



# Development of novel simple sequence repeat markers from a genomic sequence survey database and their application for diversity assessment in *Jatropha curcas* germplasm from Guatemala

R.S. Raposo<sup>1\*</sup>, I.G.B. Souza<sup>2\*</sup>, M.E.C. Veloso<sup>2</sup>, A.K. Kobayashi<sup>3</sup>,  
B.G. Laviola<sup>3</sup> and F.M. Diniz<sup>2</sup>

<sup>1</sup>Núcleo de Biologia Experimental, Universidade de Fortaleza,  
Fortaleza, CE, Brasil

<sup>2</sup>Embrapa Meio-Norte, Teresina, PI, Brasil

<sup>3</sup>Embrapa Agroenergia, PqEB, Brasília, DF, Brasil

\*These authors contributed equally to this study.

Corresponding author: F.M. Diniz

E-mail: [fabio.diniz@embrapa.br](mailto:fabio.diniz@embrapa.br)

Genet. Mol. Res. 13 (3): 6099-6106 (2014)

Received July 29, 2013

Accepted December 9, 2013

Published August 7, 2014

DOI <http://dx.doi.org/10.4238/2014.August.7.25>

**ABSTRACT.** The last few years have seen a significant increase in the number of large-scale sequencing projects generating whole genome databases. These sequence databases can be surveyed (genome sequence survey) for tandem repeats as an alternative means to develop microsatellites for monitoring and selecting natural populations and cultivars of *Jatropha curcas*. A total of 100 tandem repeats were revealed from mining 368 genomic surveyed sequences available in the Kazusa DNA Research Institute database. Twenty microsatellite sequences were successfully amplified, resulting in repeatable and scorable polymerase chain reaction products. Genotyping of *J. curcas*

accessions from the Guatemalan population revealed 18 polymorphic loci. The average number of alleles per locus was 6.9, and allelic sizes ranged from 94 to 299 bp. Expected and observed heterozygosities ranged from 0.118 to 0.906 and from 0.082 to 0.794, respectively. Polymorphic information content values ranged from 0.114 (JcSSR-34) to 0.886 (JcSSR-33) with an average of 0.627. Analysis with Micro-Checker indicated few null alleles for locus JcSSR-37 in Guatemalan populations, which may be a possible cause of its deviation from Hardy-Weinberg equilibrium, even after Bonferroni's correction. No loci showed significant linkage disequilibrium. These microsatellite loci are expected to be valuable molecular markers in *J. curcas* because they show high levels of polymorphism and heterozygosity.

**Key words:** Microsatellites; Energy crop; Physic nut; Biofuel; Euphorbiaceae; Genetic diversity

## INTRODUCTION

*Jatropha curcas* is a plant with possible origin in Central America and is currently distributed in almost all regions of the world (Fairless, 2007; Dias et al., 2012). Its commercial interest has been sparked by its strong potential in the production of alternative fuels, which is mainly due to its high oil yield (Berchmans and Hirata, 2008). The species exists in climates that are unfavorable for most food crops, and it occurs spontaneously in very poor fertility soils (Arruda et al., 2004). Because of this, *J. curcas* can be considered as one of the most promising alternative oil sources to replace petroleum diesel. The plant is very resistant to disease, insect attack, mollusks, and fungi because it segregates caustic latex that drips from damaged or torn leaves (Emeasor et al., 2005; Thomas et al., 2008).

Along with the growing interest in the exploitation and commercial cultivation of *J. curcas*, the development of molecular technologies for monitoring and selecting natural populations and cultivars is of paramount importance. The sustainability of the activity must go through the implementation of molecular techniques based on DNA analysis.

In this context, molecular markers such as microsatellites have been widely used (De-Woody et al., 1995; Ellegren, 2004; Guo et al., 2013; Qi et al., 2013). Microsatellites are repeated sequences of nucleotides that are distributed throughout the genome and occur mainly in non-coding regions (Ellegren, 2004). These sequences are used as markers because of their highly polymorphic pattern. Microsatellites can be used not only in the conservation of genetic resources of commercial interest but also in the selection of plants with desirable features for cultivation and identification of quantitative trait loci related to diseases or reproduction (Nambisan, 2007).

