1  **CONCATENATOR: SEQUENCE DATA MATRICES HANDLING MADE EASY**

2  **Authors:** F. Pina-Martins[1] and O. S. Paulo[1]

3  [1]Centro de Biologia Ambiental, Departmento de Biologia Animal, Faculdade de Ciências da
4  Universidade de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal.
5

7  **Corresponding author:**

8  Name: Francisco Rente de Pina Martins;

9  Address: Centro de Biologia Ambiental, Departmento de Biologia Animal, Faculdade de Ciências da

10  Universidade de Lisboa, Campo Grande, 1749-016 Lisboa, Portugal;

11  E-mail: f.pinamartins@gmail.com

17  **ABSTRACT**

18  *Concatenator* is a simple and user friendly software that implements two very useful functions for

19  phylogenetics data analysis. It concatenates Nexus files of several fragments in a single NEXUS file

20  ready to be used in phylogenetics softwares, such as PAUP and MrBayes and it converts FASTA

21  sequence data files to NEXUS and *vice-versa*. Additionally, concatenated files can be prepared for

22  partition tests in PAUP. It is freely available in http://cobig2.fc.ul.pt/.

23  **THE PROGRAM**

24     Sequence data files can be organized in many different formats. Different sequence analysis

25  software require differently formatted input files. The FASTA format has become very popular due to its

26  simplicity and the capacity to quickly compare sequences (Pearson & Lipman 1988); these

27  characteristics made this format one of the NCBI default outputs. The Nexus format became popular

28  due to its modular format which is at the same time flexible and standardized (Maddison *et al.* 1997).

29    The existence of different file formats for the same data types require investigators to know how to

30    handle them since they are not shared by some of the most common phylogenetic analysis software.

31    *Concatenator*'s main purpose is to turn data matrix handling into a simple task, by allowing intuitive

32    format conversions and concatenations of data matrices.

33    *Concatenator* is written in Perl using the Perl/TK module in order to give it a GUI for simplifying

34    usage. The software was compiled for Windows using the PAR module. It is available in the Win32

35    binary version and source code version at the authors' group website.

36    The only requirements are either a system with a Perl interpreter and the Tk module installed

37    (source code version) or a system running Microsoft Windows XP (not tested on other versions). The

38    software was developed with a very specific aim – the simple handling of data matrices from one

39    program to another and the concatenation of several of these data matrices. All the functions that the

40    program performs can be accomplished manually provided the user has some knowledge about the

41    involved file formats; however, even in such case this process is very error prone due to complex data

42    organization such as in the interleave Nexus format.

43    The user interface is very simple (**Fig. 1**) and consists of a window that accommodates essentially

44    the input and output entry boxes; these files can be selected from the File menu, a browse button

45    located on the right of every entry box or by entering the path and filename directly on the entry box.

46    *Concatenator* can be used to accomplish 2 essential tasks chosen from the welcome window buttons

47    or from it's "File" menu.

48    (1)    Fasta-Nexus-Fasta Converter – It converts files from FASTA to Nexus format and from Nexus

49            to FASTA format. When converting from Nexus to FASTA, there are no options available to

50            chose from, however, when converting from FASTA to Nexus the user can choose whether to

51            include a Taxa block, a leave or interleave organization, the type of data, the character for

52            missing data and the gap character. File comments are ignored when converting.

53    (2)    Matrix Concatenator – This function takes 2 to 5 Nexus formatted matrices and concatenates

54            them into a single file. Two output formats are possible, one formatted to be used with *PAUP\**

55            (Swofford 2003), and the other prepared to input to *MrBayes* (Ronquist & Huelsenbeck 2003).

56            Several parameters are customizable such as the inputs' data type, the gap character, the

57            missing character, whether or not to include the Taxa block and a pre input for performing a

58            "partition test" in *PAUP\** excluding constant characters. Should the input files contain different

59       taxa, the program will prompt the user about what to do with the "unpaired" taxa – whether to

60       ignore them in the final concatenated matrix, or to fill the missing parts with the character used

61       in the file for "missing data".

62 Each function has a help file. The whole program is simple to use, but the help files are nevertheless

63 as descriptive as possible.

64 An example of the program usage could be the following:

65 The user downloads two arrays of sequences (e.g. two different genes from the same species) from

66 the NCBI database using a program such as *BioEdit* (Hall 1999). After a proper alignment session, the

67 program outputs two FASTA files – one for each gene.

68 The user then wants to analyze these files using *PAUP\**, *MrBayes*, *TCS* (Clement *et al*. 2000) or

69 Network (Bandelt *et al.* 1999). *Concatenator* is useful in this step, because it provides a simple way to

70 convert these FASTA files into Nexus files, ready to use in the analysis programs.

71 If after this first analysis the user decides to analyze both genes as a single block, *Concatenator* can

72 join the two Nexus files in a single data matrix, ready to input on software such *MrBayes* or *PAUP\**; the

73 built in function for data partitioning will automatically add the required commands for partitioning data

74 required for a "Partition Test" in *PAUP\**.

79 **REFERENCES**

80 Bandelt HJ, Forster P, Röhl A (1999) Median-joining networks for inferring intraspecific phylogenies. In:

81 Molecular Biology and Evolution 16:37-48.

82 Clement M, Posada D, Crandall K (2000) TCS: a computer program to estimate gene genealogies. In:

83 Molecular Ecology 9(10): 1657-1660.

84 Hall TA (1999)  BioEdit: a user-friendly biological sequence alignment editor and analysis program for

85 Windows 95/98/NT. In: Nucleic Acids Symposium Series 41:95-98.

86 Maddison DR, Swofford DL, Maddison WP (1997) NEXUS: An Extensible File Format for Systematic

87 Information. In: Systematic Biology 46 4: 590-621.

88    Ronquist F, Huelsenbeck JP (2003) MrBayes 3: Bayesian phylogenetic inference under mixed models.

89    In: Bioinformatics 19:1572-1574.

90    Swofford DL (2003) PAUP*. Phylogenetic Analysis Using Parsimony (*and Other Methods). Version 4.

91    Sinauer Associates, Sunderland, Massachusetts.

92    Pearson WR, Lipman DJ (1988) Improved tools for biological sequence comparison. In: Biochemistry
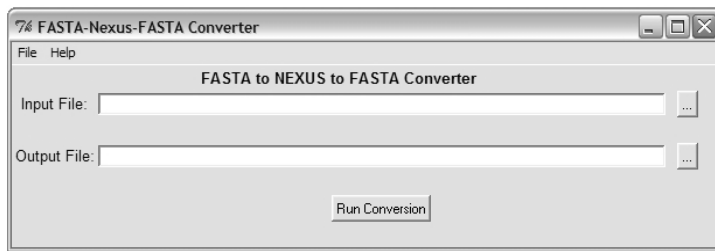
93    85: 2444-2448.

94

95



96

97

98

99

100 **Fig. 1.** *Concatenator* interface on Win XP when converting matrices from Nexus to FASTA. Should the

101 input be a FASTA file, the interface changes to accommodate the extra options to outputting a Nexus

102 file.