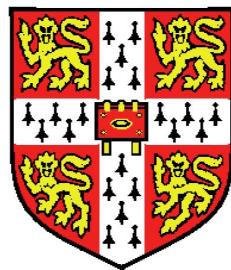


Sitewise error and constraint in mammalian comparative genomics



Gregory Jordan
European Bioinformatics Institute
University of Cambridge

A thesis submitted for the degree of
Doctor of Philosophy
September 29, 2011

Contents

| | |
|---|----------|
| Contents | 1 |
| 1 Introduction | 2 |
| 1.1 Introduction | 2 |
| 1.2 Biological background | 2 |
| 1.2.1 The evolutionary history of vertebrates and mammals | 2 |
| 1.2.2 Mammalian population structure, adaptation and natural selection | 2 |
| 1.3 Mathematical methods for genomic analysis | 3 |
| 1.3.1 Codon models of evolution | 3 |
| 1.3.2 Hypothesis testing with likelihood ratio tests | 3 |
| 1.3.3 Detecting purifying and positive selection in protein-coding sequence | 3 |
| 1.3.3.1 The Sitewise Likelihood Ratio test | 3 |
| 1.3.4 Identifying biological trends in sets of genes and proteins | 4 |
| 1.3.5 Correcting for multiple testing in genome-scale datasets | 4 |
| 2 The effects of alignment error and alignment filtering on the sitewise detection of positive selection | 5 |
| 2.1 Introduction | 6 |
| 2.1.1 Methods for detecting sitewise positive selection | 6 |
| 2.1.2 Substitution and indel processes in simulating protein-coding sequence evolution | 6 |
| 2.2 Models and parameters for simulating the evolution of mammalian genes . | 6 |
| 2.2.1 Distribution of selective pressures | 6 |
| 2.2.2 Phylogenetic tree size and shape | 6 |
| 2.2.3 Frequency and size distribution of insertions and deletions | 6 |
| 2.3 Analysis of the alignment error simulation results | 6 |

CONTENTS

| | | |
|----------|--|-----------|
| 2.4 | Methods for filtering alignments | 6 |
| 2.5 | Analysis of the alignment filtering simulation results | 6 |
| 3 | The effects of alignment error and alignment filtering on detecting positive selection in genes | 7 |
| 3.1 | Introduction | 7 |
| 3.1.1 | Existing methods for detecting gene-wide positive selection | 7 |
| 3.2 | The application of sitewise estimates to the genewise detection of positive selection | 7 |
| 3.3 | Analysis of the genewise detection results | 7 |
| 3.4 | Conclusions and further work | 7 |
| 4 | Curating a set of genome-wide orthologous gene trees in mammals | 8 |
| 4.1 | Introduction | 8 |
| 4.2 | Methods for ortholog identification | 8 |
| 4.3 | Low-coverage genomes in the Ensembl database | 10 |
| 4.4 | The Ensembl Compara gene tree pipeline | 12 |
| 4.5 | Quantifying paralogous relationships within Ensembl gene trees | 15 |
| 4.6 | Using taxonomic coverage to extract largely orthologous mammalian subtrees | 18 |
| 4.7 | Analysis of the set of root Compara gene trees | 23 |
| 4.8 | Analysis of sets of subtrees defined by taxonomic coverage and orthology annotation | 29 |
| 4.9 | Gene duplication and loss in the set of Eutherian largely orthologous trees | 34 |
| 4.10 | Comparison to gene trees from the Optic database of amniote orthologs . . | 39 |
| 4.11 | Conclusions | 39 |
| 5 | Patterns of sitewise selection in mammalian genomes | 43 |
| 5.1 | Introduction | 43 |
| 5.1.1 | The Mammalian Genome Project | 43 |
| 5.1.2 | The Sitewise Likelihood Ratio test | 44 |
| 5.1.3 | Data quality concerns: alignment and sequencing error | 45 |
| 5.2 | Preparing mammalian alignments for sitewise analysis | 49 |
| 5.2.1 | Filtering out low-quality genome sequence | 50 |
| 5.2.2 | Removing recent paralogs | 51 |
| 5.2.3 | Realigning coding sequences | 54 |

CONTENTS

| | | |
|----------|---|-----------|
| 5.2.4 | Filtering out clusters of non-synonymous substitutions | 54 |
| 5.3 | Genome-wide analysis of sitewise selective pressures in mammals | 55 |
| 5.3.1 | Mammalian species subsets for sitewise analysis | 55 |
| 5.3.2 | Evaluation of the bulk distributions and the design of a filtering approach | 57 |
| 5.3.3 | The global distribution of sitewise selective pressures in mammals . | 65 |
| 5.3.3.1 | Site patterns and ω_{ML} values reveal the prevalence of purifying selection in mammalian proteins | 65 |
| 5.3.3.2 | Sitewise confidence intervals and LRT statistics identify sites with significant evidence for purifying and positive selection | 67 |
| 5.3.3.3 | Estimating the proportion of negatively-selected sites . . . | 70 |
| 5.3.3.4 | Estimating the proportion of positively-selected sites . . . | 70 |
| 5.3.3.5 | Correlations between branch length, effective population size and sitewise summary statistics | 71 |
| 5.3.4 | Modeling the global distribution of sitewise selective pressures . . . | 71 |
| 5.3.5 | Simulations to evaluate the fit of empirical data to a simulated poisson process at varying branch lengths and ω ratios | 74 |
| 5.3.6 | Simulations to evaluate the power to detect positive selection and estimate selective pressures | 77 |
| 5.3.7 | Evaluation of the effect of GC content, recombination rate, and codon usage on sitewise dN/dS estimates and the detection of positive selection | 78 |
| 5.4 | Conclusions | 81 |
| 6 | The use of sitewise selective pressures to characterise the evolution of genes and domains in mammals | 82 |
| 6.1 | Introduction | 82 |
| 6.2 | Comparison of sitewise results to previously described sets of positively selected genes | 83 |
| 6.3 | Using sitewise selective pressures to characterise the evolution of genes . | 83 |
| 6.3.1 | Identifying genes subject to positive selection | 83 |
| 6.3.2 | Identifying genes subject to strong or weak purifying selection . . | 83 |
| 6.4 | Using sitewise selective pressures to characterise the evolution of protein domains | 83 |

