

1 **Phylogenetics analysis of TP53 gene in humans and its use in biosensors for breast**
2 **cancer diagnosis**

3 Sara da Silva Nascimento; Pierre Teodósio Félix*

4
5 Laboratory of Population Genetics and Computational Evolutionary Biology - [LaBECOM](#), UNIVISA,
6 Vitória de Santo Antão, Pernambuco, Brazil.

7 *Corresponding author/ **Contact:** pierrefelix@univisa.edu.br

8
9 **Keywords:** TP53; Biosensors; Diagnosis of Breast Cancer; Phylogeny; AMOVA

10 **Abstract**

11 Biosensors are small devices that use biological reactions to detect target analytes. Such
12 devices combine a biological component with a physical transducer, which converts bio-
13 recognition processes into measurable signals. Its use brings a number of advantages, as
14 they are highly sensitive and selective, relatively easy in terms of development, as well
15 as accessible and ready to use. Biosensors can be of direct detection, using a non-catalytic
16 ligand, such as cell receptors and antibodies, or indirect detection, in which there is the
17 use of fluorescently marked antibodies or catalytic elements, such as enzymes. They also
18 appear as bio-affinity devices, depending only on the selective binding of the target
19 analyte to the ligative attached to the surface (e.g., oligonucleotide probe). The objectives
20 were to evaluate the levels of genetic diversity existing in fragments of the TP53 gene
21 deposited in molecular databases and to study its viability as a biosensor in the detection
22 of breast cancer. The methodology used was to recover and analyze 301 sequences of a
23 fragment of the TP53 gene of humans from GENBANK, which, after being aligned with
24 the MEGA software version 6.06, were tested for the phylogenetic signal using TREE-
25 PUZZLE 5.2. Trees of maximum likelihood were generated through PAUP version
26 4.0b10 and the consistency of the branches was verified with the bootstrap test with 1000
27 pseudo-replications. After aligning, 783 of the 791 sites remained conserved. The
28 maximum likelihood had a slight manifestation since the gamma distribution used 05
29 categories + G for the evolutionary rates between sites with (0.90 0.96, 1.00, 1.04 and
30 1.10 substitutions per site). To estimate ML values, a tree topology was automatically
31 computed with a maximum Log of -1058,195 for this calculation. All positions containing
32 missing gaps or data were deleted, leaving a total of 755 sites in the final dataset. The
33 evolutionary history was represented by consensus trees generated by 500 replications,

34 which according to neighbor-join and BioNJ algorithms set up a matrix with minimal
35 distances between haplotypes, corroborating the high degree of conservation for the TP53
36 gene. GENE TP53 seems to be a strong candidate in the construction of Biosensors for
37 breast cancer diagnosis in human populations.

38 **1. Introduction**

39 Biosensors are small devices that use biological reactions to detect target analytes
40 (WANG, 2000). Such devices combine a biological component, which interacts with a
41 target substrate, to a physical transducer, which converts bio-recognition processes into
42 measurable signals (WANG, 2000; PATHAK *et al*, 2007). Its use brings a number of
43 advantages, as they are highly sensitive and selective, relatively easy in terms of
44 development, as well as accessible and ready to use. However, there are certain
45 limitations, such as electrochemically active interferences in the sample, little long-term
46 stability, and electron transfer problems (MEHRVAR; ABDI, 2004; SONG *et al*, 2006).

47 Biosensors can be direct detection (direct detection sensor or non-reticulated
48 system), in which biological interaction is measured directly, using a non-catalytic ligand,
49 such as cell receptors and antibodies, or indirect detection (marked sensor or reticulated
50 system), in which there is the use of fluorescently marked antibodies or catalytic elements,
51 such as enzymes. The crosslinked system has greater stability and is simpler to use, but
52 the non-reticulated system has better sensitivity, shorter operating time and lower costs
53 (MEHRVAR; ABDI, 2004; PATHAK *et al*, 2007; LIU *et al*, 2009). There are two types
54 of biosensors, depending on the nature of the recognition event. Bio affinity devices,
55 which depend on the selective binding of the target analyte to the ligand attached to the
56 surface (e.g., antibody or oligonucleotide probe) and bioanalytical devices, in which an
57 immobilized enzyme is used for target substrate recognition (WANG, 2000). Based on
58 this information, the objective of this work was to present a review of bibliography
59 describing the structure, functioning and applicability of biosensors in various
60 technological areas.

