

**Development of fusion and duplication finder BLAST (fdfBLAST):
a systematic tool to detect differentially distributed gene fusions
and resolve trifurcations in the tree of life**

Submitted by **Guy Leonard**

to the **University of Exeter** as a thesis for the degree of **Doctor of Philosophy in**

Biological Sciences in **December 2010**

This thesis is available for Library use on the understanding that it is copyright material and that no quotation from the thesis may be published without proper acknowledgement.

I certify that all material in this thesis, which is not, my own work has been identified and that no material has previously been submitted and approved for the award of a degree by this or any other University.

Signature:

Guy Leonard 20/07/2011

Abstract

The construction of a tree of life and the placing of taxa into their correct phylogenetic context is the underpinning of modern evolutionary biology. However, many parts of the tree are currently unresolved due to conflicts within the sequence data. These sources of conflict include: horizontal gene transfer (HGT), hidden paralogy, and the effects of methodological artefacts such as Long Branch attraction (LBA). These limitations are further compounded by absence of key taxa that are yet to be sampled. Therefore, whilst phylogenetic methods are fundamentally useful for the reconstruction of the tree of life, given their current limitations, additional strategies are needed in order to fully resolve the tree of life. Gene fusions represent a potential source of evolutionary synapomorphies useful for resolving contentious branching relationships in the tree of life. I therefore, built a program to analyse whole genome datasets for the presence of differentially distributed gene fusion events (shared derived characters - SDCs). These putative SDCs can then be polarised with the help of traditional phylogenetic techniques and used as synapomorphies on the tree of life. Having constructed this program and tested it on established fusion datasets, I analysed five sets of four genomes from across the tree of life (the Deuterostomia, Fungi, Vertebrata, Viridiplantae and Discicristata). I used this data to identify the relative rates of gene fusion events. Previous studies have suggested that fission events occurred more often than gene fusion events. However, our analysis broadly suggests the opposite (albeit with a higher rate of fissions in the Deuterostomia). This result has direct implications for the use of gene fusions as evolutionary informative synapomorphies because the identification of a lower rate of reversion suggests that

these characters are less likely to be homoplasious and therefore represent useful tools for polarising evolutionary relationships. Six phylogenetically informative synapomorphies were recovered, three in the Discicristata which resolve the monophyly of the Kinetoplastida and four in the Fungi, one of which represented a HGT event and was independently discovered and previously published. Thus, this thesis reports the development and testing of a new tool to identify differentially distributed gene fusion events. The datasets analysed demonstrate that the program can be used to find phylogenetically informative gene fusion characters that can help resolve the tree of life in conjunction with traditional phylogenetic methods.

Table of Contents

| | |
|---|----|
| Acknowledgements | 2 |
| Abstract | 3 |
| Table of Contents | 5 |
| Table of Figures | 10 |
| 1 Introduction..... | 22 |
| 1.1 Early Interpretations of the Tree of Life | 22 |
| 1.1.1 Single-Gene Phylogenies | 24 |
| 1.1.2 Multi-Gene Phylogenies | 25 |
| 1.1.3 Super-matrices (Concatenation) | 27 |
| 1.1.4 Super-trees | 30 |
| 1.2 Molecular Sequence Based Phylogeny..... | 31 |
| 1.2.1 Sequence Data..... | 32 |
| 1.3 Phylogenomics..... | 40 |
| 1.3.1 Recent Interpretations of the Tree of Life..... | 40 |
| 1.3.2 Rooting Three Branched Trees | 45 |
| 1.4 Synapomorphies or Shared Derived Characters..... | 48 |
| 1.4.1 Gene Fusion/Fission Events..... | 50 |
| 1.4.2 Differential Relative Rates of Fusion and Fission across the Tree of Life. 54 | |
| 1.4.2.1 Example of Gene Fusion: Viridiplantae | 54 |
| 1.4.2.2 Example of Gene Fusion: Fungi | 55 |
| 1.4.3 Conserved Functional Domains..... | 56 |
| 1.5 Current Gene Fusion Detection Methods..... | 57 |

| | | |
|---------|--|----|
| 1.6 | Aims of this Thesis | 59 |
| 2 | Methods | 60 |
| 2.1 | Homologous BLAST Searches..... | 60 |
| 2.2 | REFGEN and TREENAMER | 62 |
| 2.3 | Alignment and Manual Masking..... | 64 |
| 2.4 | Substitution Model Prediction..... | 65 |
| 2.4.1 | Model Parameters | 66 |
| 2.5 | Tree Construction Methods..... | 67 |
| 2.5.1 | Maximum Likelihood Analysis | 67 |
| 2.5.1.1 | PHYML | 68 |
| 2.5.1.2 | RAxML..... | 69 |
| 2.5.2 | Bayesian Analysis..... | 70 |
| 2.5.2.1 | MrBayes..... | 71 |
| 2.5.3 | Bootstrapping | 72 |
| 2.5.4 | Approximate Likelihood-Ratio Tests (aLRT) | 73 |
| 2.6 | Drawing Trees | 73 |
| 2.6.1 | Automatic Tree Construction Pipeline | 74 |
| 2.6.2 | Tree Files..... | 75 |
| 3 | Fusion and Duplication Finder BLAST (fdfBLAST) – a tool to predict differentially distributed putative fusion and duplication events between proteomes. | 77 |
| 3.1 | Introduction | 77 |
| 3.2 | Aims | 79 |
| 3.3 | Materials and Methods..... | 80 |
| 3.3.1 | Step 1: Automated Serial BLASTp Comparisons..... | 82 |