CONTENTS

| | | |
|----------|--|-----------|
| 6.4.1 | Identifying protein domains subject to positive selection | 84 |
| 6.4.2 | Identifying protein domains subject to strong or weak purifying selection | 84 |
| 6.5 | Identifying genes under unusual selective pressures in mammalian superorders | 84 |
| 7 | Evolution of protein-coding genes in gorilla and the African apes | 85 |
| 7.1 | Introduction | 86 |
| 7.1.1 | The gorilla and other primate genome projects | 86 |
| 7.1.2 | Incomplete lineage sorting | 86 |
| 7.1.3 | Effective population sizes of extant and ancestral primate populations | 86 |
| 7.1.4 | Measuring shifts in selective pressures using branch-specific likelihood ratio tests | 86 |
| 7.1.5 | Data quality concerns: sequencing, assembly and alignment error . | 86 |
| 7.2 | Constructing codon alignments of one-to-one orthologous genes in six primate species | 86 |
| 7.2.1 | Identification of genes with one-to-one homology | 86 |
| 7.2.2 | Collection of homologous DNA sequences from genome- or transcript-based multiple alignments | 86 |
| 7.2.3 | Filtering sequence regions with low sequence quality | 86 |
| 7.2.4 | Filtering sequence regions with high substitution counts | 86 |
| 7.2.5 | Filtering sequence regions with evidence of incomplete lineage sorting | 86 |
| 7.3 | Analysis of patterns of duplication and deletion in primate gene families . | 86 |
| 7.4 | Analysis of the likelihood ratio test results | 86 |
| 7.4.1 | Genes with evidence for acceleration and deceleration in the human, chimpanzee and gorilla terminal lineages | 86 |
| 7.4.2 | Genes with evidence for acceleration in the African great ape ancestral branch | 86 |
| 7.4.3 | Genes with evidence for positive selection based on the branch-site test | 86 |
| 8 | Gorilla part 2 | 87 |
| 8.1 | Analysis of incomplete lineage sorting in the African great apes within and nearby protein-coding genes | 87 |
| 8.2 | Analysis of dN/dS levels in six primate genomes | 87 |
| 8.2.1 | Genome-wide dN/dS in six primates and their ancestors | 87 |

CONTENTS

| | | |
|-------|--|-----------|
| 8.2.2 | Genome-wide dN/dS in regions of differing sitewise constraint | 87 |
| 8.2.3 | Analysis of the impact of sequence and alignment filtering on primate dN/dS estimates | 87 |
| 8.3 | Conclusions and future work | 87 |
| | Bibliography | 88 |

CONTENTS

Todo list

| | |
|--|----|
| ■ Write an introductory page or two. Keep it broad, cite some old stuff. | 2 |
| ■ n=XYZ | 53 |
| ■ n=XYZ | 53 |
| ■ 0.983 ± 0.005 | 53 |
| ■ Use the sites data and 'seq' table to do a more quantitative analysis here? | 60 |
| ■ Refer to Capra and Pollard 2011, who showed that primates and carnivores have lots of BDS (biased divergent substitutions) while glires don't! | 71 |

Chapter 1

Introduction

1.1 Introduction

Write an i
tory page
Keep it br
some old s

1.2 Biological background

[All of evolutionary biology is about understanding the history of evolution, and is thus tied to the circumstances under which such evolution occurred. Thus, a brief overview of that history is relevant and useful in this thesis.]

1.2.1 The evolutionary history of vertebrates and mammals

[Write a quick summary of the genome evolution of vertebrates and mammals. Mention 2R duplication, genome size growth, transposable elements.]

1.2.2 Mammalian population structure, adaptation and natural selection

[Introduce the concept of adaptation (molecular vs. morphological / ecological), the varied behavioral characteristics of modern day mammals (focusing on mammalian superorders and great apes, as expanded in mammals and gorilla chapters), and the impact of population structure / population size on the efficacy of natural selection.]

1.3 Mathematical methods for genomic analysis

[Introduce the importance of imperfect replication / copying as the substrate of evolution and a very convenient phenomenon for mathematical analysis. The balance between randomness and structure in evolutionary models.]

1.3.1 Codon models of evolution

[Introduce the idea of modeling protein evolution as a markov process acting on codon sequences: the incorporation of mechanistic parameters for Ts:Tv bias (κ), dN/dS ratio (ω), or empirical models a la . Talk about heterogeneity the idea that real data may strongly violate certain models.]

1.3.2 Hypothesis testing with likelihood ratio tests

[Briefly introduce the idea of nested models and likelihood ratio tests (used for PAML in Slrsim and Gorilla chapters, and for SLR in Mammals chapters)]

1.3.3 Detecting purifying and positive selection in protein-coding sequence

[Briefly run through the history of detecting purifying / positive selection in genes and sites. Mention history of PAML models, alternative approaches, and fully describe SLR's approach.]

1.3.3.1 The Sitewise Likelihood Ratio test

SLR implements a method specifically designed for sitewise estimates which has been shown in simulations to perform as well as or better than PAMLs sitewise random sites models (Massingham and Goldman, 2005). SLR models codon evolution as a continuous-time Markov process where substitutions at one site are independent of substitutions at all other sites. No assumptions are made regarding the distribution of ratios within the alignment. The value of ω is considered to be an independent parameter at each site: after first optimizing shared parameters using the whole alignment, SLR uses the shared parameters and the data at each alignment site to calculate a sitewise statistic for non-neutral evolution. This statistic is based on a likelihood-ratio test where the null model is

neutral evolution ($\pi = 1$) and the alternative model is either purifying or positive selection ($\pi < 1$ or $\pi > 1$, respectively). The raw statistic measures the strength of evidence for non-neutral evolution at each site; following Massingham and Goldman (2005) we use a signed version of the SLR statistic (created by negating the statistic for sites with $\pi < 1$) as the test statistic for positive selection.

1.3.4 Identifying biological trends in sets of genes and proteins

[Introduce the Gene Ontology and Pfam databases, which annotate genes or components of genes with structured ontologies of functions or domains, respectively. Introduce the methods for detecting enriched GO terms. Note problems and biases involved in the basic methodology and describe algorithms / corrections introduced to correct for certain biases: topGO for hierarchical GO structure and goseq for element length.]