61 **3. Objective**

62 **3.1 General**

63 To evaluate the levels of genetic diversity existing in fragments of the TP53 gene
64 deposited in molecular databases.

65 **3.2 Specifics**

66 Evaluate the levels of polymorphism in the gene encoding the TP53 protein and
67 develop methodologies that allow the investigation of patterns of genetic variability for
68 this gene.

69 **4. Methodology**

70 **4.1. Dataset**

71 Initially, 301 sequences of a fragment of the human TP53 gene recovered from
72 GENBANK (<https://www.ncbi.nlm.nih.gov/popset/430765060>) and participated in a
73 PopSet made available by Hao, X.D and collaborators in 2013 (PopSet: 430765060) were
74 recovered and analyzed.

75 **4.2. Analyses**

76 After alignment with the mega software version 4.0 (KUMAR *et al.*, 2007), the
77 phylogenetic signal will be tested using TREE-PUZZLE 5.2 (SHIMIDT, 2002). Trees of
78 maximum likelihood will be generated through PAUP version 4.0b10 (SWOFFORD,
79 2002) and to evaluate the consistency of the branches, the bootstrap test (FELSENSTEIN,
80 1985) with 1000 pseudo-replications will be used. For the visualization of variable sites,
81 logos will be generated through the Weblogo3 program (CROOKS, 2004). The analysis
82 of the number of populations will be performed with the Structure 2.3 program
83 (PRITCHARD, 2000) and two different methods are tested: a posteriori probability and
84 ad hoc (k). The “*a posteriori*” probability will be calculated using an ancestry model with
85 mixed alleles for 20,000 interactions in the burn-in period, followed by 200,000 Monte
86 Carlo interactions via Markov Chain, increasing only the K value (number of
87 populations), which will be from 1 to 10 according to Pritchard's methodology (2000).

88 The Evanno method (2005) will be used to determine the most appropriate number
89 of populations for the dataset, using an ad hoc amount based on the second-order rate of
90 the likelihood function between the successive values of K. Posteriori and k probability
91 tests will initially be applied to the dataset in isolation. For the analysis of genetic
92 variability, a project will be created with the Arlequin Software 3.1 (EXCOFFIER *et al.*,
93 2005). which aims to measure molecular diversity using standard estimators such as Theta
94 (Θ , S, k, Π), Tajima Neutrality test, paired and individual F_{ST} values, in addition to
95 temporal divergence and demographic expansion indices (mismatch and Tau values) by

96 molecular variance analysis (AMOVA) (EXCOFFIER, 1992). In this method, the
97 distance matrix between all haplotype pairs will be used in a hierarchical variance analysis
98 scheme producing estimates of variance components analogous to Wright's F statistics
99 involving nonlinear transformations of the original information in estimates of genetic
100 diversity. Mantel's Z statistic will be used to represent the divergence between possible
101 microhabitats using the MULTIVAR (Mantel for Windows) program (MANTEL, 1967).

102 5. Results

103 After being aligned, 783 of the 791 sites remained conserved. The maximum
104 likelihood had a discrete manifestation for the gamma distribution with 05 categories +
105 G for the evolutionary rates between sites with 0.90 0.96, 1.00, 1.04 and 1.10 substitutions
106 per site. Nucleotide frequencies were A = 24.37%, T/U = 22.12%, C = 23.58% and G =
107 29.93%. For ML values, a tree topology was automatically computed with a maximum
108 Log of -1058,195 for this calculation (Figures 1a and 1b). All positions containing
109 missing gaps or data were deleted, leaving a total of 755 sites in the final dataset.

