

1 **Seqs-Extractor: Automated sequences extraction to**  
2 **reduce tedious manual corrections of large datasets.**

3 Patrick D. C. Pereira<sup>1</sup>, Cleysian Dias<sup>2</sup>, Mauro A. D. Melo<sup>1</sup>, Nara G. M. Magalhães<sup>3</sup>, Cristovam  
4 Guerreiro-Diniz<sup>1</sup>, Cristovam W. Picanço-Diniz<sup>3</sup>

5 <sup>1</sup> Laboratório de Biologia Molecular e Neuroecologia, Instituto Federal de Educação, Ciência e  
6 Tecnologia do Pará, Bragança, Pará, Brasil.

7 <sup>2</sup> Grupo de Genética e Conservação, Universidade Federal do Pará, Bragança, Pará, Brasil.

8 <sup>3</sup> Laboratório de Investigações em Neurodegeneração e Infecção no Hospital Universitário João  
9 de Barros Barreto, Instituto de Ciências Biológicas, Universidade Federal do Pará, Belém, Pará,  
10 Brasil.

11

12 Corresponding Author:

13 Patrick Pereira<sup>1</sup>

14 Rua da Escola Agrícola, S/N, Vila Sinhá, Bragança, Pará, Brasil, 68600-000.

15 Email address: patrick@ufpa.br

16

17

18 **Abstract**

19 The analysis of large numbers of sequences requires the reduction of ambiguities during  
20 the analytical work to ensure that the effort will focus only on confirmed sequences.

21 Performing this work automatically may help to minimize potential errors associated with  
22 tedious manual correction, allowing more effective results. Basic local alignment search  
23 tool (BLAST) seems to be the most widely used sequence analysis program. It is free, but  
24 commercial parties enhanced BLAST applications and charge a fee for their uses. There  
25 are some tools of public domain that can perform the search of microsatellites in the next  
26 generation sequencing (NGS) data, as the microsatellite identification tool (MISA), which  
27 has some features to discover microsatellites in large datasets. Here, we developed a basic  
28 shell script (BASH script) to be ran under Linux environment that can be used to extract  
29 from a sequence dataset only confirmed (BLASTed) sequences from both nucleotide

30 (BLASTN) and protein (BLASTX) databases and extract sequences that contains  
31 microsatellites using MISA tool, using a friendly interface and no fees charged. Seqs-  
32 Extractor is a helpful tool that may enhance the analysis of large datasets in BLAST+ and  
33 MISA by minimizing the time of management, reducing potential errors caused by  
34 manipulating data and no fees charged. Seqs-Extractor is available at  
35 <https://github.com/patrick-douglas/Seqs-Extractor/wiki>.

36 **Subjects:** Bioinformatics, Computational Biology, Genomics.

37 **Keywords:** Sequences analysis, Next-generation sequencing, Databases, Bash.

38

### 39 **Introduction**

40 Evaluation of sequences homology is the most common way to help understand the  
41 biology of specific organisms (Donkor et al. 2014). Indeed, from sequence homology  
42 studies we may recognize for example, the homology between protein or DNA  
43 sequences, that may shed light on shared ancestry in the evolutionary history of life. The  
44 main program used to achieve this goal is the popular Basic Local Alignment Search Tool  
45 (BLAST), that uses a heuristic algorithm which performs comparisons between pairs of  
46 sequences, searching for regions with some similarity (Altschul et al. 1990).

47 BLAST+ is the standalone BLAST suite used in an offline environment, when is  
48 not possible to work in online-BLAST, normally due to the large size of datasets to be  
49 processed (Camacho et al. 2009). To perform a BLAST search of large datasets using  
50 standalone tools, a minimum command line knowledge is required, because the  
51 standalone application relies on that. Indeed, the command line interface introduces

52 technical difficulties ranging from installation to generation of results (Camacho et al.  
53 2009). Free easier to use software, with graphic interface may help users to circumscribe  
54 operational difficulties of command line interface, but this free facility is limited to  
55 smaller datasets (Excoffier & Lischer 2010). Alternatively, the users may pay for  
56 software with friendly user interface, like Blast2GO (Conesa et al. 2005) that perform  
57 BLAST search and many others analyses of larger datasets, however the cost to use paid  
58 programs can be high.