The last few years have seen a significant increase in the number of large-scale sequencing projects generating whole genome databases. These sequence databases can be surveyed (genome sequence survey) for tandem repeats as an alternative means to develop microsatellites, which may represent a useful and low-cost procedure for discovering novel microsatellite markers in a genome. Now that the whole genome of *J. curcas* has been sequenced (Sato et al., 2011), it is valuable to use the database for mining microsatellite loci in the species.

Thus, this study aimed at mining part of the *J. curcas* genome database to increase the availability of microsatellites, which are hereafter intended for the genetic analysis of *J. curcas* worldwide, and provide a genetic snapshot of the Guatemalan population diversity.

## MATERIAL AND METHODS

### Identification of repeated regions and primer design

A total of 368 contig consensus sequences were downloaded from the Kazusa DNA Research Institute database ([www.kazusa.or.jp/jatropha/](http://www.kazusa.or.jp/jatropha/)). These surveyed sequences were screened for tandem repeats using the Tandem Repeat Occurrence Locator software (Castelo et al., 2002). Primer pairs for microsatellite loci were designed based on the unique flanking regions of each repeat (dinucleotide, trinucleotide, tetranucleotide, and pentanucleotide) by using PRIMER 3 with the default settings (Rozen and Skaletsky, 2000) but with a fragment length limit of 300 bp (Table 1). The presence of multiple tandem repeats at distant sites in the same contig sequence allowed the design of different primer pairs targeted to each microsatellite. Only one pair of primers was designed for contiguous microsatellites, even if they were separated by a few bases.

### Plant material and DNA isolation

Fresh young plant leaves of *J. curcas* were collected from the Biocombustibles de Guatemala S.A. germplasm bank and stored in a plastic bag with silica gel to dehydrate. Approximately 20 mg plant tissue was ground in 2.0-mL microtubes containing ceramic grinding beads (CK28, BioAmerica, USA) with two 16-s pulses at 5200 rpm with 10-s intervals on a Precellys®24 Tissue Homogenizer (Bertin Technologies, France). These homogenized materials were used for DNA isolation with DNeasy Plant Mini Kit (Qiagen, Germany), according to manufacturer instructions, and checked for quality and quantity on a 1% agarose gel with ethidium bromide. The DNA concentration was measured using a NanoDrop 2000 Spectrophotometer (Thermo Scientific, USA). DNA was maintained at -20°C until further analyses.

### Polymerase chain reaction (PCR) amplification and microsatellite marker polymorphism

Forty-nine accessions of *J. curcas* from Guatemala were used for the characterization of microsatellite markers. PCR amplifications and optimization of primers were performed in a total reaction volume of 10 µL consisting of 1.0 µL template DNA (20-100 ng), 1.5-2.5 mM MgCl<sub>2</sub>, 50 µM of each dNTP, 0.5 U Taq DNA polymerase (RBC), 0.3-0.5 µM forward and reverse primers, and 1X PCR buffer (10 mM Tris-HCl, pH 8.3, 50 mM KCl) in a Veriti® 96-well thermalcycler (Applied Biosystems, USA). The cycling conditions included an initial denaturation cycle at 95°C for 2 min; 30 cycles of 95°C for 45 s, 30 s for annealing (temperature depended on each primer pair; Table 2), and 45 s for extension at 72°C; and a final extension at 72°C for 10 min.

Microsatellite markers were screened on denaturing 6% polyacrylamide gels that were stained with silver nitrate for genotype scoring. Allele sizes were initially determined against a 10-bp DNA ladder (Invitrogen, USA) and by comparison with the expected size of the cloned fragment, and then they were scored manually. If necessary, two or more runs were performed to verify the allele typing by re-ordering the samples.

## Statistical analysis

The numbers of alleles at each microsatellite locus, the proportion of individual samples that are heterozygous (observed heterozygosity), the estimate of heterozygosity (expected heterozygosity), and tests for departure from Hardy-Weinberg equilibrium (HWE) were performed using the probability test in GENEPOP version 3.3 (Raymond and Rousset, 1995). Allele frequencies obtained from the microsatellite genotypes were also used to calculate the polymorphic information content (PIC) of each locus using Cervus 3.0 (Kalinowski et al., 2007) in order to measure the degree of polymorphism obtained by a microsatellite. The occurrence of linkage disequilibrium between loci and the allelic richness for each locus were calculated using FSTAT 2.9.3 (Goudet et al., 1995). Significance levels were adjusted using sequential Bonferroni's corrections. Micro-Checker v2.2.3 was used to determine the most probable cause of any deviation from HWE (Van Oosterhout et al., 2004).

