# Caffe-IN-e or OUT: The phylogenetic analysis and comparative genomics of *C. canephora, C. arabica and C. humblotiana*

Anne Valerie Bao[1,*], Juan Antonio Deña [1], Robert Andre Nullas[1,*], Alessandra Franchesca Pastrana[1], and Jane Abigail Santiago[1]

[1] *Department of Biology, De La Salle University, Manila*
*\*Corresponding Authors: anne_bao@dlsu.edu.ph, robert_nullas@dlsu.edu.ph*

**Abstract:** Caffeine is a metabolite that plays an important role in coffee, influencing the sensory and physiological effects of the beverage. With the extensive usage of *C. arabica*, the *Coffea* genus has expanded into a diverse group. In contrast, caffeine-free coffee may be found in *Coffea humblotiana*, which lacks the caffeine synthase gene. A comparative investigation was carried out to determine the evolutionary relationship between the variations in caffeine production among these three *Coffea* species: *Coffea arabica, Coffea canephora, and Coffea humblotiana*. Phylogenetic analysis was used to determine the relationship of these three species using the maturase K multiple sequence analysis and the genome comparisons. Phylogenetic trees, maximum likelihood trees, and dot plots show a consensus when using various workflows in Galaxy. The results show that *C. humblotiana* has a closer evolutionary relationship to *C. canephora* compared to *C. arabica*. In addition, the prospects of commercialization and GMO potential given its unique characteristic of being caffeine-free should be further studied.

**Key Words:** Caffeine synthesis, phylogenetic analysis, genome, comparative genome analysis

## 1. INTRODUCTION

Almost everyone's daily ritual requires coffee, and without coffee beans, coffee would still be able to be made. The source of coffee, a coffee bean, is a seed from the *Coffea* plant. It is the seed that is within the purple or crimson fruit. This fruit, which has a pip like a cherry, is frequently called a coffee cherry. The coffee beans are called beans despite not being beans because of their similarity to actual beans. Usually, the fruits have two flat-sided stones inside of them.

Similar to all other plants, coffee beans too have several types, each with a unique taste and purportedly originating from distinct places. Most people in the world consume *C. arabica*, which is one of the most well-known coffee beans. Ethiopia is the natural home of *Coffea arabica*, which possesses the majority of the species' genetic diversity. Historians think that coffee seeds were originally transported to Yemen, where they were cultivated as a crop, from the coffee woods of southwest Ethiopia. Farmers and breeders have chosen and developed dozens of widely cultivated varieties of Arabica coffee from these early plants, each with its distinct performance and adaptation to regional circumstances. The primary seeds carried from Ethiopia to Yemen were related to the Bourbon and Typica varieties, according to recent genetic testing. Originating in Yemen, descendants of Bourbon and Typica migrated all over the world, serving as the foundation for the majority of contemporary *C. arabica* coffee plantations. Although *C. arabica* is the most well-known kind of coffee bean, Coffea

eugenioides is one of its ancestors. *C. canephora* is native from Central Africa to Gulf of Guinea and Uganda, within equatorial forests at elevations of 50 to 1500 meters above sea level. Compared to the previously described coffee beans, *Coffea humblotiana* is a lesser-known coffee bean. It is the only species of Coffee native to the Comoro Islands. It was most likely eaten and maybe even planted on Grande Comore, one of the Mayotte archipelago's adjacent islands, however there is still a dearth of information on this topic. Because of the local agricultural area growth, this species is currently considered endangered. Less than 110 trees are thought to still exist on Mayotte Island, and it is unknown if these trees can be found on the other islands in the archipelago. *C. humblotiana* is a member of the 124 recognized species in the Coffee genus. Its natural range includes tropical Africa, Madagascar, Comoros, Mauritius, and the Reunion Islands. It also extends to Australasia and southern and southeast Asia. Together with other Madagascan Coffee species, it forms a sizable monophyletic group that is predicted to be 11 times larger than the African Coffee species.