59 In the next generation sequencing (NGS) the BLAST search is a good way to  
60 validate the assembled sequences by comparing it with a valid database. Moghadam et al.  
61 (2013) used this approach to validate their transcripts of *Charadrius vociferous* (obtained  
62 in *De novo* RNA-Seq assembly) BLASTing it against some datasets of genome sequences  
63 from Chicken, Turkey and Zebra finch, using only the top hit sequences (sequences that  
64 matched 80% to 100% with a valid database). This procedure may avoid considering  
65 contaminant derived sequences during analysis. However, get only the top hit sequences  
66 are challenger because a BLAST search generates output files containing only details  
67 about alignment scores of entered sequences and the subject database, but does not  
68 generated a file containing sequences that match with a specific percentage value.

69 Simple sequence repeats (SSR) or microsatellite markers are known as a good tool  
70 to identification of genetic distances among organisms, and is commonly used in the  
71 population genetic studies (Fernandez-Silva et al. 2013; Koohi-Dehkordi et al. 2006). In  
72 the NGS studies the use of SSR is very common the production of a large amount of  
73 sequences containing microsatellites (Castoe et al. 2012; Donkor et al. 2014; Fernandez-  
74 Silva et al. 2013; Guichoux et al. 2011).

75 There are some tools of public domain that can perform the search of microsatellites

76 in the NGS data, as the microsatellite identification tool (MISA), which has some features  
77 to discover microsatellites in large datasets (Thiel et al. 2003) and has already been used  
78 in several other studies (Khlestkina et al. 2004; Varshney et al. 2005; Varshney et al.  
79 2002; Yu et al. 2004).

80 Similar to BLAST+, the MISA tool provides the microsatellites results in a table  
81 containing the information about length, nucleotides repeats, motifs but does not allow  
82 the access of sequences that contains microsatellites. A common run of MISA in a NGS  
83 dataset can produce hundreds of microsatellites results, and get manually only the  
84 sequences that contains these results is usually unfeasible due to the tedious and the long  
85 time spent on it.

86 Here we provided a bash script, named as Seqs-Extractor that can run natively under  
87 Linux, which can extract in isolation from .FASTA dataset, the sequences that match in  
88 a BLASTX or a BLASTN search (with a match percentage defined by the user), as well  
89 as positive sequences of a MISA search.

90 Our script can also perform an independent BLAST+ search, using a friendly  
91 interface and no fees charged, considering cases when the user did not run BLAST+ yet.  
92 In addition, Seqs-Extractor can extract sequences from a .FASTA dataset using only a  
93 simple text file containing only the sequence IDs.

## 94 **Methods**

95 *Seqs-Extractor* is implemented in BASH command language and run natively in Linux  
96 systems, which are based in the Debian version. This script was tested under Linux Mint  
97 17.3 Xfce and Ubuntu 16.4LTS. At time *Seqs-Extractor* does not work in CentOS/Red  
98 Hat systems. The software is freely available to be run locally in personal computer.

99 *Seqs-Extractor* uses the following free and open source third-party software: BLAST+  
100 (Altschul et al. 1990), MISA (Thiel et al. 2003) and SAMTOOLS (Li et al. 2009).

## 101 **Results**

102 To easily perform a BLAST search in the command line version, *Seqs-Extractor* provides  
103 a friendly interface where all users should inform some minimum required parameters.

