1 **TreeRipper: towards a fully automated optical tree**

2 **recognition software**

3

4 **Joseph Hughes[1]***

5

6 [1]Division of Ecology and Evolutionary Biology, Faculty of Biomedical and Life Sciences,

7 University of Glasgow, Graham Kerr Building, University Avenue, Glasgow, G12 8QQ, UK

8

9 *Corresponding author

10

11 Email addresses:

12       JH: j.hughes@bio.gla.ac.uk

13

# 14 **Abstract**

## 15 **Background**

16 Relationships between species, genes and genomes have been printed as trees for over a

17 century. Whilst this may have been the best format for exchanging and sharing phylogenetic

18 hypotheses during the 20[th] century, the worldwide web now provides faster and automated

19 ways of transferring and sharing phylogenetic knowledge. However, novel software is needed

20 to defrost these published phylogenies for the 21[st] century.

## 21 **Results**

22 TreeRipper is a command line c++ program for the fully-automated recognition of

23 multifurcating phylogenetic trees. The program accepts a range of input image formats  (PNG,

24 JPG/JPEG, GIF, TIFF or PDF ). Then follows a number of cleaning steps to detect lines,

25 remove node labels, patch-up broken lines and corners and detect line edges. The edge contour

26 is then determined to detect the branch length, tip label positions and the topology of the tree.

27 Optical Character Recognition (OCR) is used to convert the tip labels into text with the freely

28 available tesseract-ocr software. 32% of images meeting the prerequisites for TreeRipper were

29 successfully recognised, the largest tree had 115 leaves.

2

30**Conclusions**

31Despite the diversity of ways phylogenies have been illustrated making the design of a fully

32automated tree recognition software difficult, TreeRipper is a step towards automating the

33digitization of past phylogenies.  We also provide a dataset of 100 tree images and associated

34tree files for training and/or benchmarking future software.

35

36# Background

37In 1859, Darwin produced one of the first illustrations of a phylogenetic tree, notably this was

38the only figure included in *The Origin of Species* . Since, biologists have used trees to depict

39the relationships between organisms, genes and genomes. The number of studies depicting

40phylogenies exploded (see Figure 1) with the development of the polymerase chain reaction

41technique and journals were created specifically for publishing the molecular phylogenies

42generated by researchers  (e.g., Molecular Phylogenetics and Evolution established in 1992).

43Whilst in the early years of molecular phylogenetics, embedding illustrations into manuscripts

44might have been the most appropriate way to disseminate knowledge, this has resulted in the

45locking up of phylogenetic hypotheses into the pages of journals and books without an easy

46way to access this information.

47Currently, the construction of the relationships between the 1.8 million currently estimated

48species largely depends on the unprecedented growth of molecular sequence data  and this

3

49makes GenBank the most accessible source of comparative data for most taxa in the tree of life

50. Whilst more sequence data, more powerful computers and improved phylogenetic

51reconstruction algorithms will enable researcher to generate up-to-date phylogenies from the

52raw data in the future, past phylogenetic inferences will remain central to guiding researchers

53towards studying poorly supported relationships and under-sampled lineages. They are also

54central for studying the effects of new phylogenetic methodologies and new and larger datasets

55.

56Not all phylogenetically informative data are confined to sequence databases. TreeBASE is a

57very valuable repository in that respect as it holds morphological or genetic data with the

58associated published phylogeny . However, as few publishers require submission to TreeBASE

59as a pre-requisite for publication, a large number of phylogenies remain embedded as images in

60published articles. Indeed, the rapid growth of published phylogenies is not matched by the

61availability of those trees in databases (see Figure 1 in ).

62The idea of using a program to convert a tree image into a computer-readable representation of

63that tree was first implemented in TreeThief which required the user to trace a tree by clicking

64on each of its nodes in turn. The latter program is only available for the discontinued operating

65system Mac OS 9. Laubach and von Haeseler provided a conceptual advance with a semi-

66automatic program called TreeSnatcher that has recently been updated .

67 Here, we will review the way researchers present their phylogenies, demonstrate the feasibility

68 of a fully automated tree recognition software and provide a dataset of tree images and

69 associated tree files for training and/or benchmarking future programs.

70

## 71 Implementation

72 The current version of TreeRipper opens tree-image files in the formats PNG, JPG/JPEG, GIF,

73 TIFF or PDF.

74 • The tree needs to have the root on the left and leaves on the right.

75 • Horizontal branches.

76 • The tree constitutes a dark foreground on a light homogenous background (no

77 background boxes or shading).