Caffeine, being the most consumed stimulant in the world and provided by coffee beans, is missing in one of the *Coffea* species: *Coffea humblotiana*. A wild, endangered species belonging to the Comoro archipelago, has the unique feature of having complete absence of caffeine in seeds and leaves (Raharimalala et al., 2021). This absence is attributed to the lack of the caffeine synthase gene, particularly the DXMT gene responsible for the synthesis of theobromine, a precursor of caffeine. This discovery poses the possibility of utilizing CRISPR-Cas9 to knock out the DMXT gene to produce caffeine-free coffee as discussed in Leibrock et al.'s study (2022). Additionally, investigating the reasons for loss of caffeine for *Coffea humblotiana* could potentially branch into the many ecological factors affecting this such as oxygen, heat, predators, etc.

Conducting phylogenetic analysis to compare the genetic relationships and evolutionary history of *Coffea humblotiana* with other caffeine-containing *Coffea* species, could shed light on establishing its phylogenetics, and potential genetic factors underlying the absence of caffeine in *Coffea humblotiana*. Insights that stem from this can produce broad implications beyond coffee production. The significance of this study is that it may aid in better understanding caffeine synthesis and its mechanism (Zhao et al., 2023).

Moreover, it establishes the genetic basis of presence and absence of caffeine. This could also benefit the agricultural aspects of coffee production wherein various factors may be explored which affects the caffeine content of the *Coffea* species, i.e., manipulating caffeine content of coffee plants (Duygu Ağagündüz et al., 2023). Economically, this research holds promise in the coffee industry as it can provide information towards optimizing caffeine content as this offers more coffee varieties.

With understanding the need to investigate upon *Coffea humblotiana*'s unique feature of not producing caffeine, understanding its phylogenetic history is crucial. The study investigates the phylogenetic analysis of the species, focusing on gene sequence comparison between *Coffea humblotiana* and its related species, specifically, *Coffea canephora* and *Coffea arabica*. By exploring its evolutionary background, the study aims to identify the notable genes concerning caffeine synthesis and pave the way for caffeine content-manipulated coffee plants, a deeper understanding of caffeine, and a more holistic Coffea phylogeny. To accomplish this, the following objectives will be addressed. Firstly, gather the gene sequences and genomes of each species of *Coffea humblotiana, Coffea canephora, and Coffea arabica*. Afterwards, analyze the phylogeny of the gene sequences through various workflows in Galaxy. Lastly, to identify the relationship between the ancestry of the variations and the genomic assembly between the species

## 2. METHODOLOGY

### 2.1 Data Retrieval

The sequences of Maturase K (matK) and RefSeq genome assemblies from these specific *Coffea* species; *Coffea arabica, Coffea canephora, and Coffea humblotiana* were gathered from various papers through NCBI GenBank.

Table 1. Retrieved matK sequences for phylogenetic analysis

| Species | NCBI GenBank Accession Number | Author/Source |
|---|---|---|
| *Coffea arabica* | OP321037.1 | Tapaca, 2022 |
| | OP321036.1 | Tapaca, 2022 |
| | OP321035.1 | Tapaca, 2022 |
| | OP321034.1 | Tapaca et al., 2022 |

| Species | | |
|---|---|---|
| | OP321033.1 | Tapaca et al., 2022 |
| | OP321032.1 | Tapaca et al., 2022 |
| | OP321031.1 | Tapaca et al., 2022 |
| | OP321030.1 | Tapaca et al., 2022 |
| | OP321028.1 | Tapaca et al., 2022 |
| | OP321027.1 | Tapaca et al., 2022 |
| *Coffea* | AB973198.1 | Nakagawa, 2014 |
| *canephora* | AB973197.1 | Nakagawa, 2014 |
| | MK722264.1 | Panaligan et al., 2019 |
| | MK722263.1 | Panaligan et al., 2019 |
| | MK722262.1 | Panaligan et al. 2019 |
| | MK722258.1 | Panaligan et al., 2019 |
| | MF350105.1 | Zuniga et al., 2017 |
| | KC758282.1 | Constantino, 2013 |
| | KC758283.1 | Constantino, 2013 |
| | MK722259.1 | Panaligan et al., 2019 |
| *Coffea* | KJ815707.1 | Kainulainen & |
| *humblotiana* | | Bremer, 2014 |