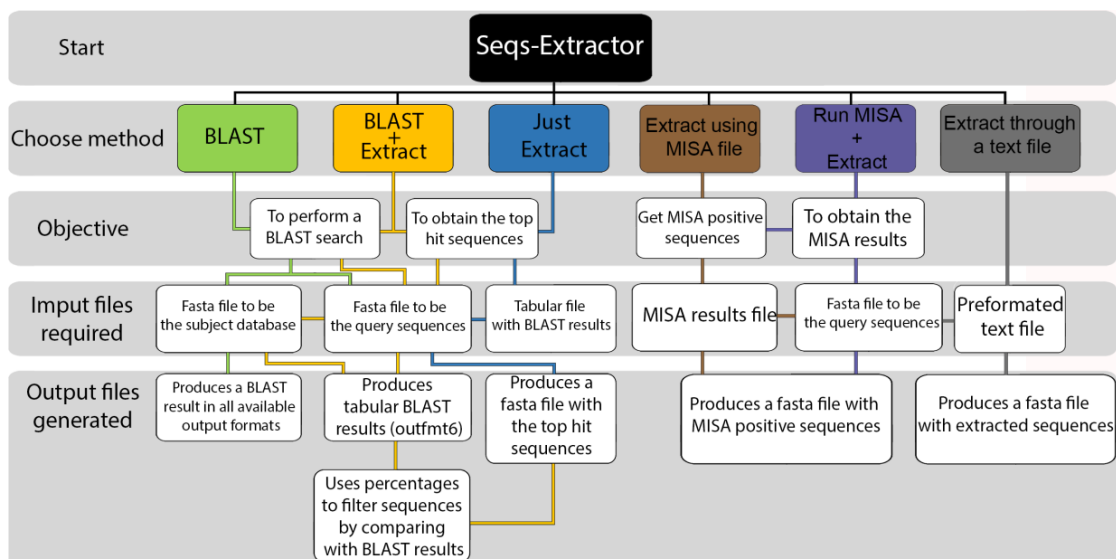
104 To test the script operation we used a .FASTA file that contains 107,185 sequences  
105 of the entire mouse genome (*Mus musculus*) obtained in NCBI database (NCBI  
106 Annotation Release 106) and run a BLASTX search against 168,031 sequences of revised  
107 and manually annotated mouse protein database obtained in (The universal protein  
108 resource (UniProt)). After BLASTX ran the resulting file contained 123,371 lines with  
109 the search results (considering expected value of  $1^{-3}$ ), where each line represents a result  
110 of alignment. To get only the sequences that show 100% homology with the subject we  
111 used *Seqs-Extractor* generates a new .FASTA file containing only 47,184 sequences with  
112 100% confirmed results aligned with the subject database.

113 Thus, *Seqs-Extractor* uses the query ID provided in the tabular file, to search inside  
114 the .FASTA file and extract only sequences that match in a specific percentage level  
115 defined by user. *Seqs-Extractor* will then generate two files: the extracted sequences and  
116 the tabular results of BLAST search. All methods use the BLASTN or BLASTX.

117 We also used MISA (Thiel et al. 2003) through *Seqs-Extractor* to perform a search  
118 (using default parameters) and extraction of microsatellites inside the .FASTA file that  
119 contains entire mouse genome. This run generated three output files: the MISA results  
120 file, the default MISA statistics file and the .FASTA file containing the positive MISA  
121 sequences.

122 In the extraction process *Seqs-Extractor* uses SAMTOOLS application (Li et al.  
 123 2009). This script works by comparing the data in the tabular BLAST format and  
 124 automatically search for sequences in the .FASTA file that match a percentage of  
 125 alignment specified by the user, creating a new .FASTA file containing only the  
 126 sequences filtered by percentage criteria. Similarly, SAMTOOLS extract sequences with  
 127 microsatellites by comparing data in the .FASTA file with IDs of these sequences in the  
 128 MISA file results. All tests performed and example files can be found in the *Seqs-*  
 129 *Extractor* webpage on Git-hub website ([https://github.com/patrick-douglas/Seqs-](https://github.com/patrick-douglas/Seqs-Extractor/wiki)  
 130 *Extractor/wiki*). The flowchart showing the *Seqs-Extractor* workflow is in Figure 1.

131 **Figure 1 – Flowchart.** Working steps of the six possible methods of *Seqs-Extractor*. The steps of each  
 132 method can be followed according to the color of the fluxes.



133

## 134 Conclusions

135 *Seqs-Extractor* is an automated tool that enhances BLAST analysis using a friendly  
 136 interface, besides extracting from a sequence dataset only 100% confirmed (BLASTed)  
 137 sequences. It also allows an extraction of positive microsatellites sequences from a MISA

138 search. It is a free of charge program that provides automatically search analysis in larger  
139 datasets that may help to minimize potential errors associated with tedious manual  
140 correction, thus allowing more effective results.

141

#### 142 **Acknowledgements**

143 We acknowledge Coordenação de Aperfeiçoamento de Pessoal de Nível Superior  
144 (CAPES), Instituto Federal do Pará (IFPA) and Universidade Federal do Pará (UFPA).

#### 145 **Competing interests**

146 The authors declare that they have no competing interests.

#### 147 **Authors' contributions**

148 Patrick D. C. Pereira: Development, source code analysis, manuscript conception  
149 writing

150 Cleysian Dias: Source code analysis, manuscript conception and writing, testing

151 Mauro A. D. Melo: Manuscript conception and writing, organization of data

152 Nara G. M. Magalhães: Manuscript conception and writing, organization of data, data  
153 examples building.

154 Cristovam Guerreiro-Diniz: Source code analysis, manuscript conception and writing,  
155 testing

156 Cristovam W. Picanço-Diniz: Manuscript conception, writing and final revision, testing,  
157 organization of data, bug fixes.