1.3.5 Correcting for multiple testing in genome-scale datasets

[Note the issue of correcting for multiple testing in genome-scale datasets. Clarify the differences between nominal p-value, family-wise error rate, FDR. Provide examples of when / where certain methods may be more applicable than others.]

Chapter 2

The effects of alignment error and alignment filtering on the sitewise detection of positive selection

2.1 Introduction

- 2.1.1 Methods for detecting sitewise positive selection
- 2.1.2 Substitution and indel processes in simulating protein-coding sequence evolution

2.2 Models and parameters for simulating the evolution of mammalian genes

- 2.2.1 Distribution of selective pressures
- 2.2.2 Phylogenetic tree size and shape
- 2.2.3 Frequency and size distribution of insertions and deletions

2.3 Analysis of the alignment error simulation results

2.4 Methods for filtering alignments

2.5 Analysis of the alignment filtering simulation results

Chapter 3

The effects of alignment error and alignment filtering on detecting positive selection in genes

3.1 Introduction

3.1.1 Existing methods for detecting gene-wide positive selection

3.2 The application of sitewise estimates to the ge- newise detection of positive selection

3.3 Analysis of the genewise detection results

3.4 Conclusions and further work

Chapter 4

Curating a set of genome-wide orthologous gene trees in mammals

4.1 Introduction

This chapter considers the identification of groups of orthologous genes in vertebrates and mammals. I begin with a discussion of the detection of orthologs within groups of vertebrate and other eukaryotic genomes, focusing on the tree-based approach to ortholog and paralog identification as used by the Optic, Compara, and Treefam databases. In preparation for the sitewise analysis of mammalian orthologs to be presented in the subsequent two chapters, I construct a simple method for isolating subtrees with largely orthologous characteristics from the set of Compara gene trees and analyze the taxonomic distribution of gene duplication and loss within those largely orthologous trees (LOTs), demonstrating along the way how strongly the history of genome duplications has affected the structure of orthologous and paralogous relationships between and within vertebrate and invertebrate genomes.

4.2 Methods for ortholog identification

Central to any evolutionary sequence analysis is the assumption that the sequences being analyzed, or some parts therein, share a common evolutionary origin. Thus, the first step in any such analysis is the collection of homologous sequence data. Starting with a source sequence, homologs are usually identified by searching a sequence database for sequences with a minimum overall similarity according to some evolutionary model. In some cases,

such as when sequences are highly divergent or subject to domain shuffling or horizontal gene transfer, a more localized or more specialized measure of similarity may be desirable [Koonin *et al.*, 2001; Sjölander *et al.*, 2011], but within closely related groups of organisms such as vertebrates and mammals, the process of gene duplication and divergence dominates patterns of relatedness between protein sequences [Ohno, 1970], making overall sequence similarity a useful method by which to identify homologous sequences in the organisms of interest.

Heuristic algorithms have been developed for performing quick and sensitive sequence homology searches within databases of protein and nucleic acid sequences [Altschul *et al.*, 1997; Eddy, 2009]. The power of these methods is high enough that homology within vertebrates, even for fast-evolving genes, can be readily detected. However, the prevalence of historical and ongoing gene duplication and loss in vertebrates complicates the problem, as orthologous genes (e.g., homologs between species, sharing a common ancestral population and related through a speciation event) and paralogous genes (e.g., homologs within a single species, sharing a common ancestral genomic sequence and related through a duplication event) are important, yet difficult, to distinguish from one another [Jun *et al.*, 2009]. In other words, sets of genes that may be confidently identified as sharing homology may still contain paralogy and orthology relationships that are more difficult to resolve. Since the evolutionary trajectories of paralogous versus orthologous genes are expected to be quite distinct [Lynch & Conery, 2000], the correct identification of orthologous versus paralogous relationships in vertebrates is critical for any detailed molecular evolutionary analysis.

Methods for distinguishing orthologous from paralogous relationships within vertebrate genomes have been in development for over a decade [Remm *et al.*, 2001; Yuan *et al.*, 1998], but the amount of overlap between vertebrate orthologous groups identified by different methods has historically been disappointingly small [Chen *et al.*, 2007; Jun *et al.*, 2009], suggesting that the problem of orthology inference is far from resolved. Still, one might expect the accuracy and power of orthology inference methods to improve with time, given the steady increase in available computing power and the sequencing of many complete vertebrate genomes over the past decade. Indeed, the growing number of available sequenced genomes, greater available computational power, and improved understanding of patterns of gene duplication and loss have led to the growing popularity of phylogeny-based approaches, which were once considered computationally impractical and too difficult to automate [Remm *et al.*, 2001]. In general, the phylogenetic approach involves estimating

a phylogenetic tree from an entire cluster of homologous genes and inferring duplication events based on the discordance between the gene tree and the species tree. Several variants of this approach have been developed and applied to gene tree reconstruction in insects, fungi and vertebrates [Cepas *et al.*, 2007; Datta *et al.*, 2009; Hahn *et al.*, 2007; Muller *et al.*, 2010; Ruan *et al.*, 2008; Vilella *et al.*, 2009], and validation against manually-curated gene trees has shown these phylogeny-based methods to be more sensitive and accurate than pairwise or graph-based methods [Datta *et al.*, 2009].

One particular phylogeny-based approach, which uses a known species tree to guide the resolution of orthologs and paralogs within a gene tree in a Bayesian framework, has been implemented in a software package called TreeBeST and used extensively to infer genome-wide vertebrate and eukaryotic gene trees for the Optic, TreeFam and Compara databases [Heger & Ponting, 2008; Ruan *et al.*, 2008; Vilella *et al.*, 2009].

I describe in this chapter an initial analysis of the gene trees within Compara which I undertook in preparation for the sitewise evolutionary analysis presented in Chapters 5 and ?. The major goals of this analysis were to better understand the distribution of orthologous genes within vertebrate genomes, identify which taxonomic constraints would best define largely-orthologous groups of mammalian genes, and evaluate whether low-coverage genomes, which are unique to the Compara database of orthologs, showed high enough annotation quality and consistency to be included in the mammalian evolutionary analysis. Throughout the chapter, attention will be paid to those aspects of the gene tree identification pipeline which may contribute to errors in the downstream evolutionary analysis.