Table 2. Retrieved NCBI/RefSeq reference genome assemblies

| Species | NCBI Datasets Accession Number | Author/Source |
|---|---|---|
| *Coffea canephora* | GCA_036785865.1 | Coffee Consortium, 2024 |
| *Coffea arabica* | GCF_003713225.1 | Johns Hopkins University, 2018 |
| *Coffea humblotiana* | GCA_023065735.1 | Institut de Recherche pour le Developpement, 2022 |

## 2.2 Molecular Evolutionary Genetics Analysis (MEGA)

Due to the discrepancy in the difference of the nucleotide sequence length of the matK gene of the *C. humblotiana* data, in relation to the limitation of the available data available on the available, credible repositories (NCBI), the sequences were aligned via the Molecular Evolutionary Genetics Analysis or MEGA to create a multiple sequence alignment (MSA) via Muscle alignment and deletion/trimming of the gaps and excess ends.

## 2.3 Phylogenetic Analysis

*Galaxy Australia* (https://usegalaxy.org.au/) was utilized for the phylogenetic analysis. A feature of the Galaxy hub is the "*Galaxy Training!*" wherein online training materials (e.g., tutorials and workflows) are deposited and readily available. One of the features is the *GTN Pan-Galactic Workflow Search* which is a search browser within the *Galaxy Training*, allowing access to publicly available workflows from all available Galaxy servers spanning from UseGalaxy.eu, UseGalaxy.org.au, UseGalaxy.fr., and WorkflowHub.eu. Already existing phylogenetic analysis workflow was then retrieved as a basis and customized according to the available data. The resulting workflow is as follows:
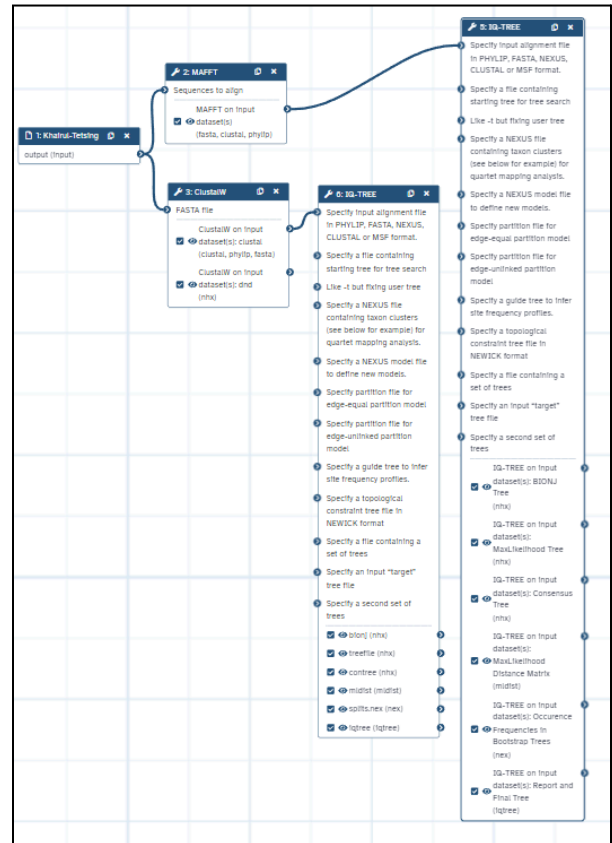


Figure 1. UseGalaxy.au phylogenetic analysis workflow

Workflow tools consist of MAFFT and ClustalW, two different multiple alignment (MSA) programs for sequences (amino acid or nucleotide), IQ-TREE, for a

phylogenomic construction of an evolutionary tree from the extracted MSA sequences, and Newick Display, for the visualization of phylogenetic trees. However, the study also utilized the use of iTOL (interactive tree of life) (https://itol.embl.de/), a different website from Galaxy which allows annotation and various displays for the phylogenetic trees.

## 2.4 Genome Comparison

Genome comparison was performed using Galaxy Australia. The Galaxy Australia method utilizes a public workflow. The workflow involved the tool 'Chromeister', a pairwise chromosome comparison tool for small and large-scale genomes.



