

1 **Brief Communication**

2 **Whole Chloroplast Genomes reveals the uniqueness of Bolivian native cacao**

3 **(*Theobroma cacao*) from the northern part of Bolivia**

4 Gumiel M<sup>1a</sup>, Rollano-Peñaloza OM<sup>2</sup>, Peralta-Rivero C<sup>1\*</sup>, Tejeda L<sup>2</sup>, Palma V<sup>2</sup>, Cartagena P<sup>1</sup>,

5 Mollinedo P<sup>2</sup>, Peñarrieta JM<sup>2\*</sup>

6 1. c, Calle Claudio Peñaranda # 2706 Esq. Vincenti, Plaza España – Sopocachi, La Paz, Bolivia.

7 2. Department of Chemistry, Instituto de Investigaciones Químicas, Universidad Mayor de

8 San Andrés UMSA, Calle 27 y Andrés Bello s/n Cota Cota, La Paz, Bolivia.

9 a. Current address: Universidad Privada Franz Tamayo, Calle Heroes del Acre, 1855, esquina

10 Landaeta. La Paz, Bolivia.

11 \* Corresponding authors: [jmpenarrieta1@umsa.bo](mailto:jmpenarrieta1@umsa.bo) and [cperalta@cipca.org.bo](mailto:cperalta@cipca.org.bo)

12 **ABSTRACT**

13 We report the complete chloroplast sequences of two varieties of *Theobroma cacao* collected in  
14 the Bolivian Amazonia using Next-Generation Sequencing. Comparisons made between these  
15 two chloroplast genomes and the Belizean reference plastid genome identified 19 and 22  
16 nucleotide variants. The phylogenetic analysis reported three main *T. cacao* clades belonging  
17 to the Forastero, Criollo and Trinitario groups. The Bolivian Native Cacao varieties were  
18 located inside the Trinitario group forming their unique branch. The Bolivian Native Cacao  
19 branch reveals a possible new subpopulation different from the well-characterized *T. cacao*  
20 subpopulations. The phylogenetic trees showed that the relationships among the *T. cacao*  
21 varieties were consistent with their geographical locations placing the Cacao Center of Origin  
22 in Western Amazon. The data presented here will contribute to the usage of ultrabarcoding to  
23 distinguish different *T. cacao* varieties and to identify native cacaos from introduced cacaos.  
24 Thus helping in the conservation of local native varieties of *T. cacao*.

25 **Keywords:** *Theobroma cacao*, Genome, chloroplast, evolution, ultrabarcoding

## 26 INTRODUCTION

27 *Theobroma cacao* is a tree cultivated in tropical and subtropical regions of the world. In Bolivia,  
28 *T. cacao* is cultivated in the humid regions of Beni, Pando, La Paz, Cochabamba and Santa Cruz  
29 departments and is a source of economical sustainability for families. The final product of the  
30 processed cacao is chocolate, cocoa butter or cocoa powder. In Bolivia as well as other Latin  
31 American countries, numerous commercial cacao varieties were introduced by the  
32 governments, and mixed with the Bolivian native cacao varieties endangering such native  
33 species (Bazoberry Chali and Salazar Carrasco, 2008). The high biodiversity of cacao varieties  
34 found in Bolivia is little studied, and therefore more research efforts needs to be done. The  
35 main interest lies into characterize and identify the phylogenetic relationships between these  
36 subpopulations and aid in the accurate subspecies identification.