## RESULTS AND DISCUSSION

A total of 100 (27.17%) tandem repeats were revealed from mining 368 genomic surveyed sequences (contigs) available in the Kazusa DNA Research Institute database for *J. curcas*. From these microsatellite-containing sequences, 71% were dinucleotide repeats, 28% were trinucleotides, and 1% were tetranucleotides. Because primer designing was not considered for flanking regions of microsatellites with a small number of repeats, as low levels of variability or fixation is expected (Amos, 1999), our efforts focused on provisional markers with more than 8 tandem repeats. Additionally, other repeat regions were also found that were too short or had GC content that was too low to have primers designed on their flanks. Twenty microsatellite sequences were successfully amplified, resulting in repeatable and scorable PCR products.

Of 20 loci, 18 (90%) contained dinucleotide microsatellites, of which, 16 (88.9%) were pure, with only two other sequences that were characterized as compound dinucleotides. Two loci were trinucleotide repeats (Table 1).

Genotyping of *J. curcas* accessions from the Guatemalan population revealed 18 polymorphic loci; however, primers for locus JcSSR-38b were merely designed as an amplification alternative. The other two loci were monomorphic (Table 2).

The average number of alleles per locus was 6.9, ranging from 3 (JcSSR-38a/b and JcSSR-42) to 16 (JcSSR-33). The allelic sizes ranged from 94 to 299 bp. Expected and observed heterozygosities ranged from 0.118 to 0.906 and from 0.082 to 0.794, respectively. Distributions of allele frequencies in the *J. curcas* population are shown in Figure 1. Mean expected and observed heterozygosities were 0.675 and 0.593. PIC values ranged from 0.114 (JcSSR-34) to 0.886 (JcSSR-33) with an average of 0.627. This genetic parameter is an important estimate of the extent of polymorphism of the marker (DeWoody et al., 1995). PIC values higher than 0.5 are considered to be highly informative. PIC values between 0.25 and 0.50 indicate a reasonably informative locus (Botstein et al., 1980). Therefore, all the analyzed microsatellites were informative in the Guatemalan population of *J. curcas* except marker JcSSR-34 (PIC: 0.114). Analysis with Micro-Checker indicated few null alleles for locus JcSSR-37 in Guatemalan populations, which may be a possible cause of its deviation from HWE, even after Bonferroni's correction for multiple comparisons at the 5% significance level (critical value for  $P > 0.0029$ ). No loci showed significant linkage disequilibrium after Bonferroni's correction.

These additional microsatellite loci described here are expected to be valuable molecular markers because they show high levels of polymorphism and heterozygosity.



**Table 1.** Continued.

Marker name	Contig code	Contig sequence	Motif	Fragment length (bp)
JcSSR-39	JCCA0001061	tttaacatattgggtgaattttactctatataatataatataaagttattact acagatgaaatgataaacctagtcctctatacctctgttattatgcaaaa agtgggttgctctagtgctattctcataatgttcctaagaatgttg ataatataatataatataaatttgcaagaaaataaatcaatttaatt aatttagcttfaataaattcaatgagaatattcaattgatatataaatgat cacacatagatatggatagcgtgacccac	(TA) <sub>11</sub>	115
JcSSR-40	JCCA0107961	gtacgagcgtggactaacctcactgatgtgtcgggtgctgctaaaatcat gcatgacgtgggcccattgatgctccaagccgtcacattaccgccatata ataatataatataatataatataatataatataatataatataatataat aatgcaatttttaaaaacacattaattgtattctcataatataaaaatata tattttctaaaattgctctctttaa	(AT) <sub>12</sub>	196
JcSSR-41	JCCA0010831	tgcatcaacacatccattcttccaccaataatataatataatataatata aacctgagataaattctgcaataacttctcaattctgctgaattagaatca agttgaaattccaatgctacttcccactc	(TA) <sub>21</sub>	169
JcSSR-42	JCCA0108601	tgcatcaacacatccattcttccaccaataatataatataatataatata aacctgagataaattctgcaataacttctcaattctgctgaattagaatca agttgaaattccaatgctacttcccactc	(AT) <sub>10</sub>	138