Figure 2. UseGalaxy.au chromosome comparison workflow

Two genome sequences in FASTA format are compared to each other. The sequence of interest was used as the query sequence and was compared to the reference sequence, which was determined to be *Coffea humblotiana.* The workflow will generate a comparison dot plot, detected events dot plot, and comparison score.

## 3.  RESULTS AND DISCUSSION

*Coffea arabica* and *Coffea canephora (Robusta)* are the two primary *coffea* species cultivated in the coffee market worldwide. The global trade is

significantly divided between the two species, with *C. arabica* at 60% and the remaining 40% attributed to the *C. canephora (Robusta)* (Davis, et al. 2019). *Coffea arabica* is described to have a more flavor profile with variations in the scented notes and aroma, while *Coffea canephora (Robusta)* is more known for its succinct bitterness (Wang et al., 2020). Moreover, a distinct difference between the two lies as well on its growth and cultivation such as the required elevations and temperature for the respective species' growth, hence a trade-off between the qualities in consideration for the global market. *Coffea humblotiana* or also referred to as "*Caféier de Humblot"* is a Coffea species endemic to the Comoro archipelago in Africa, it is one of the recorded caffeine-free *Coffea* species alongside, *Coffea pseudozanguebariae* found in Kenya, and the *Coffea charrieriana* found in Cameroon, all of which are found in the African continent. *Coffea humblotiana* particularly was studied due to the lack of identified caffeine synthase genes (Raharimalala et al., 2021). Unfortunately, the International Union for Conservation of Nature (IUCN), as of its latest assessment (2017), *Caféier de Humblot* is categorized as an "endangered" species under the criteria B1 ab(iii,v) (Chadburn & Davis, 2019). Such that limited data and documentation (e.g., data in repositories) are available for public use.
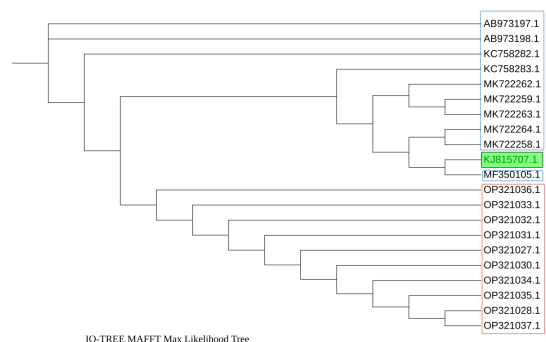
3.1. Phylogenetic Analysis



Figure 3. MAFFT MaxLikelihood Tree visualized via iTOL  (Mode: branch lengths = off) (*C. canephora* identified in blue, and *C. arabica* in red)
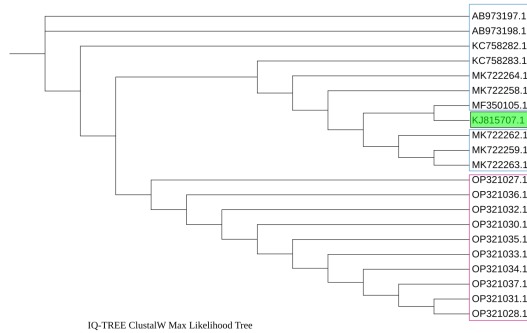
Figure 4. ClustalW MaxLikelihood Tree visualized via iTOL (Mode: branch lengths = off) (*C. canephora* identified in blue, and *C. arabica* in red)