4.3 Low-coverage genomes in the Ensembl database

A major feature that distinguishes Compara from the Treefam and Optic databases is the inclusion in Compara of several mammalian genomes which have been sequenced to low coverage as a part of the Mammalian Genome Project (MGP). The inclusion of these additional genomes was a major reason why I performed the analysis described in the following chapters using Ensembl as a source of gene trees and alignments, representing a major advantage over otherwise similar orthology databases due to the greater species sampling density. With respect to the actual sequencing of the low-coverage mammalian genomes, the motivation and goals of the MGP itself will be outlined in the next chapter. For the purpose of the present discussion, I will outline the way low-coverage genomes are

handled within the Ensembl annotation pipeline, as certain features of the low-coverage genomes cause characteristic anomalies in the orthology analyses, as will be shown in the ensuing analysis.

The prevalence of missing sequence data and fragmented contigs in low-coverage genomes presented a unique set of problems for the generation of transcript annotations in Ensembl. In recognition of these differences, the procedure used by Compara to annotate genomes assembled from low-coverage data is distinct from the usual gene-building pipeline [Hubbard *et al.*, 2007]. Briefly, a whole-genome alignment is produced between the human genome and each low-coverage target genome, and gene models are projected from human to the target genome. Small frame-disrupting insertions or deletions within orthologous exons are corrected, and missing or incomplete exons are padded with Ns in order to produce a transcript with a length equal to that of the human reference transcript. The inclusion of these error-correcting features allows for a set of intact, if not complete, coding transcripts to be generated for low-coverage genomes, leveraging the high level of sequence similarity between human and other Eutherian mammals to project genes and transcripts from the high-quality human genome to the unannotated, highly fragmented low-coverage genome assemblies. Still, in many cases the Ensembl pipeline could not map complete genes or transcripts from human to the target genome, causing difficulty in identifying duplications. On one hand, the lack of completely assembled chromosomes meant that segmental duplications in low-coverage genomes were often unresolved or unidentifiable, making it difficult or impossible to confidently identify recently duplicated genes. On the other hand, the shorter length of assembled fragments caused genes to occasionally be split between two sequence fragments; the Ensembl pipeline currently annotates such genes as two separate shortened gene fragments, resulting in an excess of shortened apparent paralogs in resulting gene trees (see Section ?? for more detail on this artifact).

The Compara gene family pipeline, which I describe in more detail below, uses the set of annotated transcripts from each species as its input [Vilella *et al.*, 2009], so the quality of gene annotation from each source genome has a direct impact on the overall quality and accuracy of the resulting gene trees. Although the reliance on genome-wide alignments to, and gene annotations from, a reference genome could be criticised for potentially causing a bias towards the genomic properties of the reference, this approach is a reasonable workaround in the absence of higher-coverage sequence data or a painstakingly curated assembly. Furthermore, the gene model error-correcting features of the Ensembl pipeline are especially beneficial, making more complicated methods for correcting errors from low-

coverage genomes such as those described by [Hubisz *et al.*, 2011] seem less necessary.

4.4 The Ensembl Compara gene tree pipeline

All genomic data and gene trees used for this analysis were sourced from version 63 of the Ensembl Compara database [Flicek *et al.*, 2011; Vilella *et al.*, 2009]. Although a complete description of the design, implementation, and validation of the pipeline behind the Ensembl database would be beyond the scope of this chapter, I will briefly outline the major aspects of the approach used by the Compara pipeline, focusing on a few details which are relevant to the current analysis and ensuing discussion.

The Compara pipeline begins with a set of protein-coding transcripts collected from each individual species' annotation database. This step is not exactly straightforward, however, as the prevalence of alternative splicing in Eutherian mammals makes it common for a single gene to harbor many different transcript structures.

In terms of biology and evolution, alternative splicing is a very interesting phenomenon. Tightly linked to the evolutionary innovation of regulatory control and tissue-specific gene expression, the existence of multiple transcripts per gene is one of the likely substrates of biological and developmental complexity within vertebrates and mammals as compared to single-celled eukaryotes, which show less developmental complexity but largely similar numbers of genes [Csuros *et al.*, 2011]. Further evidence of the unique evolutionary characteristics of alternatively-spliced exons comes from molecular evolutionary studies which have shown such exons to show, on average, higher levels of evolutionary constraint, possibly owing to the importance of exonic splice enhancers in modulating the inclusion or exclusion of their associated exons [Parmley *et al.*, 2006].

In terms of organizing biological data, however, pervasive alternative splicing is somewhat burdensome. Roughly 34% of human genes contain at least two, and up to several dozen, transcripts per gene [Mironov *et al.*, 1999], showing tissue-specific and species-specific expression patterns, different levels of overall transcription, and sometimes comprising mutually exclusive exons. Within the context of the Compara database, the first problem in maintaining consistency across many source genome sequences is the fact that primary data on alternative transcript structures (e.g., data from expressed sequence tags, RNA-seq, or proteomics experiments) are largely absent from most organisms with sequenced genomes. Furthermore, the task of incorporating multiple transcripts per gene into an evolutionary framework is non-trivial, and leaves many unresolved questions open

to debate: should all transcripts be treated as independent evolutionary entities, or should some form of meta-transcript be produced, comprising all possible transcripts for a given gene? Should expression levels and tissue-specificity be taken into account (as both factors have been correlated with evolutionary rate, e.g. [Koonin & Wolf, 2006; Zhu *et al.*, 2008])? And what is the expected evolutionary impact of the loss, gain, or modulation of the prevalence or tissue-specificity of a given exon or transcript in one lineage? Even a fairly shallow consideration of the topic quickly reveals layers of complexity that would quickly hinder many large-scale evolutionary analyses, such as the definition of orthologous groups, whose main goals are to understand the evolutionary relationships of gene families within some set of species' genomes.

As a result of these difficulties, the current design of the Compara pipeline only incorporates one 'canonical' transcript per gene into the evolutionary analysis and the resulting inferred gene trees. This reflects a conscious decision to sacrifice some biological fidelity for reduced design complexity and computational load, as the inclusion of multiple transcripts would inevitably require some amount of additional processing and/or calculation. Unfortunately, this only somewhat alleviates the problem, shifting the burden from "how to deal with multiple transcripts in a comparative setting" to "how to choose the best representative transcript for each gene." In the case of a gene with many transcripts of varying sizes containing many non-overlapping exons, the negative consequences of choosing a non-optimal transcript are clear: too short of a transcript could exclude important sequence information from the dataset, while transcripts with spurious exons (resulting from misannotation or erroneous experimental evidence for a transcript) could introduce potentially large amounts of non-orthologous, nonfunctional, or nonconserved sequence into the evolutionary analysis.