37 The classification of *T. cacao* has been traditionally divided in: Criollo, Forastero and  
38 Trinitario (Cheesman, 1944). The product derived from the varieties corresponding to the  
39 Criollo is considered of being the best quality, whereas the Forastero varieties present  
40 phenotypes that are more disease resistant. On the other hand, Trinitario varieties are  
41 supposed to be originated from hybridization between Criollo and Forastero groups, and  
42 presents the best characteristics of both lineages. Vast genetic analyses using different  
43 molecular markers revealed a huge number of genetic groups (de la Cruz et al., 1995) localized  
44 in the Amazonia region and in Central America. Motamayor et al. (2008) proposed ten  
45 genetically differentiated groups: Amelonado, Contamana, Criollo, Curaray, Guianna, Iquitos,  
46 Marañon, Nacional, Nanay and Purus. Cornejo et al. (2018), corroborated with a population  
47 genomic analyses that these ten groups underwent strong domestication. Moreover, they  
48 found that the Criollo, Amelonado and Nacional varieties contributed to the individual  
49 ancestry, and that the samples of the Amazonian Basin present a higher diversity in contrast  
50 to the lower diversity found in the samples from the Atlantic side.

51 The Bolivian native cacao varieties have not being considered within the ten groups described  
52 above (Motamayor et al., 2008; Cornejo et al., 2018) and they might be a different population.  
53 These studies have been done with Population Genetics but could also be resolved with

54 Barcoding. Nowadays, as the DNA sequencing prices have gone notably down is more  
55 accessible to sequence entire genomes or whole plastid genome rather than genes or DNA  
56 fragments. Whole plastid genomes as the chloroplast genome are more conserved than the  
57 nuclear genome. Plastid genomes enclose essential information markers for phylogenetic  
58 relationships among closely related taxa due to the low rate of polymorphism in the  
59 chloroplast. Therefore whole chloroplast genome sequencing has become an interesting  
60 barcoding tool for plants called ultra-barcoding (Kane and Cronk, 2008). Ultra-barcoding is  
61 based on data derived from high-throughput sequencing also called Next-Generation  
62 Sequencing (NGS). NGS data is more sensitive than traditional molecular markers (e.g.  
63 microsatellite) because genome target regions are significantly larger.

64 In this work, we focus on *T. cacao* ultra-barcoding which was used to identify genetic variation  
65 below species level. The results obtained in this work will provide valuable DNA sequence  
66 information for taxonomic studies and the development of molecular markers for below-  
67 species level identification of *T. cacao* coming from Bolivia. Ultimately, providing tools for *T.*  
68 *cacao* germplasm conservation.

## 69 MATERIALS AND METHODS

70 **Plant material:** *Theobroma cacao* fully developed leaves were collected, from two different  
71 varieties in two different regions from Bolivia: “mir20” from an indigenous Community  
72 (Miraflores, Pando) and “naz7” from a peasant community (Nazareth, Beni). Each tree received  
73 standard agronomic practices, and the cacao samples gained prizes for the best chocolate  
74 product in 2013, 2015, 2017 and 2019 at the international contest “Salon du Chocolat” in Paris,  
75 France. The leaves were collected and dry-stored at -20°C.

76 **DNA isolation and sequencing:** 5 mg of frozen leaves were crushed in a mortar to obtain fine  
77 powder using liquid nitrogen and the powder was transferred to microcentrifuge tubes. DNA  
78 isolation was performed using the Purelink Genomic DNA Kit (Thermo, CA, USA) according  
79 to the procedure described in the manufacture’s instructions. Chloroplast genome sequencing  
80 was outsourced to Omega Bioservices, USA and sequenced on a Miseq (Illumina, CA, USA).

81 Using 2 x 150 bp paired-end reads generated with the Nextera Truseq libraries (Illumina, CA,  
82 USA).

83 **Bioinformatic analysis:** Reads were trimmed and cleaned Assembly, mapping and short read  
84 post-processing were performed using Velvet (1.2.10), Bowtie2 (bowtie-  
85 bio.sourceforge.net/bowtie2, V. 2.3.5.) and Samtools utilities (htslib.org/). The annotation of  
86 the chloroplast genes was made by GeSeq (Tillich et al., 2017). A chlorogenome map was  
87 generated using OGDraw (Greiner, Lehwark, and Bock, 2019).