Both alignment tools, MAFFT and ClustalW were able to show similar phylogenetic relationships as evidenced by the separation of the *C. canephora* clade and the *C. arabica* clade as visualized in Figures 3 and 4. In both analysis, the evolutionary relation remained the same wherein, *C. arabica* samples have diverged from the *C. canephora* lineage which is expected as the *C. arabica* is a result of a polyploidy mutation, allotetraploid from its progenitors, *C. canephora* and *C. eugenioides* (Charr et al., 2020). *C. humblotiana* (KJ815707.1) exhibited a closer phylogenetic relationship towards the *C. canephora* clade, particularly with sample MF35015.1, observed in both topologies. From the results, no significant differences between the clades were observed aside from the difference in the evolutionary lineage of some of the *C. canephora* samples which can be attributed with the difference in the algorithmic properties of either tools. Further assessments were made in relevance via genomic comparison.
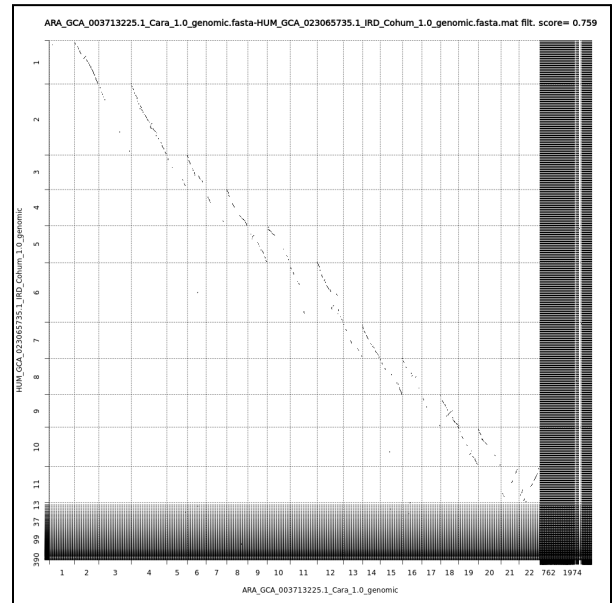
3.2. Chromosome Comparison



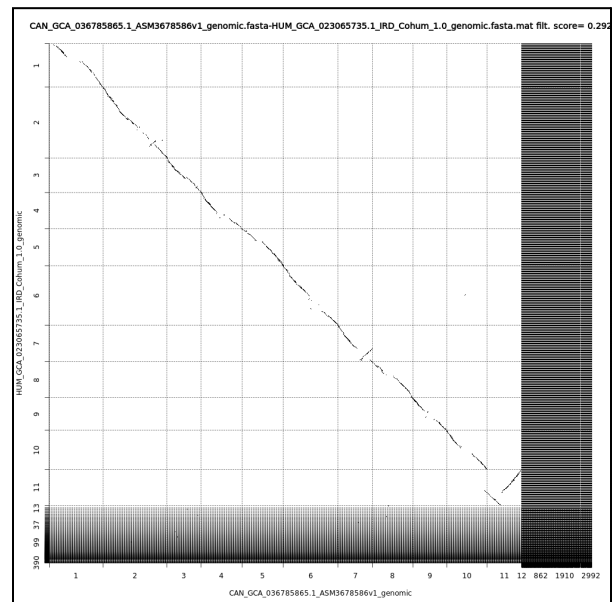Figure 5. Dot plot and Comparison Score for *Coffea arabica and Coffea humblotiana*



Figure 6. Dotplot and Comparison Score for *Coffea canephora and Coffea humblotiana*

Pairwise genome comparison was performed between the genomes of *C. arabica* and *C. canephora* with *C. humblotiana* as the reference genome sequence.

*C. humblotiana* and *C. canephora* yielded a score of 0.292 while *C. humblotiana* and *C. arabica* yielded 0.759. A score of 0 indicates similarity whilst the score of 1 or approximately 1 indicates high to absolute similarity (Pérez-Wohlfeil et al., 2019); higher collinearity between the genome sequence of *C. canephora* and *C. humblotiana*, lower collinearity between *C. arabica* and *C. humblotiana.* The results further support the earlier phylogenetic analysis via maximum likelihood tree wherein, *C. humblotiana* is closer in relation to *C. canephora,* implying a properly conserved synteny between the two species.

In addition, the dot plot for *C. canephora* appears to have a mostly continuous match. There are some frameshift mutations present with the gaps within various parts of the figure. Lastly, there are inverted repeat sequences, represented by the lines going in the opposite direction. This could indicate the presence of complementary inverted repeats. In comparison, the dot plot for *C. arabica* shows a lot of frameshift insertions and mutations.