**Table 2.** Characteristics of 20 genome sequence survey-derived microsatellites for *Jatropha curcas*.

Marker name		Ta (°C)	Size range (bp)	N <sub>A</sub>	H <sub>E</sub>	H <sub>O</sub>	HWE	F <sub>IS</sub>	PIC
JcSSR-24	F: CGCTGTAGGAATGGCATGTT R: CAAGCAATCAATCTCTCTTTCTTTT	60	100-109	4	0.691	0.714	0.0040	-0.033	0.625
JcSSR-25	F: CTGGTGTGGGCTTACATT R: GGTACACGCACTAGGAACCTGA	62	231-249	4	0.698	0.592	0.2563	0.152	0.630
JcSSR-26	F: AATCTACTATCCCGCATGAA R: TCCTACTCTTTCTCCATTCTGT	58	136-160	6	0.791	0.689	0.8493	0.129	0.751
JcSSR-27	F: CTGCATCTGGGAAACAAAT R: ATAACACCCATGCCTATCAA	57	106-128	4	0.430	0.364	0.0878	0.153	0.390
JcSSR-28	F: TATGCAAGTACGGTTGTGTA R: CAGCTAATCTTGGTTATGG	57	197-207	5	0.611	0.531	0.0053	0.132	0.556
JcSSR-29	F: AAAATAGAAAATAAACTCGCTCAA R: CAATGTATTGTATGGTTGCAATTA	56	94-128	13	0.876	0.794	0.0087	0.093	0.848
JcSSR-30	F: GTGCGCCCAAATGAAAA R: ATGGGGTTGCCAGGG	57	254-288	7	0.805	0.750	0.4681	0.069	0.768
JcSSR-31	F: ATTTGATTGATTGCCGAC R: GCAGCACCTCCCCTAAAA	59	114	1	-	-	-	-	-
JcSSR-32	F: TCTTTTGACTACCAATGACCA R: TCGTGCTAGAGTTGGTTGAAAT	60	186-192	4	0.644	0.563	0.5595	0.126	0.565
JcSSR-33	F: GAGAACGTGGAGAGTTAAAA R: TGGATGGATCATTGTCTG	55	237-297	16	0.906	0.778	0.9576	0.142	0.886
JcSSR-34	F: CTA AAAAGTTTGGACATTTTAC R: TTAATGGAAGCACTCCTGTAAT	55	174-186	4	0.118	0.082	0.1529	0.311	0.114
JcSSR-35	F: GAGTTTGGTTGGATGGATA R: TAGAACCGCAGCAGAAAA	54	172-208	10	0.889	0.769	0.9896	0.135	0.856
JcSSR-36	F: GCACTCAAAGTGTTCCCAA R: ATCACACCAAAGAAGAATGCAG	59	163-187	4	0.699	0.636	0.9487	0.090	0.620
JcSSR-37	F: TGCAGGCACACCAACAAA R: TTTCGATGGAATCTAGGTTAGTCA	58	269-299	11	0.873	0.761	0.0000*	0.128	0.847
JcSSR-38a	F: TGAAGCTGTCTTTGATTGAGGAA R: TCTTGGGATCGTCTACAAGGTCT	61	209-213	3	0.521	0.469	0.2779	0.098	0.431
JcSSR-38b	F: ATCAGCAACTCCTTAATGCCAA R: CAACTCAAAGTTTCCAGCCAAA	59	153-157	3	0.537	0.461	0.2991	0.081	0.441
JcSSR-39	F: TTTTAAACATAITGGGTTGAATTT R: TTTTGCATAAATAACAGAAGGT	58	115	1	-	-	-	-	-
JcSSR-40	F: AGTGGTGGTTGTGCTCATAGTG R: GTGGGTCACGCTATCCATATCT	61	194-218	11	0.877	0.723	0.1786	0.175	0.853
JcSSR-41	F: GTACGAGCGTGGACTAACT R: GACATTACAAATATTTGAAACG	52	157-179	9	0.800	0.700	0.9525	0.125	0.758
JcSSR-42	F: TGCATCAACACATCCATTCTTT R: GAGTGGGAAGTAGCATTGGAA	57	138-142	3	0.375	0.298	0.0379	0.206	0.338

\*Significant at the 5% significance level ( $P > 0.0029$ ). Forty-nine samples were genotyped; N<sub>A</sub> = number of alleles at each locus; H<sub>E</sub> = expected heterozygosity; H<sub>O</sub> = observed heterozygosity; HWE = deviation from Hardy-Weinberg equilibrium; F<sub>IS</sub> = inbreeding coefficient; PIC = polymorphic information content.