## 88 RESULTS

89 Trimmed and cleaned reads were further filtered out by mapping them against a *Theobroma*  
90 *cacao* reference chloroplast genome (RefSeq assembly: GCF\_000208745.1, National Center for  
91 Biotechnology Information). The chloroplast reference genome corresponds to the sample  
92 Scavina-6 from Perú (Kane et al., 2012; Argout et al., 2017). A total of approximately 47  
93 million reads (mir20 sample) and 25 million reads (naz7 sample) were used in the assembly of  
94 the two plastid genomes (Table 1). The chloroplast reference genome (HQ244500.2) has  
95 160,619 base pairs and our plastid coverage was more than a 100X for both samples. The  
96 coverage was enough to assemble the plastid genomes of both samples. The two plastids from  
97 the *T. cacao* Bolivian varieties differ in only 19 nucleotides for mir20 and 22 nucleotides for  
98 naz7 compared to the Belizean plastid genome. The two plastid genomes were deposited in  
99 GenBank under accessions: MW243993 for mir20 and MW243994 for naz7) The GC content  
100 of both samples was 36,9 % (Figure 1).

101

102 *Table 1. Illumina sequence summary statistics and observed coverage of the nuclear and chloroplast*  
103 *genome for Bolivian native cacao varieties.*

<i>Cacao</i> <i>variety</i>	<i>Read</i> <i>length</i>	<i>Total number</i> <i>of reads (PE)</i>	<i>Number of</i> <i>Reads after</i> <i>trimming</i>	<i>Total</i> <i>sequenced</i>	<i>Chloroplast</i> <i>coverage</i>	<i>Nuclear</i> <i>coverage</i>
Mir20	150 bp	48,384.722	47.505.946	7.12 Gbp	100 X	16.6 X
Naz7	150 bp	29.735.479	25.195.138	3.78 Gbp	100 X	8.8 X

104 PE: Paired-end

105 The comparison between the Belizean reference plastid genome and the Mir20 sample  
106 identified 19 different nucleotides. Comparing the reference plastid genome with the Naz7  
107 sample 22 variant nucleotide positions were observed. The Bolivian cacao plastid annotation  
108 contained 132 genes including 8 ribosomal RNA, 37 tRNA genes and 85 protein-coding genes.  
109 To explore the phylogenetic relationships we included 13 plastid genomes in our analysis  
110 (Table 2). The phylogenetic analyses revealed significant divergence between clades of *T. cacao*  
111 from the diverse varieties. The Maximum Likelihood tree showed two strongly supported *T.*  
112 *cacao* clades. The clades on both ends correspond to the Forastero and Criollo Groups. The  
113 clade in the middle belongs to the hybrid group between them, the Trinitario group (Fig. 2).  
114 The Bolivian Native Cacao forms a different group inside the hybrid clade. The tree also  
115 verifies that the Trinitario varieties (e.g. ICS01, ICS06) are hybrids between Forastero and  
116 Criollo. The *T. cacao* plastid genome accession HQ336404.2 (Jansen et al., 2011), which has no  
117 publicly available information about its geographical origin groups strongly with the  
118 Forastero plastid variety (e.g. Amelonado) (Fig. 2).

## 119 **DISCUSSION**

120 DNA barcoding is useful for several applications including identification below species level,  
121 estimating phylogenetic diversity and to identify species that are new to science. DNA  
122 barcoding allows to identify taxon through DNA sequencing using specific nuclear locus (e.g.

123 ITS region) (Bellemain et al., 2010), mitochondrial genes (e.g. *COI* and *CytB*) (Degli Esposti et  
124 al., 1993; Hebert et al., 2003) or complete genomes. The results reported in this work shows  
125 the benefits of whole chloroplast genome barcoding for organism identification below the  
126 species level as the Bolivian Native Cacao is identified as a different group from other *T. cacao*  
127 varieties (Fig. 2).

