# FUSTr: a tool to find gene families under selection in transcriptomes

Timothy J Cole [Corresp., 1] , Michael S Brewer [Corresp. 1]

[1] Department of Biology, East Carolina University, Greenville, North Carolina, United States

Corresponding Authors: Timothy J Cole, Michael S Brewer
Email address: coleti16@students.ecu.edu, brewermi14@ecu.edu

**Background:** The recent proliferation of large amounts of biodiversity transcriptomic data has resulted in an ever-expanding need for scalable and user-friendly tools capable of answering large scale molecular evolution questions. FUSTr identifies gene families involved in the process of adaptation. This is a tool that finds genes in transcriptomic datasets under strong positive selection that automatically detects isoform designation patterns in transcriptome assemblies to maximize phylogenetic independence in downstream analysis.

**Results:** When applied to previously studied spider transcriptomic data as well as simulated data, FUSTr successfully grouped coding sequences into proper gene families as well as correctly identified those under strong positive selection in relatively little time.

**Conclusions:** FUSTr provides a useful tool for novice bioinformaticians to characterize the molecular evolution of organisms throughout the tree of life using large transcriptomic biodiversity datasets and can utilize multi-processor high-performance computational facilities.

1 **FUSTr: a tool to find gene Families Under Selection in Transcriptomes**

2 T. Jeffrey Cole[1,*], Michael S. Brewer[1]

3 [1]Department of Biology, East Carolina University, Greenville, NC 27858.

4 **Correspondence:** coleti16@students.ecu.edu

5

6 **Abstract**

7     **Background:** The recent proliferation of large amounts of biodiversity transcriptomic

8 data has resulted in an ever-expanding need for scalable and user-friendly tools capable of

9 answering large scale molecular evolution questions. FUSTr identifies gene families involved in

10 the process of adaptation. This is a tool that finds genes in transcriptomic datasets under strong

11 positive selection that automatically detects isoform designation patterns in transcriptome

12 assemblies to maximize phylogenetic independence in downstream analysis.

13     **Results:** When applied to previously studied spider transcriptomic data as well as

14 simulated data, FUSTr successfully grouped coding sequences into proper gene families as well

15 as correctly identified those under strong positive selection in relatively little time.

16     **Conclusions:** FUSTr provides a useful tool for novice bioinformaticians to characterize

17 the molecular evolution of organisms throughout the tree of life using large transcriptomic

18 biodiversity datasets and can utilize multi-processor high-performance computational facilities.

19 **Background**

20    Elucidating patterns and processes involved in the adaptive evolution of genes and

21    genomes of organisms is fundamental to understanding the vast phenotypic diversity found in

22    nature. Recent advances in RNA-Seq technologies have played a pivotal role in expanding

23    knowledge of molecular evolution through the generation of an abundance of protein coding

24    sequence data across all levels of biodiversity (Todd, Black & Gemmell, 2016). In non-model

25    eukaryotic systems, transcriptomic experiments have become the *de facto* approach for functional

26    genomics in lieu of whole genome sequencing. This is due largely to lower costs, better targeting

27    of coding sequences, and enhanced exploration of posttranscriptional modifications and

28    differential gene expression (Wang, Gerstein & Snyder, 2009). This influx of transcriptomic data

29    has resulted in a need for scalable tools capable of elucidating broad evolutionary patterns in

30    large biodiversity datasets.

31    Billions of years of evolutionary processes gave rise to remarkably complex genomic

32    architectures across the tree of life. Numerous speciation events along with frequent whole

33    genome duplications have given rise to a myriad of multigene families with varying roles in the

34    processes of adaptation (Benton, 2015). Grouping protein encoding genes into their respective

35    families *de novo* has remained a difficult task computationally. This typically entails homology

36    searches in large amino acid sequence similarity networks with graph partitioning algorithms to

37    cluster coding sequences into transitive groups (Andreev & Racke, 2006). This is further

38    complicated in eukaryotic transcriptome datasets that contain several isoforms via alternative

39    splicing (Matlin, Clark & Smith, 2005). Further exploration of Darwinian positive selection in

40    these families is also nontrivial, requiring robust Maximum Likelihood and Bayesian

41    phylogenetic approaches.