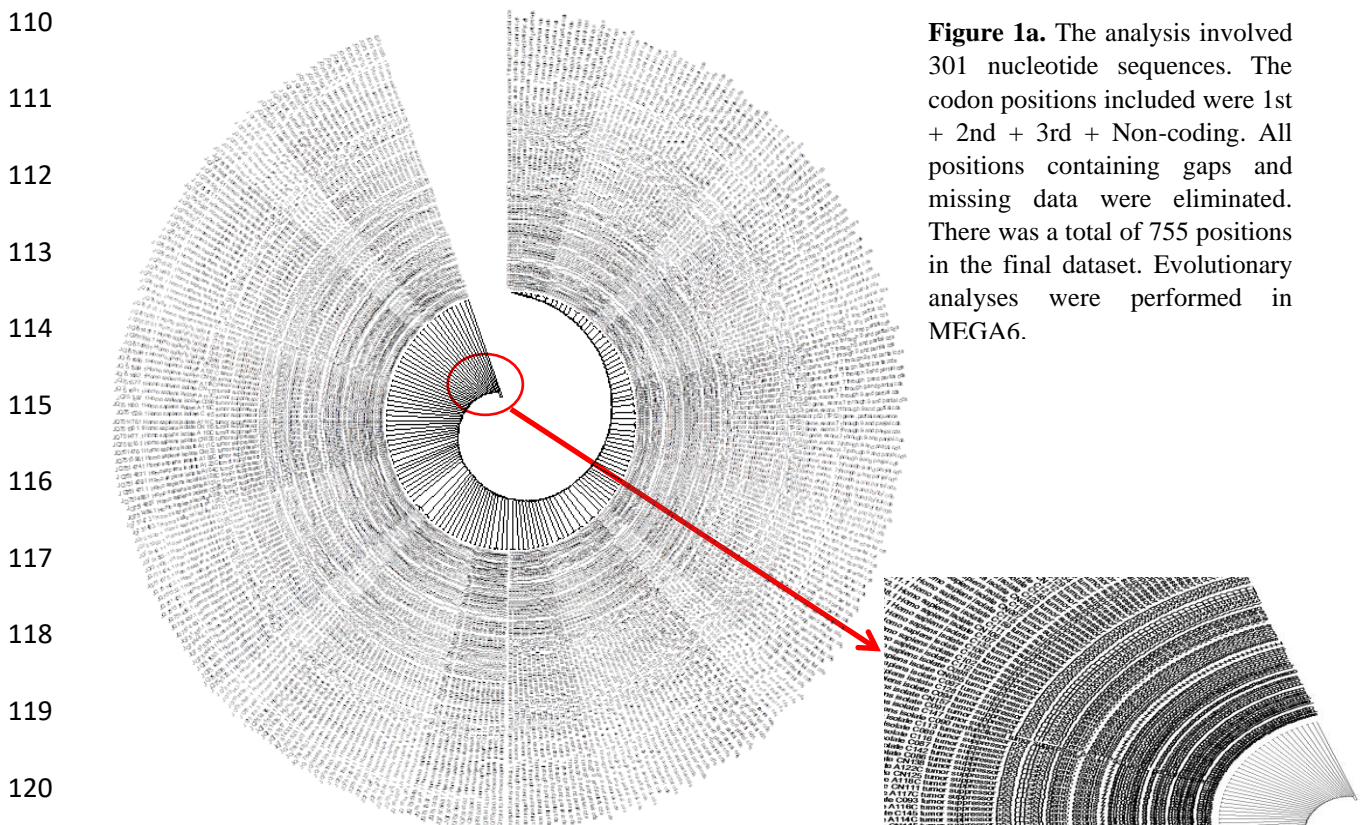