128 The *T. cacao* phylogenetic tree constructed with whole chloroplast genomes revealed three  
129 different clades below the species level (Fig. 2). The clades formed for the Criollo and  
130 Forastero group have been reported before with whole chloroplast genome sequencing as a  
131 barcoding tool (Kane et al., 2012). The Forastero group represents the cluster of cacaos from  
132 South America and the Criollo group represents the Central America group. A third clade was  
133 observed in our study and was formed mostly by the Trinitario hybrid variety (Fig. 2). This  
134 clade includes the Bolivian samples in a separate branch. The Bolivian Native Cacao forming a  
135 different cacao subpopulation has been already reported with microsatellites (Zhang et al.,  
136 2012). Bolivian Native Cacao is very likely that will form a new subgroup among the ten  
137 subpopulations described by Cornejo et al. (2018) and Motamayor et al. (2008).

138 The fact that Bolivian Native Cacao samples are inside the Trinitario group might be  
139 explained because most of the *T. cacao* varieties in this group live in the Western Amazon and  
140 the pacific coast of South America, just in the middle of the Criollo distribution zone (Central  
141 America) and the Forastero distribution zone (Atlantic coast of South America). The Western  
142 Amazon has been proposed as the Center of Origin for Cacao species through population  
143 genomics (Cornejo et al. 2018) and archeological evidence (Zarrillo et al., 2018). Thus, we  
144 suggest that the Trinitario group should be renamed as the Center of Cacao Origin group.  
145 Thus, avoiding the misconception that many of these varieties are hybrids between Forastero  
146 and Criollo.

## 147 **CONCLUSIONS**

148 Whole chloroplast genome sequencing of *T. cacao* is a useful approach to identify cacao  
149 varieties below the species level. The sequencing of more cacao varieties will deepen the

150 results obtained in this work, showing the Center of Cacao Origin in the Western Amazon  
151 also through Ultrabarcoding.

## 152 **ACKNOWLEDGEMENTS**

153 The authors would like to thankfully acknowledge the support of Centro de Investigación y  
154 Promoción del Campesinado (CIPCA), the Bolivian IDH funding at Universidad Mayor de San  
155 Andrés and the Swedish International Development Cooperation Agency (SIDA). The authors  
156 would like to express their gratitude to cacao-producing communities in Beni and Pando,  
157 Bolivia for his enormous help and support in this study.

## 158 **AUTHOR CONTRIBUTIONS**

159 PC, JMP, PM designed the study. MG, CP, OMRP, conducted the analysis, data  
160 interpretation and drafted the manuscript. LT and VP conducted part of the analysis and  
161 experiments. PC, JMP, PM, and CP supervised the work. All the authors contributed to and  
162 approved the final manuscript.

## 163 **LITERATURE CITED**

164 Argout, X., G. Martin, G. Droc, O. Fouet, K. Labadie, E. Rivals, J. M. Aury, AND C. Lanaud.  
165 2017. The cacao Criollo genome v2.0: an improved version of the genome for genetic and  
166 functional genomic studies. *BMC Genomics* 18: 730.

167 Bazoberry Chali, O., AND C. Salazar Carrasco. 2008. El cacao en Bolivia una alternativa  
168 económica de base campesina indígena. Centro de Investigación y Promoción del  
169 Campesinado.

170 Bellemain, E., T. Carlsen, C. Brochmann, E. Coissac, P. Taberlet, AND H. a. a. Kausrud.  
171 2010. ITS as an environmental DNA barcode for fungi: an in silico approach reveals potential  
172 PCR biases. *BMC microbiology* 10: 1.

173 Cheesman, E. 1944. Notes on the nomenclature, classification and possible relationships of  
174 cocoa populations. *Tropical Agriculture* 21: 144–159.