78 • The tree must be bi- or multifurcating (not a network)

79 • The inner nodes are branching points between lines and have no circles, rectangles, etc.

80 inscribed.

81 • Tip branches must have branch lengths greater than 0.

82

83 TreeRipper is written in c++ using a set of Standard Template Library algorithms provided by

84 Magick++. The image is first converted to black and white and rescaled so that horizontal lines

85 are on average 2 pixels thick. The image is cleaned by removing a series of patterns such as

86black pixels surrounded by a box of white pixels and horizontal lines that are not connected to

87vertical lines. Lines and corners are then patched up before the contour is traced and the

88topology detected. The locations of branch tips are then used to crop the tip labels from the

89original image. Tip labels are converted to text using the freely available tesseract-ocr prgram.

90The steps in the program are depicted in Figure 2.

## 91Results and Discussion

92We downloaded 322 images which had phylogen* or supertree in their caption from 249

93articles published in the Open Access journal BMC Evolutionary Biology between 1997 and

942009.  Only eleven out of these 249 articles have submitted their alignment and tree files to

95TreeBASE. All images were visually inspected to check whether the image met the

96prerequisites. Twenty-four images were not phylogenies, 26 were represented as radial tree

97layouts, 8 as polar tree layouts and 5 as cladograms. Of those represented with a rectangular

98tree layout, 40 had background boxes, 31 had lines intersecting branches or branches drawn

99with dotted or dashed lines, 32 had circles or boxes as nodes, 6 were illustrated over multiple

100pages, 4 had triangles as tip leaves, 3 had leaves with zero branch lengths. A further 29 would

101need some form of pre-processing (rotating or splitting into component images). Of the 298

102images of phylogenies downloaded only 114 (38%) would meet the prerequisites for this

103program, which are very similar to those of the semi-automatic recognition software

104TreeSnatcher . This small proportion of the total phylogenetic images illustrates the plethora of

105 ways trees are currently represented in one journal alone. Of the 114 phylogenies that meet the

106 prerequisites, the topologies of 37 trees (i.e., 32%) were successfully recognized by TreeRipper

107 without any prior processing. The largest phylogeny successfully recognised had 115 leaves.

108 We do not review the accuracy of the OCR here as it has been done elsewhere (see ).

109 The successfully recognised tree images along with a further 63 images manually converted to

110 tree files  are provided as supplementary material in NEXUS, Newick and phyloXML formats

111 (Additional file 1) for training and/or benchmarking future programs.

112

## 113 Conclusions

114 Although the program has a high failure rate, it is the first step towards an automated approach

115 for optical tree recognition and proves the feasibility of an approach which might permit us to

116 defrost published phylogenetic hypotheses. We are unlikely to ever be able to create an

117 application that recognises all possible trees due to the sheer diversity of ways phylogenies

118 have been illustrated but at the very least, this program could be used for automating tree

119 recognition of large sets of tree images before using manual conversion or semi-automated

120 programs like TreeSnatcher.

121 As phylogenetics enters a third phase of growth with the advent of next-generation sequencing,

122 one hopes that the work of future phylogenetists will be published in a format that will enable

123 the digital curation and preservation of their hard work.

124

## 125 Availability and requirements

126 Project name: TreeRipper (automated phylogeny recognition from images)

127 Project home page: https://code.google.com/p/treeripper/

128 Programming language: C++

129 **Prerequisites**

130 Tesseract-OCR  licensed with the Apache 2.0 License except the tesseractTrainer.py, which is

131 licensed with GPL: http://code.google.com/p/tesseract-ocr

132 Imagemagick, license is compatible with the GPL: http://www.imagemagick.org/

133

## 134 Authors' contributions

135 JH developed the idea, wrote the code, tested the software and drafted the manuscript.

136

## 142References

143

144

145

## 146Figures

147**Figure 1  - Percentage of articles with phylogeny* in the title**

148The percentage of articles with phylogen* in the title out of the total number of publication for

149each year since 1980 from PubMed.

150**Figure 2  - Architecture of the software design for TreeRipper**

151The input image is scaled, node labels are removed, branches are smoothed and corners

152patched-up, the contour is detected. Tips locations are used to determine leaf label boxes for

153which the text is recognised using Tesseract. TreeRipper summarizes the tree topology and

154labels in a text file and an SVG file which shows the contours.

## 155Additional files

**156Additional file 1 – Tree images and associated newick file**

157Set of images and associated nexus tree file as a zip file.

158