42    Here we present a fast tool for finding Families Under Selection in Transcriptomes

43    (FUSTr), to address the difficulties of characterizing molecular evolution in large-scale

44    transcriptomic datasets. FUStr can be used to classify selective regimes on homologous groups of

45    phylogenetically independent coding sequences in transcriptomic datasets and has been verified

46    using large transcriptomic datasets and simulated datasets. The presented pipeline implements

47    simplified user experience with minimized third-party dependencies, in an environment robust to

48    breaking changes to maximize long-term reproducibility.

49        While FUSTr fills a novel niche among sequence evolution pipeline, a recent tool, VESPA

50    (Webb et al., 2017), performs several similar functions. Our tool differs in that it can accept *de*

51    *novo* transcriptome assemblies that are not predicted ORFs. VESPA requires nucleotide data to be

52    in complete coding frames, and does not filter isoforms or utilize transitive clustering to deal with

53    domain chaining. Additionally, VESPA makes use of slow maximum likelihood methods for tests

54    of selection and provides no information about purifying selection, whereas FUSTr utilizes a Fast

55    Unconstrained Bayesian Approximation (FUBAR) (Murrell et al., 2013) to analyze both

56    pervasive and purifying regimes of selection.


57    **Implementation**

58        FUSTr is written in Python with all data filtration, preparation steps, and command line

59    arguments/parameters for external programs contained in the workflow engine Snakemake

60    (Köster & Rahmann, 2012). Snakemake allows FUSTr to operate on high performance

61    computational facilities, while also maintaining ease of reproducibility. FUSTr and all third-party

62    dependencies are distributed as a Docker container (Merkel, 2014).  FUSTr contains ten

63    subroutines that takes transcriptome assembly FASTA formatted files from any number of taxa as

64    input and infers gene families that are either under diversifying or purifying selection. A graphical

65    overview of this workflow and parallelization scheme has been outlined in Fig. 1.

66        *Data Preprocessing* The first subroutine of FUSTr acts as a quality check step to ensure

67    input files are in valid FASTA format. Spurious special characters resulting from transferring text

68    files between multiple operating system architectures are detected and removed to facilitate

69    downstream analysis.

70         *Isoform detection* Header patterns are analyzed to auto-detect whether the given assembly

71    includes isoforms by detecting naming convention redundancies commonly used in isoform

72    designations, in addition to comparing the header patterns to common assemblers such as Trinity

73    *de novo* assemblies (Haas et al., 2013) and Cufflinks reference genome guided assemblies

74    (Trapnell et al., 2014).

75         *Gene prediction* Coding sequences are extracted from transcripts using Transdecoder

76    v3.0.1 (Haas et al., 2013). Transdecoder predicts Open Reading Frames (ORFs) using likelihood-

77    based approaches. A single best ORFs for each transcript with predicted coding sequence is

78    extracted providing nucleotide coding sequences (CDS) and complementary amino acid

79    sequences. This facilitates further analyses requiring codon level sequences while using the more

80    informative amino acid sequences for homology inferences and multiple sequence alignments. If

81    the data contain several isoforms of the same gene, at this point only the longest isoform is kept

82    for further analysis to ensure phylogenetic independence. The user may customize the use of

83    Transdecoder by changing minimum coding sequence length (default: 30 codons) or strand-

84    specificity (default: off). Users also have the option to only retain ORFs with homology to known

85    proteins through a BLAST search against Uniref90 or Swissprot in addition to searching PFAM

86    to identify common protein domains.

87         *Homology search* All coding sequences are assigned a unique identifier and then

88    concatenated into one FASTA file. Homology of peptide sequences is assessed via BLASTP

89    acceleration through DIAMOND (v.0.9.10) with an e-value cutoff of $10^{-5}$.

90         *Gene Family inference* The resulting homology network is grouped into putative gene

91    families using transitive clustering with SiLiX v.1.2.11, which is faster and has better memory

92    allocation than other clustering algorithms such as MCL, and greatly reduces the problem of

93    domain chaining (Miele, Penel & Duret, 2011). Sequences are only added to a family with 35%

94    minimum identity, 90% minimum overlap, with minimum length to accept partial sequences in

95    families as 100 amino acids, and minimum overlap to accept partial sequences of 50%. These are

96    the optimal configurations of SiLiX (Bernardes et al., 2015), but the user is free to configure

97    these options.