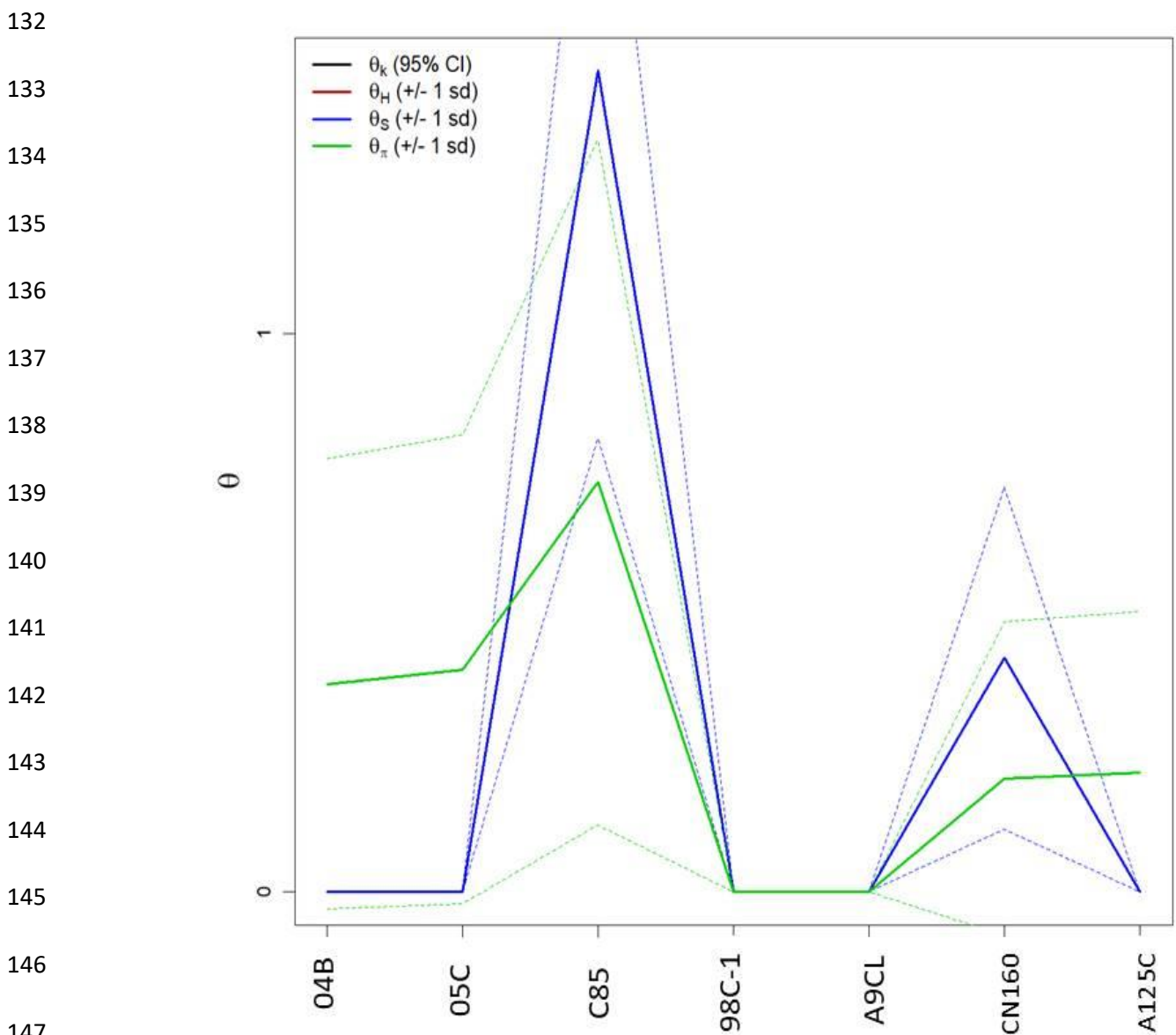