- 175 Cornejo, O. E., M.-C. Yee, V. Dominguez, M. Andrews, A. Sockell, E. Strandberg, D.  
176 Livingstone, et al. 2018. Population genomic analyses of the chocolate tree, *Theobroma cacao*  
177 *L.*, provide insights into its domestication process. *Communications biology* 1: 1-12.
- 178 de la Cruz, M., R. Whitkus, A. Gómez-Pompa, AND L. Mota-Bravo. 1995. Origins of cacao  
179 cultivation. *Nature* 375: 542-543.
- 180 Degli Esposti, M., S. De Vries, M. Crimi, A. Ghelli, T. Patarnello, AND A. Meyer. 1993.  
181 Mitochondrial cytochrome b: evolution and structure of the protein. *Biochimica et Biophysica*  
182 *Acta (BBA)-Bioenergetics* 1143: 243-271.
- 183 Greiner, S., P. Lehwark, AND R. Bock. 2019. OrganellarGenomeDRAW (OGDRAW)  
184 version 1.3.1: expanded toolkit for the graphical visualization of organellar genomes. *Nucleic*  
185 *Acids Research* 47: W59-W64.
- 186 Hebert, P. D., A. Cywinska, S. L. Ball, AND J. R. Dewaard. 2003. Biological identifications  
187 through DNA barcodes. *Proceedings of the Royal Society of London. Series B: Biological Sciences*  
188 270: 313-321.
- 189 Jansen, R. K., C. Saski, S.-B. Lee, A. K. Hansen, AND H. Daniell. 2011. Complete plastid  
190 genome sequences of three rosids (*Castanea*, *Prunus*, *Theobroma*): evidence for at least two  
191 independent transfers of *rpl22* to the nucleus. *Molecular biology and evolution* 28: 835-847.
- 192 Kane, N., S. Sveinsson, H. Dempewolf, J. Y. Yang, D. Zhang, J. M. M. Engels, AND Q. Cronk.  
193 2012. Ultra-barcoding in cacao (*Theobroma* spp.; Malvaceae) using whole chloroplast genomes  
194 and nuclear ribosomal DNA. *American Journal of Botany* 99: 320-329.
- 195 Kane, N. C., AND Q. Cronk. 2008. Botany without borders: barcoding in focus. *Molecular*  
196 *Ecology* 17: 5175-5176.
- 197 Motamayor, J. C., P. Lachenaud, J. W. d. S. e Mota, R. Loo, D. N. Kuhn, J. S. Brown, AND R.  
198 J. Schnell. 2008. Geographic and genetic population differentiation of the Amazonian chocolate  
199 tree (*Theobroma cacao* L). *PloS one* 3.



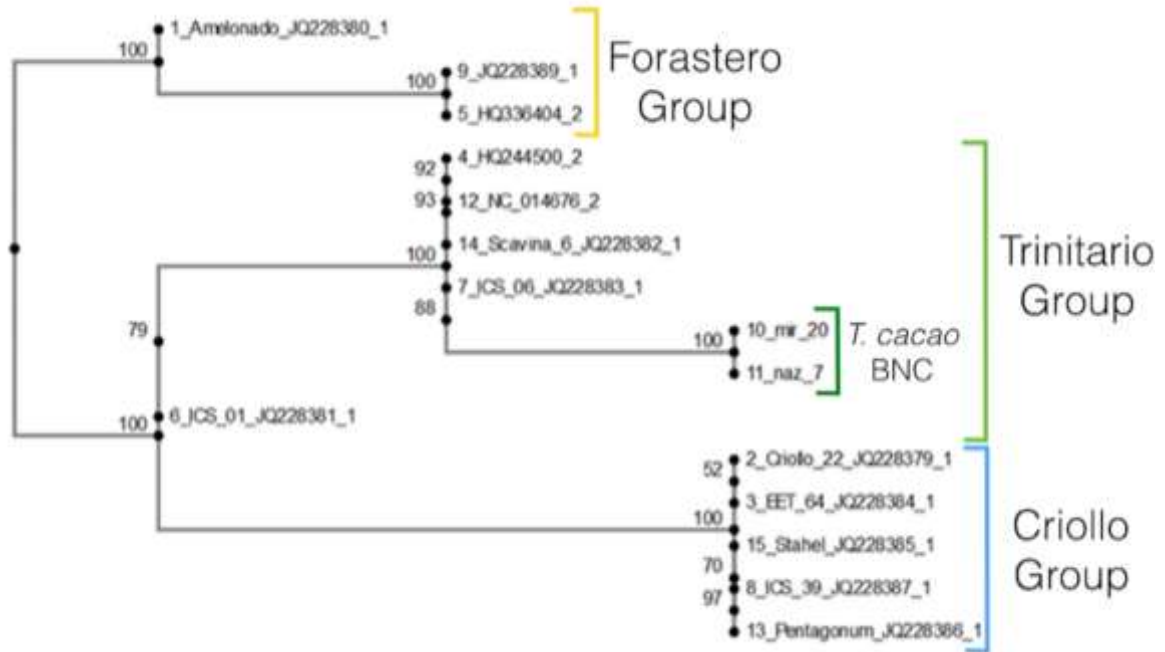
200 Tillich, M., P. Lehwark, T. Pellizzer, E. S. Ulbricht-Jones, A. Fischer, R. Bock, AND S.  
201 Greiner. 2017. GeSeq - versatile and accurate annotation of organelle genomes. *Nucleic Acids*  
202 *Res* 45: W6-w11.