98         *Multiple sequence alignment and phylogenetic reconstruction* Multiple amino acid

99    sequence alignments of each family are then generated using the appropriate algorithm

100   automatically detected using MAFFT v7.221 (Katoh & Standley, 2013). Spurious columns in

101   alignments are removed with Trimal v1.4.1's *gappyout* algorithm (Capella-Gutiérrez & Silla-

102   Martínez, 2009). Phylogenetic reconstruction of each family's untrimmed protein multiple

103   sequence alignment using FastTree v2.1.9 (Price, Dehal & Arkin, 2010). Trimmed multiple

104   sequence codon alignments are then generated by reverse translation of the amino acid alignment

105   using the CDS sequences.

106        *Tests for selective regimes* Families containing at least 15 sequences have the necessary

107   statistical power for tests of adaptive evolution (Wong et al., 2004). Tests of pervasive positive

108   selection at site specific amino acid level are implemented with FUBAR (Murrell et al., 2013).

109   Unlike codeml, FUBAR allows for tests of both positive and negative selection using an ultra-fast

110   Markov chain Monte Carlo routine that averages over numerous predefined site-classes. When

111   compared to codeml, FUBAR performs as much as 100 times faster (Murrell et al, 2013). Default

112   settings for FUBAR, as used in FUSTr, include twenty grid points per dimension, five chains of

113   length 2,000,000, with the first 1,000,000 used as burn-in, 100 samples drawn from each chain,

114   and concentration parameter of the Dirichlet prior set to 0.5.

115        Users have the option to also run tests for pervasive selection using the much slower

116   CODEML v4.9 (Yang, 2007) with the codon alignments and inferred phylogeny. Log-likelihood

117   values of codon substitution models that allow positive selection are then compared to respective

118    nested models not allowing positive selection (M0/M3, M1a/M2a, M7/M8, M8a/M8), Bayes

119    Empirical Bayes (BEB) analysis then determines posterior probabilities that the ratio of

120    nonsynonymous to synonymous substitutions ($d_N/d_S$) exceeds one for individual amino acid sites.

121            *Final output and results* The final output is a summary file describing what gene families

122    were detected, and those that are under strong selection and the average $d_N/d_S$ per family. A CSV

123    file for each family under selection is generated giving the following details per codon position of

124    the family alignment : alpha mean posterior synonymous substitution rate at a site; beta mean

125    posterior non-synonymous substitution rate at a site; mean posterior beta-alpha; posterior

126    probability of negative selection at a site; posterior probability of positive selection at a site;

127    Empiricial Bayes Factor for positive selection at a site; potential scale reduction factor; and

128    estimated effective sample site for the probability that beta exceeds alpha.


129    **Validation**

130            We tested FUSTr on six published whole body transcriptome sequences from an adaptive

131    radiation of Hawaiian *Tetragnatha* spiders (NCBI Short Read Archive Assesion numbers:

132    SRX612486, SRX612485, SRX612477, SRX612466, SRX559940, SRX559918) assembled

133    using the same methods from the original publication (Brewer et al., 2015). Spider genomes

134    contain numerous gene duplications lending to gene family rich transcriptomes. Additionally, this

135    adaptive radiation has been shown to facilitate strong, positive, sequence-level selection in these

136    transcriptomes (Brewer et al., 2015). This dataset provides an ideal case use for FUSTr.

137            A total of 273,221 transcripts from all six *Tetragnatha* samples were provided as input for

138    FUSTr, a total of 4,258 isoforms were removed leaving 159,464 coding sequences for analysis

139    after gene prediction. The entire analysis ran in 13.7 core hours, completing within an hour when

140    executed on a 24-core server. Time of completion and memory usage for each of FUSTr's

141    subroutines performance in this analysis is reported in Table 1. FUSTr recovered 134 families

142    containing at least 15 sequences, of these 46 families contained sites under pervasive positive

143    selection while all families also contained sites under strong purifying selection. This can be

144    contrasted to the analysis by Brewer et al. (2015) which found 2,647 one-to-one six-member

145    orthologous loci (one ortholog per each of the same samples), with 65 loci receiving positive

146    selection based on branch-specific analysis. The original analysis did not allow paralogs whereas

147    FUSTr does not reconstruct one-to-one orthogroups but entire putative gene families, and the

148    selection analysis utilized by FUSTr is site-specific and not branch-specific. Thus, it is not

149    expected that the results from FUSTr would perfectly match up with the original analysis,

150    however five of the 46 families FUSTr found to be under selection included loci from Brewer et

151    al.'s (2015) original 65 under selection based on branch-specific analysis.