Figure 1a. The analysis involved 301 nucleotide sequences. The codon positions included were 1st + 2nd + 3rd + Non-coding. All positions containing gaps and missing data were eliminated. There was a total of 755 positions in the final dataset. Evolutionary analyses were performed in MEGA6.

Figure 1b. Cut showing the details of the haplotypes in the ML tree.

121
122

123 The evolutionary history was represented by consensus trees generated with 1,000
124 replications, which according to the algorithms of Neighbor-Join and BioNJ, set up a
125 matrix of distance between the haplotypes that corroborated the high degree of
126 conservation for the gene. For molecular variance tests, the 301 sequences were divided
127 into 07 groups (04b, 05c, c85, 98c-1, a9cl, cn160 and a125c) that did not present levels
128 of molecular diversity (0.05) (figure 2a, 2b), as well as in the Ewens-Watterson,
129 Chakraborty, Tajima D and Fu Fs tests (table 1). In the F_{ST} tests, the only important
130 variations were found within groups c85 and 04b with 0.73 and 0.39 respectively (figure
131 3, figure 4).



148 **Figure 2a:** Graphic representation of molecular diversity indices in groups 04 B, 05 C, C
149 85, 98C-1, A9CL, CN160, A125C. *Generated by the statistical package in R language
using the output data of the Arlequin software version 3.5.1.2.

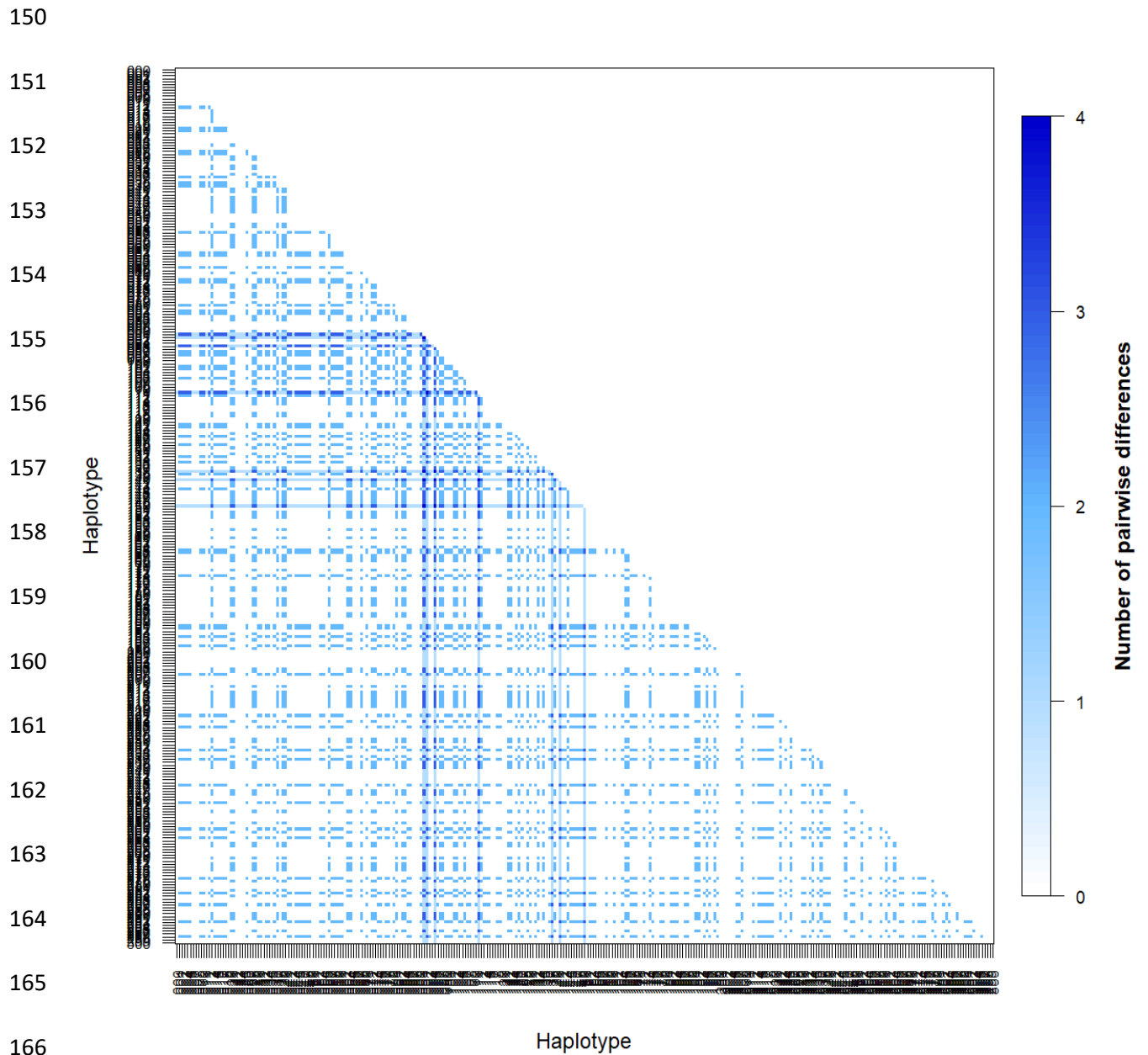


Figure 2b: Representation of the haplotypic distance matrix among the 301 sequences studied.
*Generated by the statistical package in R language using the output data of the Arlequin software version 3.5.1.2.