Fortunately, the consensus coding sequence (CCDS) project was initiated in 2005 to "identify a core set of human and mouse protein coding regions that are consistently annotated and of high quality" [Pruitt *et al.*, 2009]. Although the transcripts that satisfy these two criteria will not necessarily be the same as those which meet the desired definition of "the best representative transcript for use in an evolutionary study," the confidence that one can have in the quality and consistency of CCDS transcripts helps to reduce the prevalence of potentially damaging errors in the Compara pipeline. Thus, in the current release (version 63), the "representative" transcript used for the Compara pipeline is chosen on the basis of (a) existence within the CCDS set of transcripts and (b) the total length of the transcript's coding sequence. The combination of these two factors can be expected to

identify a reasonably representative transcript, at least for the human and mouse genomes. The situation will be similar for genomes whose Ensembl annotation is derived largely from synteny and orthology to human and mouse annotated genes, but two classes of genomes—those resulting from low-coverage sequencing and those from more distant species whose annotations are derived from largely independent data sources—will still suffer from some amount error in the form of poor transcript choice. The OPTIC database, which contains orthologous groups identified within a wide range of animal and fungal clades, uses the transcript with the maximum length as the canonical transcript [Heger & Ponting, 2008].

Once the set of canonical transcripts is chosen, an all-against-all protein BLAST search is performed using the Washington University variant of BLAST and genes are clustered into homologous groups using *hcluster_sg*, an implementation of a hierarchical clustering algorithm for sparse graphs. Sequences are aligned using MCoffee, a meta-aligner algorithm which combines the results from different aligners into one alignment using a maximum-consistency criterion [Wallace *et al.*, 2006]. The aligners used for the M-Coffee alignment include MAFFT [Katoh *et al.*, 2005], MUSCLE [Edgar, 2004], KAlign [Lassmann *et al.*, 2009], and T-Coffee [Notredame *et al.*, 2000]. Finally, the aligned sequences are input to TreeBeST, which infers a gene tree (including gene duplication and loss events) given a set of aligned sequences and a known species tree [Ruan *et al.*, 2008]. The type of the homology relationship between each pair of genes (e.g., one-to-one ortholog, one-to-many ortholog, within-species paralog) is determined using a simple set of rules based on the structure of the inferred gene tree and the annotation of ancestral nodes where a duplication event has likely occurred.

The Compara pipeline has been a part of the Ensembl ecosystem since its first introduction to Ensembl in release 42 [Birney *et al.*, 2006]. Remarkably, aside from slight tweaks to the protein clustering method and some changes in the exact aligners used, the pipeline has changed little from its original published form [Vilella *et al.*, 2009]. In part, this lack of change reflects the ease with which sets of vertebrate orthologs can be identified using the existing methodology, lying in contrast to the equivalent task in sets of insect or fungal genomes where divergence levels between extant species with sequenced genomes are much larger [Siepel *et al.*, 2005] and the shape of the underlying species tree may be uncertain and/or unknown [MacKenzie *et al.*, 2008], making the development of specialized methods or extensive manual annotation necessary [Kellis *et al.*, 2004; Rasmussen & Kellis, 2007]. This is equivalent to saying that Ensembl’s pipeline, while not perfect in its orthology predictions or tree inferences (as indicated in a series of back-and-forth papers between

Milinkovitch et al. [2010] and Villela et al. [2011]), has proved sufficiently accurate enough that an extensive reworking of the system has not yet been deemed necessary. However, it is worth noting that the recent development of a new Bayesian method for gene tree reconciliation, based on a generative model of gene family evolution as a birth-death process and incorporating species-specific and gene-specific rate variation, showed good performance in resolving fungal and invertebrate gene trees, and could easily be adapted to the vertebrate genomes in the Compara database . Additional validation of the overall approach taken by Compara comes in the form of Treefam and Optic [Ruan *et al.*, 2008], databases of animal and fungal gene trees which applied a similar set of tools to infer gene trees from a more diverse set of genomes, producing qualitatively similar results [Vilella *et al.*, 2009].

4.5 Quantifying paralogous relationships within Ensembl gene trees

The first task in preparing the Compara data for subsequent sitewise analysis was to identify and extract a set of gene trees or gene subtrees comprising largely of orthologous relationships within mammals (i.e. LOTs), avoiding as much as possible the inclusion of paralogous relationships. It was necessary to postprocess the Compara gene trees in order to achieve this goal, because many of the Compara trees contained multiple sets of complete mammalian orthologous trees linked by ancient gene duplication events, while I wished to study the evolution of each mammalian LOT in isolation. In other words, the Compara gene trees could be characterized as being over-clustered with respect to the core set of mammalian orthologous trees. This over-clustering was not necessarily inaccurate with respect to the evolutionary history of vertebrate genes, but it was not a desirable feature for my intended use of the data. This section is concerned with identifying and quantifying the level of such over-clustering within the Compara database.

I first collected the set of 18,607 Compara gene trees and analyzed their overall size and the number of human genes contained within each tree. The results of this analysis are presented in Table 4.2 and Figure 4.3; a complete description of these data will follow in Section 4.7, but the aspects relevant to characterizing the amount of paralogy contained within these trees will be discussed here.

The first row of Table 4.2 shows various summary statistics from the full set of Compara gene trees, with the columns under the “Human Content” heading showing the fraction of all gene trees containing zero, one, or two or more human genes. Evidence for the existence

of large numbers of paralogs within Compara trees came from the observation that 20% of Compara trees contained 2 or more human genes. If each Compara tree contained only one set of mammalian orthologs, then the 20% of trees with multiple human gene copies could only be explained by an unrealistically high rate of gene duplication in the lineage leading towards human. Instead, a more parsimonious explanation was that many of the Ensembl trees contained not one group of mammalian LOTs, but two or more sets of mammalian LOTs joined by one or more ancient duplication events. This explanation was further supported by Figure 4.3, which shows a histogram of total gene counts in the Ensembl root trees. A large number of trees contained more than 48 sequences (the number of vertebrate genomes in Ensembl), with clear peaks in the histogram of trees with 2, 3, and 4 times the number of vertebrate genomes in Ensembl. These patterns were highly consistent with a high number of Compara gene trees containing multiple mammalian LOTs.