152        The same 273,221 transcripts were entered as input for VESPA as a comparative analysis.

153    Because VESPA cannot filter Open Reading Frames in transcripts, it was unable to infer proper

154    coding sequences. In its first phase of cleaning input fasta files, 86,269 transcripts were wrongly

155    removed for having "internal stop codons" via improper reading frame inference, and 182,000

156    transcripts were removed due to "abnormal sequence length". Approximately 98% of the

157    transcripts were removed in the first phase of VESPA with no gene predictions, rendering further

158    analysis unnecessary for proper comparison of the performance of the two pipelines.

159        We further validated FUSTr by utilized coding sequences from simulated gene families

160    with predetermined selective regimes. We used EvolveAGene (Hall, 2007) on 3,000 random

161    coding sequences of a random length of 300-500 codons to generate gene families containing 16

162    sequences evolved along a symmetric phylogeny each with average branch lengths chosen

163    randomly between 0.01-0.20 evolutionary units. Selective regimes with a selection modifier of

164    3.0 were randomly chosen for each family so that a random 10% partition of the family receive

165    pervasive positive selection, purifying selection, or constant selection. All other settings for

166    EvolveAGene were left as their defaults: the probability of accepting an insertion was set to the

167    default 0.1, the probability of accepting a deletion defaulted to 0.025, the probability of accepting

168    a replacement was left at 0.016, no recombination was allowed. A visual schema for these

169    simulations can be found in Fig. 2.

170        The resulting 48,000 simulated sequences were used as input for FUSTr with

171    Transdecoder set to be strand-specific. FUSTr correctly recovered all 3,000 families, and all 975

172    that were randomly selected to undergo strong positive selection were correctly classified as

173    receiving pervasive positive selection. Additionally, the families selected to undergo purifying

174    selection were correctly classified, and families not selected receive constant selection were

175    classified as not having any specific sites undergoing purifying or pervasive positive selection.

176    Scripts for these simulations can be found at https://github.com/tijeco/FUSTr.


177    **Conclusions**

178        Current advances in RNA-seq technologies have allowed for a rapid proliferation of

179    transcriptomic datasets in numerous non-model study systems. It is currently the only tool

180    equipped to deal with the nuances of transcriptomic data, allowing for proper prediction of gene

181    sequences and isoform filtration. FUSTr provides a fast and useful tool for novice

182    bioinformaticians to detect gene families in transcriptomes under strong selection.  Results from

183    this tool can provide information about candidate genes involved in the processes of adaptation,

184    in addition to contributing to functional genome annotation.


185    **Availability:** FUSTr is freely available under a GNU license and can be downloaded at

186    https://github.com/tijeco/FUSTr.