158 **Data Availability**

159 The following information was supplied regarding data availability:

160 *Seqs-Extractor* is available at:

161 <https://github.com/patrick-douglas/Seqs-Extractor/wiki>

162 **References**

163 Altschul SF, Gish W, Miller W, Myers EW, and Lipman DJ. 1990. Basic local alignment  
164 search tool. *Journal of molecular biology* 215:403-410.

165 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, and Madden TL.  
166 2009. BLAST+: architecture and applications. *BMC bioinformatics* 10:1.

167 Castoe TA, Poole AW, de Koning AJ, Jones KL, Tomback DF, Oyler-McCance SJ, Fike  
168 JA, Lance SL, Streicher JW, and Smith EN. 2012. Rapid microsatellite  
169 identification from Illumina paired-end genomic sequencing in two birds and a  
170 snake. *PLoS one* 7:e30953.

171 Conesa A, Götz S, García-Gómez JM, Terol J, Talón M, and Robles M. 2005. Blast2GO:  
172 a universal tool for annotation, visualization and analysis in functional genomics  
173 research. *Bioinformatics* 21:3674-3676.

174 Consortium U. 2008. The universal protein resource (UniProt. *Nucleic acids research*  
175 36:D190-D195.

176 Donkor ES, Dayie NT, and Adiku TK. 2014. Bioinformatics with basic local alignment  
177 search tool (BLAST) and fast alignment (FASTA). *Journal of Bioinformatics and*  
178 *Sequence Analysis* 6:1-6.

179 Excoffier L, and Lischer HE. 2010. Arlequin suite ver 3.5: a new series of programs to  
180 perform population genetics analyses under Linux and Windows. *Molecular*  
181 *ecology resources* 10:564-567.

182 Fernandez-Silva I, Whitney J, Wainwright B, Andrews KR, Ylitalo-Ward H, Bowen BW,  
183 Toonen RJ, Goetze E, and Karl SA. 2013. Microsatellites for next-generation  
184 ecologists: a post-sequencing bioinformatics pipeline. *PLoS one* 8:e55990.



- 185 Guichoux E, Lagache L, Wagner S, Chaumeil P, Léger P, Lepais O, Lepoittevin C,  
186 Malausa T, Revardel E, and Salin F. 2011. Current trends in microsatellite  
187 genotyping. *Molecular ecology resources* 11:591-611.
- 188 Khlestkina EK, Than MHM, Pestsova EG, Röder MS, Malyshev SV, Korzun V, and  
189 Börner A. 2004. Mapping of 99 new microsatellite-derived loci in rye (*Secale*  
190 *cereale* L.) including 39 expressed sequence tags. *Theoretical and Applied*  
191 *Genetics* 109:725-732.
- 192 Koochi-Dehkordi M, Sayed-Tabatabaei B, Yamchi A, and Danesh-Shahraki A. 2006.  
193 Microsatellite markers in pomegranate. XXVII International Horticultural  
194 Congress-IHC2006: II International Symposium on Plant Genetic Resources of  
195 Horticultural 760. p 179-184.
- 196 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, and  
197 Durbin R. 2009. The sequence alignment/map format and SAMtools.  
198 *Bioinformatics* 25:2078-2079.
- 199 Moghadam HK, Harrison PW, Zachar G, Székely T, and Mank JE. 2013. The plover  
200 neurotranscriptome assembly: transcriptomic analysis in an ecological model  
201 species without a reference genome. *Molecular ecology resources* 13:696-705.
- 202 Thiel T, Michalek W, and Varshney R. 2003. Exploiting EST databases for the  
203 development of cDNA derived microsatellite in barley (*Hordeum vulgare* L.).  
204 *Theoretical and Applied Genetics* 106:411-422.
- 205 Varshney RK, Graner A, and Sorrells ME. 2005. Genic microsatellite markers in plants:  
206 features and applications. *TRENDS in Biotechnology* 23:48-55.
- 207 Varshney RK, Thiel T, Stein N, Langridge P, and Graner A. 2002. In silico analysis on  
208 frequency and distribution of microsatellites in ESTs of some cereal species.  
209 *Cellular and Molecular Biology Letters* 7:537-546.
- 210 Yu J-K, Dake TM, Singh S, Benscher D, Li W, Gill B, and Sorrells ME. 2004.  
211 Development and mapping of EST-derived simple sequence repeat markers for  
212 hexaploid wheat. *Genome* 47:805-818.
- 213
- 214