The prevalence of over-clustered Eutherian orthologs in the Compara database could be explained by a combination of the *hcluster_sg* algorithm used for the hierarchical clustering step, which uses only protein distances as its source of clustering information, and the wide range of protein evolutionary rates in the vertebrate genome. As I mentioned in Section 4.4, the Compara pipeline uses all-by-all protein BLAST E-value scores and the *hcluster_sg* algorithm to produce sets of sequences containing minimal average within-group E-values. No additional biological information, such as the source species of each sequence or the overall taxonomic coverage (TC) of each cluster, is used in identifying clusters, and no attempt is made to fit clusters to an expected model of orthologous gene evolution. On the one hand, the lack of additional information and assumptions allows the algorithm to remain simple and the clustering behavior to remain consistent across different groups of genomes; on the other hand, a number of technical (in the sense of non-biologically meaningful) parameters and thresholds must be tuned in order to result in the desired cluster sizes and contents. Importantly, even after these parameters were tuned to perform well on the dataset as a whole (i.e., to successfully cluster a majority of vertebrate orthologs), the reliance on protein distances alone means that fast-evolving proteins will be more likely to be under-clustered and slow-evolving proteins will be more likely to be over-clustered. The rate of protein evolution varies widely within a genome, as evidenced by a study of amino-acid substitution rates of roughly 6,000 eukaryotic orthologous genes by Koonin et al. [Koonin et al. 2004], who found that the middle 90% of genes show a nearly fourfold variation in evolutionary rate. Given this wide range of protein evolutionary

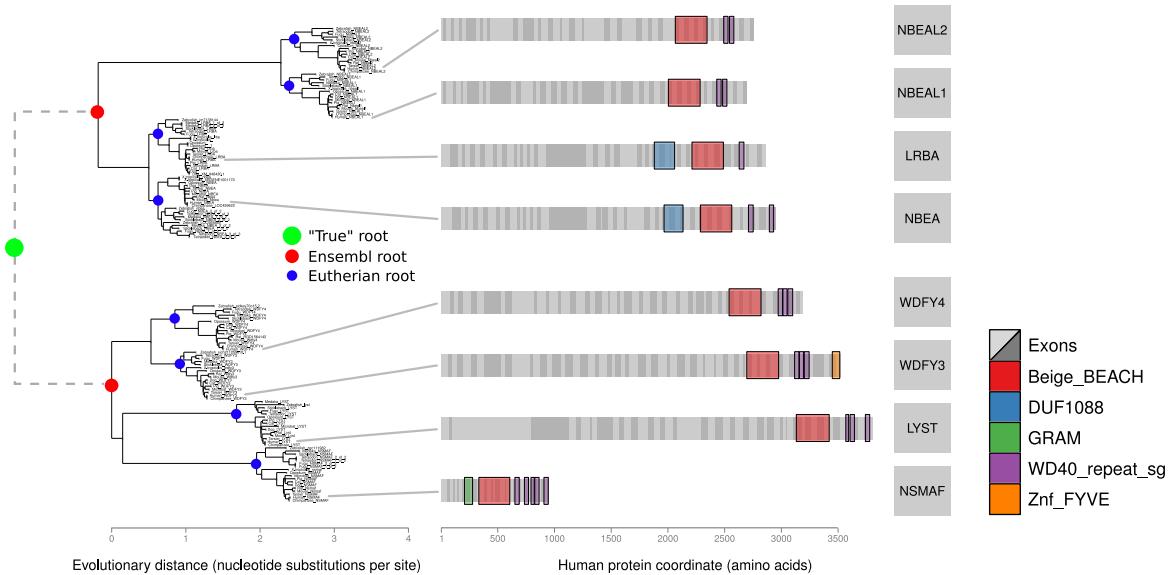


Figure 4.1: The evolutionary history of the human *neurobeachin-like 2* gene (*NBEAL2*) and its paralogs. Left, two phylogenetic trees from Ensembl Compara (release 60) are shown, summarizing the evolution of *NBEAL2* and its three paralogs (top) and *LYST*, a presumed distant paralog of *NBEAL2*, and its three paralogs (bottom) in 15 vertebrate species. The phylogeny shows that *NBEAL2* is taxonomically conserved and distinct from its paralogs. Red dots highlight the root nodes of Ensembl gene trees, blue dots highlight the root nodes of Eutherian orthologous subtrees, and a dashed line with a green dot represents the putative paralogous relationship (with a hypothetical root) between the two Ensembl gene trees. Right, the exon and domain structure of each human gene is shown: exons are displayed alternating shades of gray, and Pfam domain annotations are colored according to their Pfam identifier.

rates, the excess of over-clustered orthologs in the Compara database is understandable and even somewhat expected.

Before continuing with a description of a method to identify LOTs within Compara gene trees, I should note that my use of the phrase “over-clustered” refers only to over-clustering with respect to the current goal of analyzing independent sets of orthologous genes within mammals. Certainly these large “over-clustered” trees, which represent a more distant evolutionary history than a single mammalian orthologous group, are just as accurate with respect to the true evolutionary history of the genes as more narrow groupings would be. Furthermore, the inclusion of a deeper evolutionary context may sometimes be more useful to users of the Compara database, for whom an understanding of the overall evolutionary history of a gene may be the topic of primary interest.

As an example, take the gene *NBEAL2* and its human paralogs, whose gene trees, exon structures and domain classifications were extracted from Ensembl v62 and summa-

rized in Figure 4.1. A recent medical sequencing project identified *NBEAL2*, a gene of previously unknown function, as the putative causative gene for gray platelet syndrome, a predominantly recessive platelet disorder resulting in moderate to severe bleeding [Albers *et al.*, 2011]. I performed, with Botond Sipos, an evolutionary analysis of *NBEAL2* for the article describing the discovery, and it was important for the purpose of this study to ensure that the *NBEAL2* gene was both well-conserved across mammals and distinct from its paralogs. The Compara pipeline clustered *NBEAL2* with three of its closest paralogs into one tree (and similarly clustered four more distant *NBEAL2* paralogs into a separate tree), yielding two views which together showed both the full taxonomic coverage of the *NBEAL2* subtree and the large amount of evolutionary distance between paralogs. Had each mammalian ortholog been displayed independently in Ensembl (i.e., using the blue “Eutherian root” nodes in Figure 4.1), it would have been more difficult for a non-expert to make such claims regarding the evolutionary history of *NBEAL2* without further analysis. Conversely, had the Compara pipeline been even more inclusive in its clustering approach and identified the hypothetical deeper root connecting these two sets of trees (represented by the green node in Figure 4.1), the connection between these eight genes would have been even more immediately apparent.