187    **Acknowledgements**
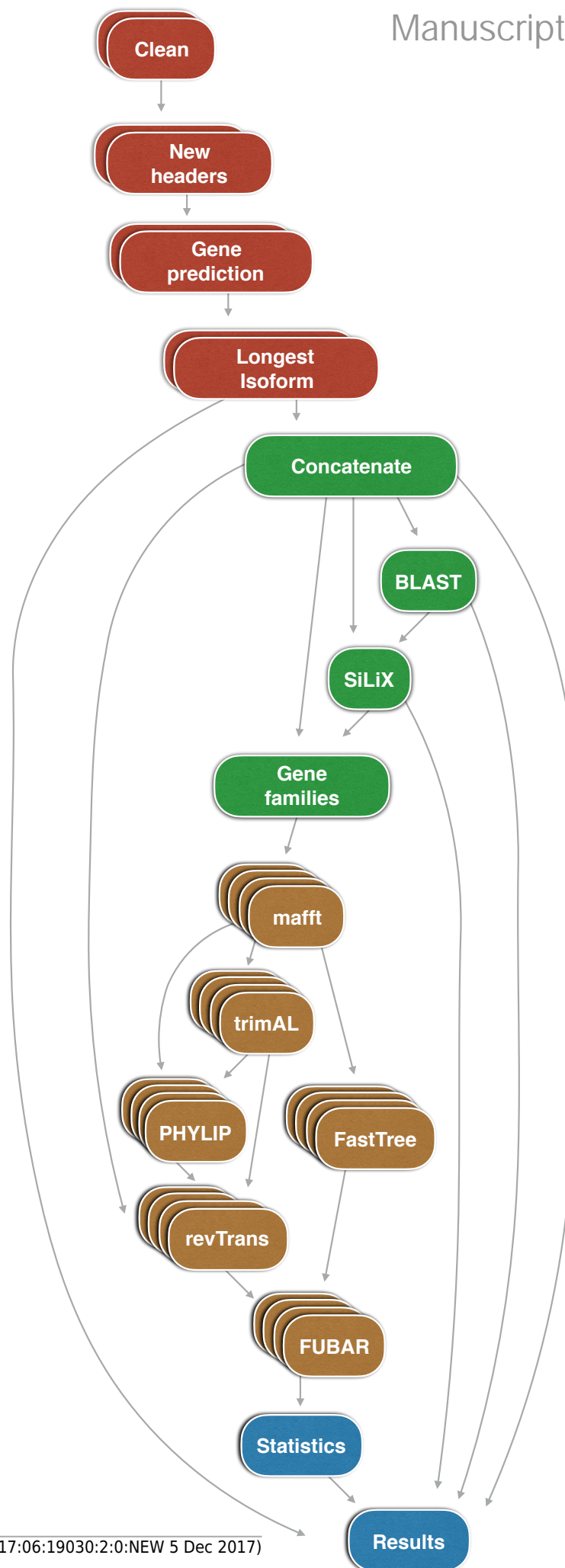
190    **References**

191    Andreev K, Racke H. 2006. Balanced Graph Partitioning. *Theory of Computing Systems* 39:929–
192          939. DOI: 10.1007/s00224-006-1350-7.
193    Benton R. 2015. Multigene Family Evolution: Perspectives from Insect Chemoreceptors. *Trends
194          in Ecology & Evolution* 30:590–600. DOI: 10.1016/j.tree.2015.07.009.
195    Bernardes JS., Vieira F.R., Costa L.M., Zaverucha G. 2015. Evaluation and improvements of
196          clustering algorithms for detecting remote homologous protein families. *BMC
197          Bioinformatics* 16: 34.
198    Brewer MS, Carter RA, Croucher PJP, and Gillespie, RG. 2015, Shifting habitats, morphology,
199          and selective pressures: Developmental polyphenism in an adaptive radiation of Hawaiian
200          spiders. *Evolution*, 69: 162–178. doi:10.1111/evo.12563
201    Capella-Gutiérrez S, Silla-Martínez  JM. 2009. trimAl: a tool for automated alignment trimming
202          in large-scale phylogenetic analyses.  ….
203    Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D,
204          Li B, Lieber M, Macmanes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N,
205          Westerman R, William T, Dewey CN, Henschel R, Leduc RD, Friedman N, Regev A.
206          2013. De novo transcript sequence reconstruction from RNA-seq using the Trinity
207          platform for reference generation and analysis. *Nature protocols* 8:1494–512. DOI:
208          10.1038/nprot.2013.084.
209    Hall B. 2007. EvolveAGene 3: A DNA coding sequence evolution simulation program. *Nature
210          Precedings*. DOI: 10.1038/npre.2007.1230.1.
211    Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
212          improvements in performance and usability. *Molecular biology and evolution*.
213    Köster J, Rahmann S. 2012. Snakemake—a scalable bioinformatics workflow engine.
214          *Bioinformatics* 28:2520–2522.
215    Matlin A, Clark F, Smith C. 2005. Understanding alternative splicing: towards a cellular code.
216          *Nature Reviews Molecular Cell Biology*:386–398. DOI: 10.1038/nrm1645.
217    Merkel D. 2014. Docker: lightweight linux containers for consistent development and
218          deployment Linux J. 239 2
219    Miele V, Penel S, Duret L. 2011. Ultra-fast sequence clustering from similarity networks with
220          SiLiX. *BMC Bioinformatics* 12:1–9. DOI: 10.1186/1471-2105-12-116.
221    Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Pond S, Scheffler K. 2013. FUBAR: A
222          Fast, Unconstrained Bayesian AppRoximation for Inferring Selection. *Molecular Biology
223          and Evolution* 30:1196–1205. DOI: 10.1093/molbev/mst030.
224    Price MN, Dehal  PS, Arkin  AP. 2010. FastTree 2–approximately maximum-likelihood trees for
225          large alignments. *PloS one*. DOI: 10.1371/journal.pone.0009490.
226    Todd E, Black M, Gemmell N. 2016. The power and promise of RNA‑seq in ecology
227          and evolution. *Molecular Ecology*:1224–1241. DOI: 10.1111/mec.13526.
228    Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley D, Pimentel H, Salzberg S, Rinn J,
229          Pachter L. 2014. Differential gene and transcript expression analysis of RNA-seq
230          experiments with TopHat and Cufflinks. *Nature Protocols*:562–578. DOI:
231          10.1038/nprot.2012.016.
232    Wang Z, Gerstein M, Snyder M. 2009. RNA-Seq: a revolutionary tool for transcriptomics.
233          *Nature Reviews Genetics*:57–63. DOI: 10.1038/nrg2484.
234    Webb, Andrew E., Thomas A. Walsh, and Mary J. O'Connell. 2017. VESPA: Very large-scale