203 Zarrillo, S., N. Gaikwad, C. Lanaud, T. Powis, C. Viot, I. Lesur, O. Fouet, et al. 2018. The use  
204 and domestication of *Theobroma cacao* during the mid-Holocene in the upper Amazon. *Nature*  
205 *Ecology & Evolution* 2: 1879-1888.

206 Zhang, D., W. J. Martínez, E. S. Johnson, E. Somarriba, W. Phillips-Mora, C. Astorga, S.  
207 Mischke, AND L. W. Meinhardt. 2012. Genetic diversity and spatial structure in a new  
208 distinct *Theobroma cacao* L. population in Bolivia. *Genetic Resources and Crop Evolution* 59: 239-  
209 252.

210





223

224 *Fig 2. The cacao phylogenetic tree constructed with chloroplast genomes reveals a unique group formed*  
225 *by Bolivian native cacao (BNC) varieties. Bootstrap values are given in percentage (%). Details of the*  
226 *chloroplast genome accessions are in Table 2.*

227 **Table 2.** List of the 13 chloroplast genomes used in this study. The accessions from the National Center for Biotechnology Information (NCBI) are listed  
 228 for easy reference..

Variety	Genbank Accession	Origin	Groups according to our barcoding analysis	Traditional variety classification
Amelonado	JQ228380.1	USDA, TRAS, Puerto Rico .	Forastero	Forastero <sup>1</sup>
-	HQ336404.2	-	Forastero	Unknown <sup>2</sup>
-	JQ228389.1	GI328924764	Forastero	Unknown <sup>1</sup>
Scavina-6	HQ244500.2*	Peru	Unknown	Forastero <sup>1</sup>
ICS-01	JQ228381.1	USDA, TARS, Puerto Rico	Trinitario?	Trinitario <sup>1</sup>
ICS-06	JQ228383.1	USDA, TARS, Puerto Rico	Trinitario?	Trinitario <sup>1</sup>
Scavina 6	NC_014676.2	Peru	Unknown	Forastero <sup>1</sup>
Scavina-6	JQ228382.1	Peru	Unknown	Forastero <sup>1</sup>

Stahel	JQ228385.1	Suriname	Criollo	Trinitario <sup>1</sup>
Criollo 22	JQ228379.1	USDA, SPCL, Beltsville, MD	Criollo	Criollo <sup>1</sup>
EET-64	JQ228384.1	USDA, TARS, Puerto Rico	Criollo	Forastero – Trinitario Hybrid <sup>1</sup>
ICS-39	JQ228387.1	USDA, TARS, Puerto Rico	Criollo	Trinitario <sup>1</sup>
Pentagonum	JQ228386.1	USDA, TARS, Puerto Rico	Criollo	Criollo <sup>1</sup>
Mir20	MW243993	Bolivia	Bolivian Native Cacao	Unknown <sup>3</sup>
Naz7	MW243994	Bolivia	Bolivian Native Cacao	Unknown <sup>3</sup>

229 \*Reference plastid Genome

230 1. Kane et al., 2012

231 2. Jansen et al., 2010

232 3. This paper

233