166
167
168
169
170
171
172
173

Table 1. Neutrality test for the seven groups studied

Estatísticas		04B	05C	C85	98C-1	A9CL	CN160	A125C	Mean	s.d.
Teste de Ewens-Watterson										
	Sample size	43	43	66	1	1	67	80	43.00	31.62
	No. of alleles(unchecked)	43	43	66	1	1	67	80	43.00	31.62
	Observed F value	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.
	Expected F value	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.
	Watterson test: Pr(rand F <= obs F)	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.
	Slatkin's exact test P-value	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.
Teste de Chakraborty										
	Sample size	43	43	66	1	1	67	80	43.00	31.62
	No. of alleles(unchecked)	43	43	66	1	1	67	80	43.00	31.62
	Obs. homozygosity	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Exp. no. of alleles	2.43	2.52	3.90	0.00	0.00	1.90	1.99	1.82	1.40
	P(k or more alleles)	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.	N.A.
Teste D de Tajima										
	Sample size	43	43	66	1	1	67	80	43.00	31.62
	S	0	0	7	0	0	2	0	1.28	2.62
	PI	0.37	0.39	0.73	0.00	0.00	0.20	0.21	0.27	0.25
	Tajima's D	0.00	0.00	-1.25	0.00	0.00	-0.86	0.00	-0.30	0.53
	Tajima's D-p-value	1.00	1.00	0.09	1.00	1.00	0.19	1.00	0.75	0.41
Teste FS de Fu										
	No. of alleles(unchecked)	43	43	66	1	1	67	80	43.00	31.62