235          evolutionary and selective pressure analyses. *PeerJ Computer Science* 3: e118.

236    Wong W. S. W., Yang Z., Goldman N., & Nielsen R. 2004. Accuracy and Power of Statistical

237          Methods for Detecting Adaptive Evolution in Protein Coding Sequences and for

238          Identifying Positively Selected Sites. *Genetics*, 168(2), 1041–1051.

239          http://doi.org/10.1534/genetics.104.031153

240    Yang Z. 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Molecular Biology and*

241          *Evolution* 24:1586–1591. DOI: 10.1093/molbev/msm088.

**Figure 1**(on next page)

Parallelization scheme and workflow of FUSTr.

Color coding denotes functional subroutines in the pipeline: preparation and open reading frame prediction (red); homology inferenece and gene family clustering (green); multiple sequence alignment, phylogenetics, and selection detection (brown); and model selection and reconciliation (blue).

**Figure 2**(on next page)

Schematice used for EvolveAGene.

Randomly generated sequence is evolved along a symmetric phylogeny of a given selective regime.

1. Seed sequence (300-500 codons)

2. Evolve gene along 16 node phylogeny
($\bar{x}$ branch length of 0.01-0.20)

3. Predefine selective regime
(positive, purifying, constant)
across sequences

# Table 1(on next page)

Benchmarks for each subroutines' time and memory used for the *Tetragnatha* transcriptome assembly analysis.

Red highlighted row represents subroutine consuming the most memory and time per task, blue highlighted row represents subroutine consuming the most memory and time in total.

1
2
3
4
5
6
7

| subroutine | tasks | $\bar{x}$ seconds per task | total seconds | $\bar{x}$ RAM per task (MiB) | total RAM (MiB) |
|---|---|---|---|---|---|
| Clean fastas | 6 | 1.40 | 8.38 | 46.5 | 278.9 |
| New headers | 6 | 1.65 | 9.90 | 43.6 | 261.5 |
| Long isoform | 6 | 0.512 | 3.07 | 51.5 | 309.13 |
| **Transdecoder** | **1** | **10,436.7** | **10,436.7** | **3,249.8** | **3,249.8** |
| Diamond | 1 | 32.1 | 32.1 | 234.0 | 234.0 |
| SiLiX | 1 | 4.51 | 4.51 | 22.8 | 22.8 |
| Mafft | 135 | 3.24 | 437.8 | 18.3 | 2,466.5 |
| FastTree | 135 | 3.09 | 417.4 | 18.5 | 2,491.3 |
| TrimAL | 135 | 1.87 | 252.2 | 17.9 | 2,415.6 |
| **FUBAR** | **135** | **278.6** | **37,605.5** | **28.8** | **3,886.2** |

8