1–37.7–31.4	–	BEAST	Chloroplast DNA (<i>ndhF</i> , <i>matK</i> , <i>rbcl</i>) and mitochondrial DNA (<i>matR</i> , <i>rps3</i>)	three fossils ^b
Mohammadin et al. ⁸¹	58.9–48.0–37.5	35.4	BEAST	plastomes	one secondary calibration
Guo et al. ⁸²	41.8–34.9–29.0	29.8–25.1–21.3	MCMCTree	plastomes	14 fossils ^c
Mandáková et al. ³⁹	54.7–40.1–29.4	30.6	BEAST	plastomes	four fossils
Huang et al. ²³	33.2–29.9–26.8	22.9–21.3–19.6	BEAST	plastomes	four fossils
Ramírez-Barahona et al. ⁸³	52.7–41.9–30.5	–	BEAST	<i>rbcl</i> , <i>atpB</i> , <i>matK</i> , <i>ndhF</i> , 18S, 26S, 5.8S	238 fossils, angiosperm-wide
Walden et al. ²	35.7–29.9–24.3	29.6–25.1–20.9	BEAST	plastomes	four fossils
Legalov et al. ⁸⁴	–	>36.4	indirect calibration via phytophagous beetles	–	<i>Ceutorhynchus</i> beetle fossils
This study (nuclear dataset)	25.7–24.5–23.1	22.4–21.1–19.9	treePL	phylogeny using superstrict routine; calibration using subset of 20 most clock-like genes	one fossil validated by nine secondary calibration points
This study (plastome dataset)	29.0–20.2–13.0	24.3–16.9–10.2	treePL	plastomes	one fossil validated by eight secondary calibration points

Comparison of estimates from past studies and the current study (taken and updated from Huang et al.²³).

^aCore Brassicaceae are all Brassicaceae, excluding basal tribe Aethionomeae.

^bDating results excluding *Thlaspi primaevum*.

^cOnly results that exclude Brassicales fossils.

Merci de votre attention

Des questions?



I- Reconstruction phylogénétique des Brassicacées

2) Sélection de gènes orthologues putatifs :

Méthodes de Sélection des Gènes Orthologues putatifs

- **Méthode n°1** : Sélection de **28 gènes** issus d'analyses précédentes et validés par phylogénie monogénique dans 13 Brassicacées.
- **Méthode n°2** : Sélection de **45 gènes** parmi deux bases de données d'orthologues couvrant 9 Angiospermes, en vérifiant leur conservation dans les Brassicacées.
- **Méthode n°3** : Sélection de **40 gènes** orthologues à copie unique conservés entre 5 espèces de Brassicacées qui sont également conservés chez les Angiospermes.

113 gènes combinés issus des trois groupes, formant une base robuste pour des analyses phylogénétiques...

Méthodes de Sélection des Gènes Orthologues putatifs

Groupe de gènes	Origine des données	Critères de sélection	Résultat
Premier groupe	Études précédentes (Zhang et al. 2012 ; Zeng et al. 2014) utilisant HaMStR et OrthoMCL sur 9 génomes d'angiospermes.	- Reconstruction de phylogénies monogéniques avec RAxML.- Regroupement des espèces par lignées.- Exclusion des duplications dues aux polyploïdisations fréquentes dans les Brassicaceae.	28 gènes retenus.
Deuxième groupe	Chevauchement entre 2 bases de données : 1. 4 180 OG (HaMStR).2. 1 989 OG (OrthoMCL) chez 7 angiospermes.	- Partagés par les deux bases de données.- Couverture > 80 % des 13 espèces.- Protéines > 300 acides aminés.- Regroupement cohérent des lignées dans les arbres monogéniques.	45 gènes retenus.
Troisième groupe	Chevauchement entre :1. 4 180 OG (HaMStR).2. 6 552 OG (OrthoMCL) à copie unique dans 5 Brassicaceae.	- Copie unique conservée dans les angiospermes et Brassicaceae.- Couverture > 80 % des 13 espèces.- Protéines > 500 acides aminés.- Regroupement cohérent selon les trois lignées.	40 gènes retenus.
Total (3 groupes)	Combinaison des trois approches.	- Analyse rigoureuse des phylogénies monogéniques.- Maximisation de la qualité et de la pertinence des marqueurs génétiques retenus.	113 marqueurs

Méthode 1 : Gènes orthologues conservés chez les angiospermes

Objectif : Identifier des gènes nucléaires à faible nombre de copies conservés chez les angiospermes et les adapter aux Brassicacées.

Procédure :

1. Les gènes ont été récupérés à l'aide de deux outils :
 - **HaMStR** (Ebersberger et al. 2009).
 - **OrthoMCL v1.4** (Li et al. 2003).
2. Ces outils ont permis d'identifier 106 gènes partagés par neuf espèces d'angiospermes ayant un génome séquencé (*Z. mays*, *S. or. bicolor*, *O. sativa*, *G. max*, *Me. truncatula*, *P. o. trichocarpa*, *V. vinifera*, *A. thaliana* et *Sol. lycopersicum*).
3. Les homologues de ces gènes ont ensuite été récupérés dans les génomes de 13 espèces de Brassicacées.
4. Chaque gène a été évalué à l'aide de phylogénies monogéniques (avec RAXML), en vérifiant que les espèces de la même lignée (LI, LII, ou LIII) étaient correctement regroupées, sans exiger des relations spécifiques entre ou au sein des lignées.