174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197

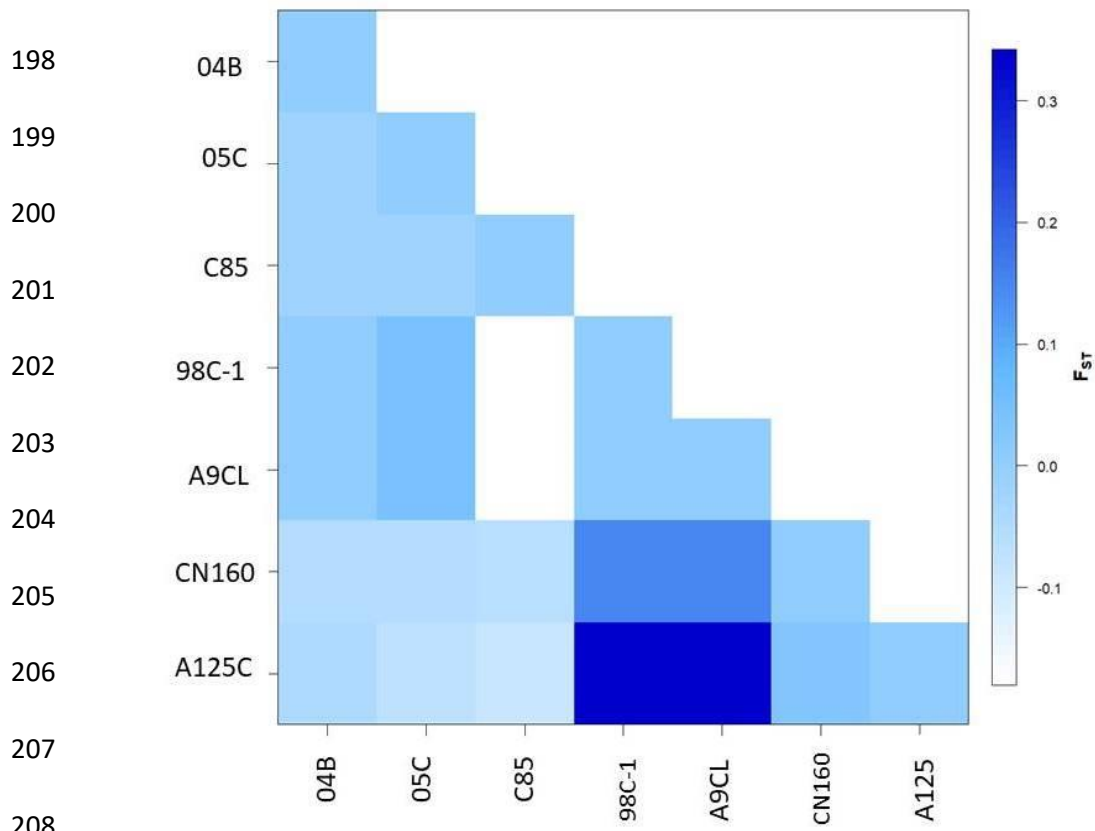


Figure 3. Matrix of genetic distance based on F_{ST} among the seven populations. * Generated by the statistical package in R language using the output data of the Arlequin software version 3.5.1.2.

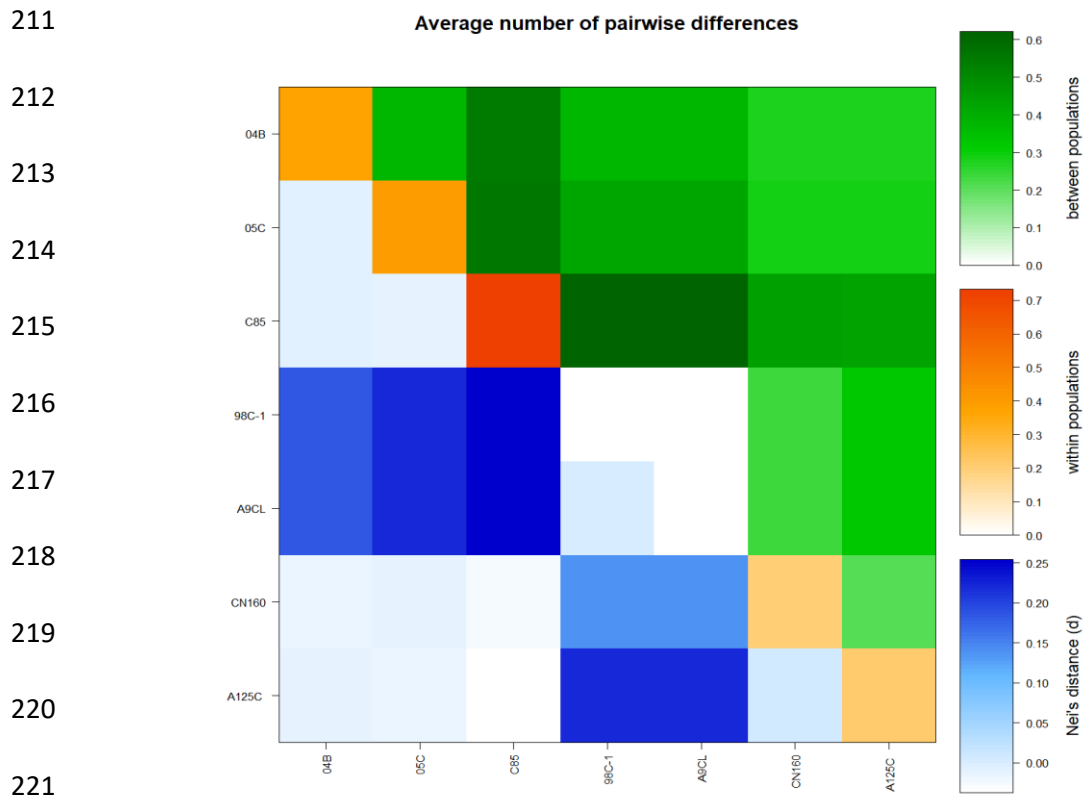


Figure 4. Matrix of paired differences between the populations studied: between the groups; within groups; and Nei distance for the seven groups. *Generated by the statistical package in R language using the output data of the Arlequin software version 3.5.1.2.