The production of GMO or Genetically Modified Organisms generally incorporates the process wherein it is initiated through the identification of a target gene of interest to be isolated and constructed with the chosen recipient genome. Hybridization also follows the same processes which leans more towards the amplification of a target gene. However, in propagating a possible new variant/cultivar or GMO with a decaffeinated status in mind does not follow the same process. Instead of amplification of a target gene, the ideal application is through the utilization of knock-out methods in gene suppression or simply, removing the gene region that transcribes for the required biological compounds, a good example of this technology is the CRISPR/Cas9 genome editing technology (Movahedi et al., 2023)[1]. The caffeine biosynthesis of the *Coffea* species requires the presence of several genes which generally code for the enzymes, specifically the NMT genes (N-methyltransferases). The absence of these genes in the *C. humblotiana* is the reference in regards to a possible, feasible, knock-out application (Raharimalala et al., 2021). A comparison between a well conserved reference and query genome can offer to shed some light on the said application especially on possible after-effects or challenges. While the advancements in molecular technology continue to update, challenges in relation to its effects in the long run continue to plague. Concerns regarding a potential off-target application resulting in an unwanted mutagenesis; unwanted alteration may affect the product not only in regards to its quality but possibly affect the development of the GMO product as a whole (Movahedi et al., 2023). As such, biosafety concerns as well as the social-perception are important factors to be considered as well, besides the biological data.

## 4. CONCLUSIONS

The phylogenetic analysis of matK sequences and comparative genomics of various *Coffea* species with *C. humblotiana* has exhibited results that there is a close ancestry between *C. humblotiana and C. canephora.* However, since there is a limitation to the availability of matK sequences for *C. humblotiana,* further sequencing, and available data in the future could indicate a more detailed phylogenetic relationship within these species.

Knowing the limitation of only using genomes for phylogenetic analysis, it is highly recommended for future studies to isolate its RNA and the desired gene sequences related to caffeine synthase. Conducting RNA-Seq on *C. humblotiana*, for example, would allow for a more complete examination of its genetic expression and pathways. Furthermore, not only does it provide valuable data on this wild species, but it also offers the potential to fully understand the reasons for its unique property of caffeine absence and to successfully compare it with other species for a more comprehensive phylogenetic analysis. With that in mind, *C. humblotiana* presents itself as feasible as an ideal GMO for commercialization when given further studies with the lack of genes for caffeine biosynthesis.

## 5. REFERENCES

Chadburn, H. & Davis, A.P. 2017. Coffea humblotiana. The IUCN Red List of Threatened Species 2017: e.T108652718A108665565. https://dx.doi.org/10.2305/IUCN.UK.2017-3.RLTS.T108652718A108665565.en. Accessed on 31

March 2024.

Clayton, L. (2021, October 27). *What Is Eugenioides Coffee?* Sprudge Coffee. https://sprudge.com/what-is-eugenioides-coffee-181142.html

Cuff, W. R., Venkata R.S.K. Duvvuri, Liang, B., Duvvuri, B., Wu, G. E., Wu, J., & Raymond S.W. Tsang. (2010). A Novel Interpretation of Structural Dot Plots of Genomes Derived from the Analysis of Two Strains of Neisseria meningitidis. *Genomics, Proteomics & Bioinformatics*, *8*(3), 159–169. https://doi.org/10.1016/s1672-0229(10)60018-6

Davis AP, Chadburn H, Moat J, O'Sullivan R, Hargreaves S, Lughadha EN. 2019. High extinction risk for wild coffee species and implications for coffee sector sustainability. Science Advances 5:eaav3473.

Duygu Ağagündüz, Gülsüm Deveci, Elif Çelik, Özge Cemali, Teslime Özge Şahin, Ayşe Derya Bayazıt, Özbay, S., & Makbule Gezmen Karadağ. (2023). Determination of Caffeine Amount of Organic and Conventional Coffee. *Records of Agricultural and Food Chemistry, S.* https://doi.org/10.25135/rfac.16.2303.0004

*FIGURE 3 | Coffee (Coffea arabica and C. canephora). (A) Coffee (C. . ..* (n.d.). ResearchGate. https://www.researchgate.net/figure/Coffee-Coffea-arabica-and-C-canephora-A-Coffee-C-arabica-berries-B-Roasted_fig3_347256910.