For the purposes of this project, however, it was important to isolate individual mammalian gene trees for further processing and sitewise analysis. To this end, I designed a simple scheme for splitting gene trees into non-overlapping subtrees based on flexible TC criteria; the remainder of this chapter presents design and results of this process when applied to the Compara gene trees.

4.6 Using taxonomic coverage to extract largely orthologous mammalian subtrees

Based on the observation that the clustering step of the Compara pipeline did not make use of any taxonomic information, I hypothesized that a relatively simple set of rules based on TC would be sufficient to identify most mammalian LOTs. Two well-established observations in mammalian genomes supported the decision to use TC in this context. First, the existence of two rounds of whole-genome duplication preceding the evolution of vertebrates [Dehal & Boore, 2005] suggested that many of the ancient duplication events contained within Ensembl gene trees occurred before the divergence of mammals, making it possible to cleanly separate out taxonomically complete mammalian subtrees in the

majority of cases. This would not be possible if duplication events were common and spread evenly throughout the mammalian tree; if that were the case, many duplication events would have occurred after the divergence of some or all of the major mammalian groups, resulting in a larger proportion of mammalian genes with “internal” duplications and, thus, fewer singly orthologous trees with high taxonomic coverage. Second, the overall low rate of gene duplication and loss in mammals [Demuth *et al.*, 2006] (excluding, of course, the aforementioned whole-genome duplication events which occurred in the ancestral vertebrate genome) predicts that few mammalian gene trees will be subject to one or more gene duplication or loss events. In other words, most mammalian gene trees should contain sequences from a majority of mammalian species, so the effectiveness of using TC to identify mammalian subtrees should be largely unaffected by continued (i.e., post whole-genome duplication) gene duplication or loss events. The potential utility of TC was further bolstered by the star-like shape of the mammalian tree [Bininda-Emonds *et al.*, 2007]: a star-like tree contains more branch length within terminal lineages than a ladder-like tree with an equivalent total branch length, making it less likely that a gene duplication or loss event (if such events occurred randomly throughout the mammalian tree) would result in a significant disruption to the TC of the gene tree.

The TC-based tree splitting process worked as follows. For every internal node N of each Compara gene tree, the TC was calculated for several vertebrate clades. The TC for node N and clade C was calculated as $TC(N, C) = \text{species}(N)/\text{species}(C)$, where $\text{species}(N)$ is the number of unique species represented by the sequences beneath node N and $\text{species}(C)$ is the number of species within the vertebrate clade C . The tree was then traversed from root to tip, and if a given set of TC constraints (TCCs) were satisfied by both subtrees below node i , then the tree was split into two subtrees at node i , with the new trees having root nodes placed at the two child nodes, i_a and i_b . The traversal continued recursively until every node was tested against the TCCs. The smallest subtrees which satisfied the TCCs were included in the resulting subtree set. If only the entire gene tree satisfied the TCCs then the entire tree was included; if the entire gene tree failed to satisfy the TCCs, it was excluded altogether from the resulting subtree set.

I chose a variety of TCCs to apply to the set of Compara trees, all of which were run against the 18,607 gene trees within the Compara database to generate several genome-wide subtree sets. Table 4.1 shows the details of the various TCCs used. The clade names (e.g., $TC(\text{Primates})$) refer to the set of species in the Ensembl database that are contained within the given subtree of the NCBI taxonomy; the NCBI classification of

species contained within Ensembl is shown in Figure 4.2, including labels on internal nodes corresponding to the clade names given in Table 4.1.

For the Ingroup and Outgroup categories of TCCs, a TC value of greater than 0.6 was required for a single taxonomic clade. This value was not arbitrarily chosen; rather, it was important to use a TC value slightly above 0.5 to achieve the desired result of identifying orthologous subtrees. A value much higher would be too restrictive: if, for example, the required TC value were set to 1, then all subtrees containing a deletion in any species within the clade of interest would not satisfy the TCC. On the other hand, a required TC value of less than 0.5 would allow a single LOT to be split into two subtrees, with one subtree having $TC < 0.5$, and the other subtree, containing the other half of the species within the clade of interest, also having $TC < 0.5$. Thus, 0.6 was deemed a sufficient TC requirement for isolating subtrees with reasonably high TC while allowing for some amount of gene deletion.

Two additional types of TCCs were designed for use in the MammalSubgroups and MammalSubgroupsPlusOutgroup methods. Inspired by the alignment filtering method from Pollard et al. [2010], which required at least one sequence to be present from each of the three major mammalian superorders (Primates, Glires, and Laurasiatheria) for a column to pass through the filter, the TC_{all} TCC required that the TC for all of the included clades was above a given minimum value. To complement the TC_{all} constraint, the TC_{any} constraint was designed to require that the TC for any of the included clades was above a given value. These more complicated TCCs were included in the analysis in order to determine whether combinations of more specialized constraints would perform as well or better than the simplest approach at isolating LOTs from the Compara gene trees.

The methods within the Orthologs category of subtree sets were implemented separately from the rest. Instead of splitting Compara trees based on TCCs, subtrees in the Orthologs category were defined from the set of genes annotated by Ensembl as orthologs to each gene from a given source species. Thus, for each gene from the source species, the Compara subtree containing all Ensembl-annotated orthologs was extracted and stored; this was guaranteed to yield exactly one subtree for every gene in the source species. Human, mouse, zebrafish, and drosophila were chosen as source species for testing. This approach differed from the tree-splitting strategy in two ways: first, it made use of the orthology annotations resulting from Ensembl's orthology pipeline, which uses the Compara gene trees as its source. Second, this approach did not guarantee that each subtree would contain a completely unique set of genes. For example, a gene which was recently duplicated in