222 The results presented suggest that the TP53 gene is a strong candidate in the
223 construction of biosensors for the diagnosis of breast cancer in human populations, since
224 its polymorphism levels are not significant and its molecular diversity indexes are
225 unimpressive. Further analyses are still underway and we will soon have even more robust
226 results corroborating our hypothesis.

227 **7. Author's contributions**

228 To authors Xiao-Dan and collaborators by the availability of sequences in the
229 public databank.

230 **8. Acknowledgments**

231 I thank my advisor Prof. Dr. Pierre Teodósio Félix; to UNIVISA for the
232 availability and disposition of resources for the development of this work; to colleagues
233 in the Lab for patience and to Laboratory of Population Genetics and Computational
234 Evolutionary Biology - LaBECOM.

235 **9. References**

236 CROOKS G.E., HON G, CHANDONIA JM, BRENNER SE WEBLOGO: A
237 **sequence logo generator**, Genome Research, 14:1188-1190, (2004).

238 EVANNO, G; REGNAUT, S; GOUDET, J (2005). **Detecting the number of**
239 **clusters of individuals using the software structure: a simulation study**. Mol. Ecol.
240 14, 2611 – 2620.

241 EXCOFFIER, L. AND H.E. L. LISCHER (2010) **Arlequin suite ver 3.5: A new**
242 **series of programs to perform population genetics analyses under Linux and**
243 **Windows**. Molecular Ecology Resources. 10: 564-567.

244 EXCOFFIER, L; SMOUSE, P; QUATTRO, J (1992) **Analysis of molecular**
245 **variance inferred from metric distances among DNA haplotypes: Application to**
246 **human mitochondrial DNA restriction data**. Genetics 131, 479-491.

247 FELSENSTEIN J (1985) **Confidence limits on phylogenies - An approach**
248 **using the bootstrap**. Evolution 39:783-791.

249 HAO XD, YANG Y, SONG X, et al. **Correlation of telomere length shortening**
250 **with TP53 somatic mutations, polymorphisms and allelic loss in breast tumors and**
251 **esophageal cancer.** *Oncol Rep.* 2013;29(1):226-236. doi:10.3892/or.2012.2098.

252 KUMAR S, STECHER G, LI M, KNYAZ C, AND TAMURA K. **MEGA X:**
253 **Molecular Evolutionary Genetics Analysis across computing platforms. (2018).**
254 *Molecular Biology and Evolution* 35:1547-1549.

255 MANTEL, N. (1967). **The detection of disease clustering and a generalized**
256 **regression approach.** *Cancer Res.* 27, 209-220.

257 MEHRVAR, M.; ABDI, M. **Recent Developments, Characteristics, and**
258 **Potential Applications of Electrochemical Biosensors.** *Analytical Sciences*, v.20,
259 p.1113-1126, ago 2004.

260 PATHAK, P.; KATIYAR, V. K.; GIRI, S. **Cancer Research - Nanoparticles,**
261 **Nanobiosensors and Their Use in Cancer Research.** AZojono, 2004.

262 PRITCHARD, JK; STEPHENS, P; DONNELLY, P (2000) **Inference of**
263 **population structure using multilocus genotype data.** *Genetics* 155, 945–959.

264 SCHMIDT, J, K; **Business-unit-level relationship between employee**
265 **satisfaction, employee engagement, and business outcomes: A meta-analysis.** *Journal*
266 *of Applied Psychology*, Vol 87(2), Apr 2002, 268-279 (2002).

267 SWOFFORD DL (2002) **PAUP* phylogenetic analysis using parsimony (*and**
268 **other methods) Version 4.** Sinauer Associates, Sunderland, Massachusetts.

269 WANG, G. et al. **A Living Cell Quartz Crystal Microbalance Biosensor for**
270 **Continuous Monitoring of Cytotoxic Responses of Macrophages to Single-Walled**
271 **Carbon Nanotubes.** *Particle and Fibre Toxicology*, v.8, n.4, 2011.

272 WANG, J. **From DNA Biosensors to Gene Chips.** *Nucleic Acids Research*,
273 v.28, n.16, p.3011-3016, 2000.

274 WANG, Y. et al. **Electrochemical Sensors for Clinic Analysis.** *Sensors*, v.8,
275 p.2043-2081, mar 2007.