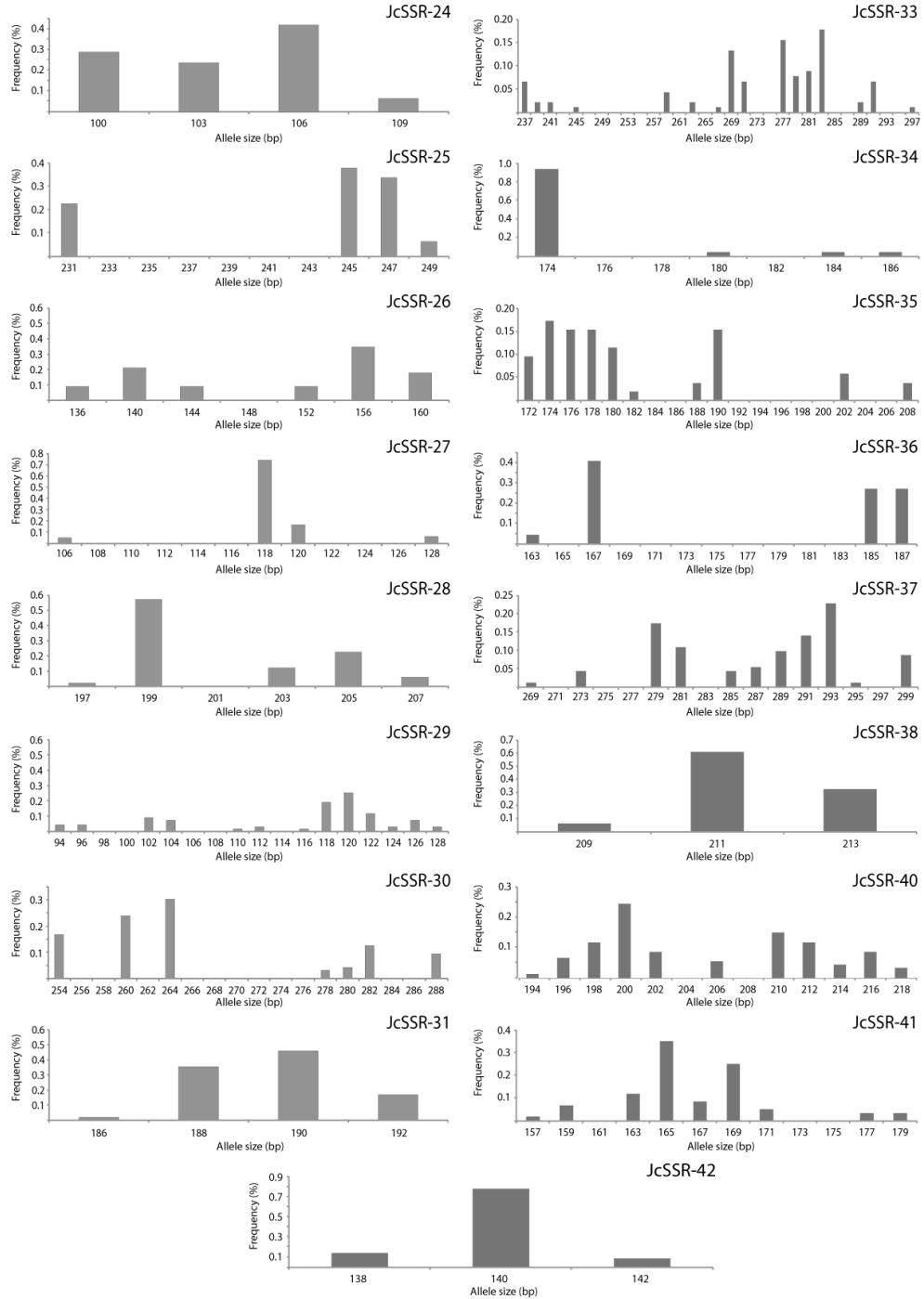


Figure 1. Distributions of allele frequencies at the 17 polymorphic microsatellite loci in *Jatropha curcas* populations.