| | | |
|---------|--|-----|
| 3.3.2 | Step 2: Comparative Hit Counts and Identification of Differential Distribution of Hit Numbers | 84 |
| 3.3.3 | Step 3: Reciprocal Hits | 88 |
| 3.3.4 | Step 4: Rank and Sort | 90 |
| 3.3.4.1 | Sorting of Data | 90 |
| 3.3.4.2 | Ranking of Data | 91 |
| 3.3.4.3 | Graphical Representation of Ranked and Sorted Data | 93 |
| 3.3.5 | Step 5: PFAM and CDD – Identification of Discrete Functional Domains to Confirm Fusions..... | 94 |
| 3.4 | fdfBLAST Program Overview..... | 96 |
| 3.5 | Discussion | 100 |
| 4 | Field-Testing the fdfBLAST Program with the Nakamura et al. (2006) Dataset ... | 102 |
| 4.1 | Introduction | 102 |
| 4.1.1 | A Re-evaluated Nakamura et al. (2006) Dataset..... | 103 |
| 4.1.2 | Arabidopsis thaliana-composite genes and Oryza sativa-split genes representing candidate gene fusions..... | 104 |
| 4.1.3 | Oryza sativa-composite genes and Arabidopsis thaliana-split genes representing candidate gene fusions..... | 108 |
| 4.2 | Results..... | 110 |
| 4.2.1 | fdfBLAST and a Two Plant Genome Dataset | 110 |
| 4.2.2 | fdfBLAST Results Vs the Re-evaluated Nakamura et al. (2006) Dataset | 111 |
| 4.3 | Discussion | 115 |
| 5 | Inferring the phylogeny of the kinetoplastids: a comparative genomics approach using whole genome datasets with low taxon sampling | 119 |

| | | |
|---------|--|-----|
| 5.1 | Introduction | 119 |
| 5.2 | Methods..... | 122 |
| 5.3 | Results..... | 125 |
| 5.3.1 | Multi-gene Phylogenetic Analysis of the Kinetoplastids | 125 |
| 5.3.2 | Paralogue Mirror-Tree Analysis..... | 129 |
| 5.3.3 | Phylogeny with Increased Taxa and Reduced Gene Sampling..... | 130 |
| 5.3.4 | Serial-Stripping of Fast Evolving Sites..... | 132 |
| 5.3.5 | Polarised Kinetoplastid Phylogeny and Gene Gain and Loss | 133 |
| 5.4 | Discussion | 134 |
| 6 | Four-way Genome Analyses using fdfBLAST on Five Calibration Datasets from across the Tree of Life | 137 |
| 6.1 | Introduction | 137 |
| 6.1.1 | Four-way Genome Dataset Selection for Calibration Purposes..... | 137 |
| 6.1.1.1 | Extended Viridiplantae Dataset | 140 |
| 6.1.1.2 | Fungi | 140 |
| 6.1.1.3 | Vertebrata | 141 |
| 6.1.1.4 | Discicristata | 141 |
| 6.1.1.5 | Deuterostomia | 142 |
| 6.2 | Methods..... | 143 |
| 6.2.1 | fdfBLAST Comparisons..... | 143 |
| 6.2.2 | Phylogenetically Informative Putative Shared Derived Characters | 143 |
| 6.3 | Results..... | 144 |
| 6.3.1 | Extended Viridiplantae Dataset fdfBLAST Analysis | 145 |
| 6.3.2 | Fungi | 147 |

| | | |
|---------|---|-----|
| 6.3.2.1 | Fungi: Phylogenetically Informative Datasets..... | 149 |
| 6.3.3 | Vertebrata..... | 158 |
| 6.3.4 | Discicristata..... | 159 |
| 6.3.4.1 | Phylogenetically Informative Datasets | 161 |
| 6.3.5 | Deuterostomia..... | 168 |
| 6.4 | Discussion | 169 |
| 6.4.1 | Comparisons Between the 4-way Datasets..... | 169 |
| 6.4.2 | Comparative Rates of Fusion and Fission..... | 171 |
| 6.5 | Conclusion..... | 174 |
| 7 | Discussion..... | 177 |
| 7.1 | Future Directions | 180 |
| 7.1.1 | fdfBLAST Applications..... | 180 |
| 7.1.2 | fdfBLAST Program Design | 181 |
| 8 | Appendix: A List of the 795 Taxa included in ‘Darren’s Orchard’ Automatic Phylogeny Pipeline | 184 |
| 9 | Appendix: fdfBLAST Perl Code Listing | 192 |
| 10 | Appendix: Putative Gene Fusions and Fissions as Predicted by fdfBLAST | 233 |
| 10.1 | Extended Viridiplantae | 233 |
| 10.2 | Fungi..... | 235 |
| 10.3 | Deuterostomia | 239 |
| 10.4 | Vertebrata..... | 243 |
| 10.5 | Discicristata..... | 247 |
| 11 | Phylogenetically Informative Tree Topologies | 249 |
| | Bibliography..... | 253 |