Griffiths, R. R., Bigelow, G. E., Liebson, I. A., O'Keeffe, M., O'Leary, D., & Russ, N. (1986). HUMAN COFFEE DRINKING: MANIPULATION OF CONCENTRATION AND CAFFEINE DOSE. *Journal of the Experimental Analysis of Behavior*, *45*(2), 133–148. https://doi.org/10.1901/jeab.1986.45-133

Goldemberg, D. C., Antônio, A. G., Farah, A., & Maia, L. C. (2015, January 1). *Coffea canephora*. Elsevier eBooks. https://doi.org/10.1016/b978-0-12-409517-5.00069-3

*History of Arabica*. (2023, June 9). World Coffee Research. https://varieties.worldcoffeeresearch.org/arabica-2/history-of-arabica.

Jean-Claude Charr, Garavito, A., Christophe Guyeux, Crouzillat, D., Descombes, P., Fournier, C., Ly, S. N., Raharimalala, E. N., Jean-Jacques Rakotomalala, Piet Stoffelen, Janssens, S., Hamon, P., & Guyot, R. (2020). Complex evolutionary history of coffees revealed by full plastid genomes and 28,800 nuclear SNP analyses, with particular emphasis on Coffea canephora (Robusta coffee). *Molecular Phylogenetics and Evolution*, *151*, 106906–106906. https://doi.org/10.1016/j.ympev.2020.106906

Leibrock, N. V., Santegoets, J., Mooijman, P. J. W., Yusuf, F., Zuijdgeest, X. C. L., Zutt, E. A., Jacobs, J. G. M., & Schaart, J. G. (2022). The biological feasibility and social context of gene-edited, caffeine-free coffee. *Food Science and Biotechnology*, *31*(6), 635–655. https://doi.org/10.1007/s10068-022-01082-3

Raharimalala, N., Rombauts, S., McCarthy, A., Garavito, A., Orozco-Arias, S., Bellanger, L., Morales-Correa, A. Y., Froger, S., Michaux, S., Berry, V., Metairon, S., Fournier, C., Lepelley, M., Mueller, L., Couturon, E., Hamon, P., Rakotomalala, J.-J., Descombes, P., Guyot, R., & Crouzillat, D. (2021). The absence of the caffeine synthase gene is involved in the naturally decaffeinated status of Coffea humblotiana, a wild species from Comoro archipelago. *Scientific Reports*, *11*(1), 8119. https://doi.org/10.1038/s41598-021-87419-0

Maurin, O., Davis, A. P., Chester, M., Mvungi, E. F., Jaufeerally-Fakim, Y., & Fay, M. F. (2007). Towards a Phylogeny for Coffea (Rubiaceae): Identifying Well-supported Lineages Based on Nuclear and Plastid DNA Sequences. Annals of Botany, 100(7), 1565–1583. https://doi.org/10.1093/aob/mcm257

Movahedi, A., Soheila Aghaei-Dargiri, Li, H., Qiang Zhuge, & Sun, W. (2023). CRISPR Variants for Gene Editing in Plants: Biosafety Risks and Future Directions. *International Journal of Molecular Sciences*, *24*(22), 16241–16241. https://doi.org/10.3390/ijms242216241

Selvaraj, D., Rajeev Kumar Sarma, & Ramalingam Sathishkumar. (2008). Phylogenetic analysis of chloroplast matK gene from Zingiberaceae for plant DNA barcoding. *Bioinformation*, *3*(1), 24–27. https://doi.org/10.6026/97320630003024

Toparslan, E., Karabag, K., & Bilge, U. (2020). A workflow with R: Phylogenetic analyses and visualizations using mitochondrial cytochrome b gene sequences. *PLOS ONE*, *15*(12), e0243927–e0243927. https://doi.org/10.1371/journal.pone.0243927

Xiuju Wang, Loong-Tak Lim, Yucheng Fu, Review of Analytical Methods to Detect Adulteration in Coffee, Journal of AOAC INTERNATIONAL, Volume 103, Issue 2, March-April 2020, Pages 295–305, https://doi.org/10.1093/jaocint/qsz019

Zhao, L., Wei, J., Hu, Y., Pi, D., Jiang, M., & Lang, T. (2023). Caffeine Synthesis and Its Mechanism and Application by Microbial Degradation, A Review. *Foods*, *12*(14), 2721–2721. https://doi.org/10.3390/foods12142721