Résultat : 28 gènes ont été retenus comme marqueurs orthologues fiables.

Méthode 2 : Orthogroupes à faible nombre de copies partagés par les angiospermes

Objectif : Élaborer un deuxième ensemble de marqueurs en utilisant des bases de données indépendantes pour diversifier la sélection.

Procédure :

1. Deux bases de données ont été combinées :

- Une contenant 4 180 groupes de gènes nucléaires orthologues (OG) identifiés par HaMStR à partir de neuf espèces d'angiospermes.
- Une autre avec 1 989 OG à faible nombre de copies identifiés par OrthoMCL à partir de sept espèces d'angiospermes (*A. thaliana*, *P. trichocarpa*, etc.).

2. Les OG partagés par ces deux bases de données ont été considérés comme conservés et à faible nombre de copies.

3. Critères supplémentaires pour sélectionner des gènes fiables :

- Présence dans au moins 80 % des 13 espèces de Brassicacées.
- Codage de protéines de plus de 300 acides aminés.

4. Les arbres monogéniques de ces gènes ont été analysés pour vérifier que les espèces se regroupent correctement selon les trois lignées des Brassicacées.

Résultat : 45 gènes ont été retenus pour ce deuxième groupe.

Méthode 3 : Gènes conservés chez les angiospermes et les Brassicacées

Objectif : Identifier des orthologues à copie unique dans les Brassicacées, qui sont également conservés chez les angiospermes.

Procédure :

1. Les orthologues ont été extraits de deux ensembles de données :
 - Les 4 180 OG identifiés par HaMStR (voir Méthode 2).
 - Les 6 552 OG identifiés par OrthoMCL comme des orthologues à copie unique dans cinq Brassicacées (*A. lyrata*, *B. rapa*, etc.).
2. Critères supplémentaires pour sélectionner des gènes fiables :
 - Présence dans plus de 80 % des 13 espèces représentatives de Brassicacées.
 - Codage de protéines de plus de 500 acides aminés.
3. Les arbres monogéniques ont été vérifiés pour garantir que les espèces se regroupent correctement selon les trois lignées des Brassicacées.

Résultat : 40 gènes ont été retenus pour ce troisième groupe.

1. Préparation des séquences

- **Obtention des gènes :**

Les 113 gènes orthologues ont été récupérés à l'aide de **HaMStR**, avec une valeur seuil d'E < e-20 (gènes significativement similaires).

- **Alignement :**

Les séquences nucléotidiques ont été alignées en utilisant **MUSCLE** (un logiciel d'alignement de séquences) avec ses paramètres par défaut.

- **Nettoyage des alignements :**

- Les séquences de mauvaise qualité ont été supprimées.
- Les régions mal alignées ont été réduites à l'aide de **trimAl**, garantissant des alignements fiables.

- **Concaténation :**

Les alignements des gènes individuels ont été concaténés pour former une seule matrice à l'aide de **SeaView**.

- **Gestion des duplications :**

Pour les gènes avec plusieurs copies dans une espèce, la copie avec la branche la plus courte (la plus proche de l'orthologie supposée) a été conservée pour éviter les paralogues.

a) Inférence bayésienne (BI) - MrBayes :

- **Paramètres de l'analyse :**

- Utilisation de **MrBayes** avec 5 000 000 générations.
- Application de la méthode **MCMC (Markov Chain Monte Carlo)** pour estimer la probabilité des arbres :
 - 1 chaîne froide et 3 chaînes chaudes pour explorer l'espace des arbres phylogénétiques.
 - Un seuil d'arrêt automatique (**stopval = 0,01**) pour garantir la convergence.

- **Suivi de la convergence :**

La convergence des chaînes a été contrôlée avec **AWTY**, un outil permettant de visualiser si les chaînes ont atteint une solution stable.

- **Échantillonnage des arbres :**

- Les arbres ont été échantillonnés toutes les 100 générations.
- Les 25 % premiers arbres (phase de burn-in) ont été exclus.
- Les arbres restants ont été utilisés pour générer un **arbre de consensus**.

b) Maximum de vraisemblance (ML) - RAxML :

- **Paramètres de l'analyse :**

- Les analyses ont été réalisées avec **RAxML**, initialement avec le modèle **GTRGAMMA**, mais la topologie et les valeurs bootstrap étaient équivoques.
- Pour optimiser l'efficacité, le modèle **GTRCAT** (approximation plus rapide de GTRGAMMA) a été utilisé pour des analyses supplémentaires.

- **Bootstrap :**

Les scores de bootstrap (évaluant la robustesse des branches) ont été calculés, bien que les résultats initiaux aient montré des ambiguïtés.