## ACKNOWLEDGMENTS

Research supported by the Brazilian Government, the Brazilian Ministry of Science and Technology (MCT/FINEP, BRJatropha Project), the Bank of Northeastern Brazil (BNB/ETENE/FUNDECI), and the Brazilian Agricultural Research Corporation (EMBRAPA).

## REFERENCES

- Amos W (1999). A Comparative Approach to the Study of Microsatellite Evolution. In: Microsatellites: Evolution and Applications (Goldstein DB and Schlotterer C, eds.). Oxford University Press, Oxford, 67-79.
- Arruda FP, Beltrão NEM, Andrade AP, Pereira WE, et al. (2004). Cultivo de pinhão manso (*Jatropha curcas* L.) como alternativa para o semi-árido nordestino. *Rev. Bras. Oleag. Fibr.* 8: 789-799.
- Berchmans HJ and Hirata S (2008). Biodiesel production from crude *Jatropha curcas* L. seed oil with a high content of free fatty acids. *Bioresour. Technol.* 99: 1716-1721.
- Botstein D, White RL, Skolnick M and Davis RW (1980). Construction of a genetic linkage map in man using restriction fragment length polymorphisms. *Am. J. Hum. Genet.* 32: 314-331.
- Castelo AT, Martins W and Gao GR (2002). TROLL - tandem repeat occurrence locator. *Bioinformatics* 18: 634-636.
- DeWoody JA, Honeycutt RL and Skow LC (1995). Microsatellite markers in white-tailed deer. *J. Hered.* 86: 317-319.
- Dias LA, Missio RF and Dias DC (2012). Antiquity, botany, origin and domestication of *Jatropha curcas* (Euphorbiaceae), a plant species with potential for biodiesel production. *Genet. Mol. Res.* 11: 2719-2728.
- Ellegren H (2004). Microsatellites: simple sequences with complex evolution. *Nat. Rev. Genet.* 5: 435-445.
- Emeasor KC, Ogbuji RO and Emosairue SO (2005). Insecticidal activity of some seed powders against *Callosobruchus maculatus* (F.) (Coleoptera: Bruchidae) on stored cowpea. *J. Plant Dis. Protect.* 112: 80-87.
- Fairless D (2007). Biofuel: the little shrub that could - maybe. *Nature* 449: 652-655.
- Goudet J (1995). FSTAT (version 1.2): a computer program to calculate F-statistics. *J. Hered.* 86: 485-486.
- Guo BY, Qi PZ, Zhu AY, Lv ZM, et al. (2013). Isolation and characterization of new polymorphic microsatellite markers from the cuttlefish *Sepiella maindroni* (Cephalopoda; Sepiidae). *Genet. Mol. Res.* 12: 2376-2379.
- Kalinowski ST, Taper ML and Marshall TC (2007). Revising how the computer program CERVUS accommodates genotyping error increases success in paternity assignment. *Mol. Ecol.* 16: 1099-1106.
- Nambisan P (2007). Biotechnological intervention in *Jatropha* for biodiesel production. *Curr. Sci.* 93: 1347-1348.
- Qi PZ, Xie CX, Guo BY, Wu CW, et al. (2013). Development of new polymorphic microsatellite markers in topmouth culter (*Culter alburnus*) and determination of their applicability in *Culter mongolicus*. *Genet. Mol. Res.* 12: 1761-1765.
- Raymond M and Rousset F (1995). Genepop (version 1.2): population genetics software for exact tests and ecumenicism. *J. Hered.* 86: 248-249.
- Rozen S and Skaletsky HJ (2000). Primer3 on the WWW for General Users and for Biologist Programmers. In: Bioinformatics Methods and Protocols (Methods in Molecular Biology) (Krawetz S and Misener S, eds.). Humana Press, Totowa, 365-386.
- Sato S, Hirakawa H, Isobe S, Fukai E, et al. (2011). Sequence analysis of the genome of an oil-bearing tree, *Jatropha curcas* L. *DNA Res.* 18: 65-76.
- Thomas R, Sah NK and Sharma PB (2008). Therapeutic biology of *Jatropha curcas*: a mini review. *Curr. Pharm. Biotechnol.* 9: 315-324.
- Van Oosterhout C, Hutchinson WF, Shipley P and Wills DPM (2004). Micro-Checker: Software for identifying and correcting genotyping errors in microsatellite data. *Mol. Ecol. Notes* 4: 535-538.