| Method | | |
|------------|-----------------------------|--|
| Category | Name | Constraints |
| Ingroup | Primates | $TC(Primates) > 0.6$ |
| | Glires | $TC(Glires) > 0.6$ |
| | Laurasiatheria | $TC(Laurasiatheria) > 0.6$ |
| | Sauria | $TC(Sauria) > 0.6$ |
| | Fish | $TC(Clupeocephala) > 0.6$ |
| Outgroup | Eutheria | $TC(Eutheria) > 0.6$ |
| | Amniotes | $TC(Amniota) > 0.6$ |
| | Vertebrates | $TC(Vertebrata) > 0.6$ |
| Subgroups | Fungi/Metazoa | $TC(Fungi/Metazoa) > 0.6$ |
| | MammalSubgroups | $TC_{all}(Laur., Glires, Primates) > 0.1$ |
| | MammalSubgroupsPlusOutgroup | $TC_{all}(Laur., Glires, Primates) > 0.1 \text{ AND } TC_{any}(Sauria, Clupeo., Ciona, Marsup.) > 0$ |
| Orthologs | Human Orthologs | |
| | Mouse Orthologs | |
| | Zebrafish Orthologs | |
| | Drosophila Orthologs | |
| Root Nodes | Ensembl Roots | |

Table 4.1: Subtree constraints used for identifying Eutherian orthologous subtrees. Ensembl gene trees were split into subtrees based on taxonomic coverage (TC) requirements at internal nodes. Laur. - Laurasiatheria; Clupeo. - Clupeocephala; Marsup. - Marsupiala

the human terminal lineage would yield two subtrees, one for each human paralog, with identical sets of non-human genes in each tree. Although the orthology-based method might be useful when an evolutionary study is focused on a specific target or reference species, as is often done with human and mouse due to their finished genome sequence and high-quality annotation, I considered it to be less applicable to the current study due to the potential for introducing reference genome-specific biases such as over-representation of genes with gene family expansions in the reference species or non-representation of genes which have been deleted in the reference species. Still, I expected that the sets of subtrees resulting from the Ensembl ortholog annotations would serve as a useful reference against which to compare the methods based purely on TCCs.

The subtree splitting scheme was applied to the 18,607 gene trees from the Compara database, producing a set of subtrees for each of the TCCs shown in Table 4.1. In the next two sections I will describe the resulting sets of trees and subtrees and discuss what they reveal about the evolutionary history of vertebrates and the feasibility of using TCCs to isolate mammalian LOTs for sitewise analysis.

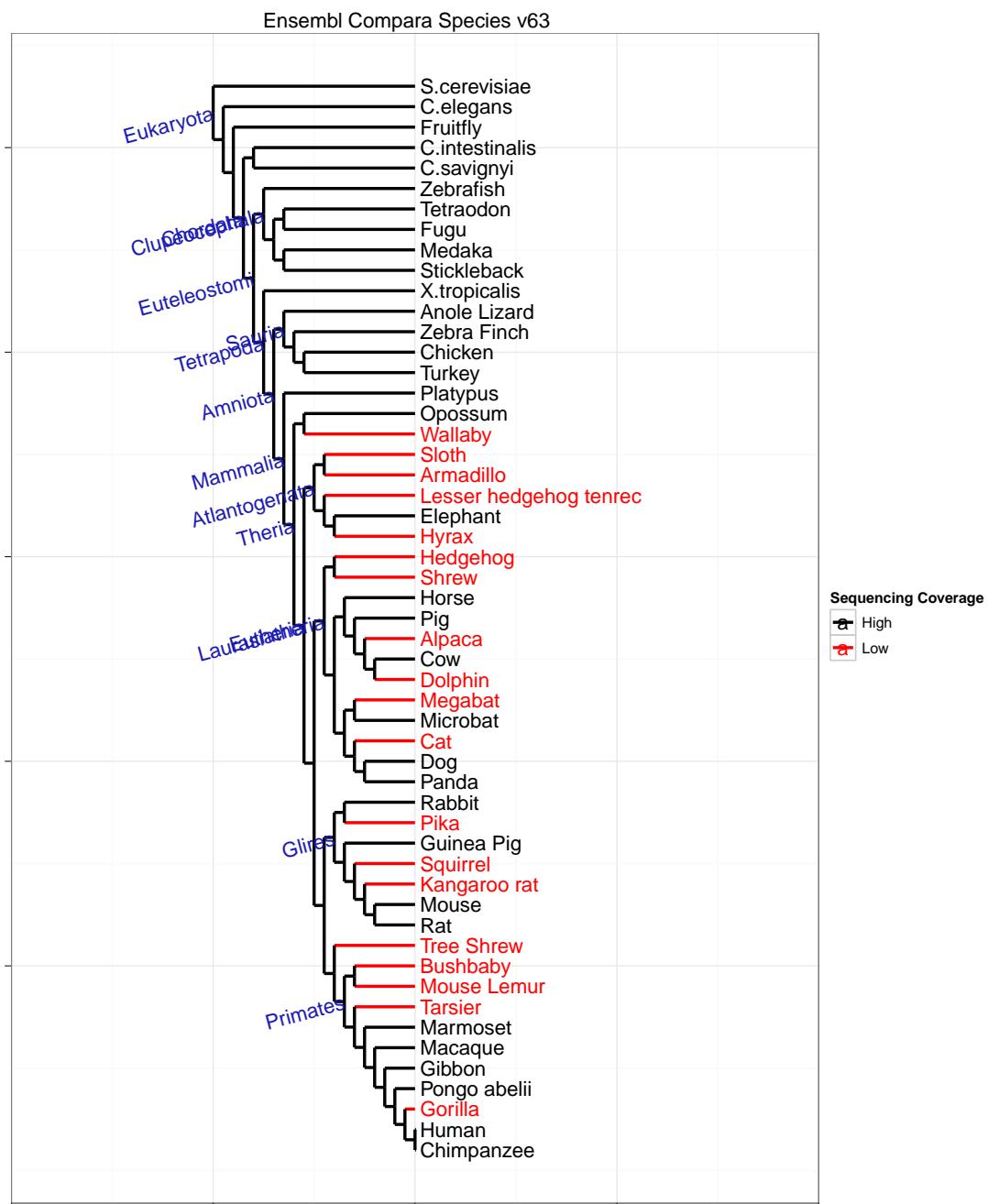


Figure 4.2: The NCBI taxonomy of species within the Ensembl Compara database. Note that branch lengths are not drawn to scale. Low-coverage genomes are labeled in red, high-coverage genomes are in black. Selected internal nodes, used are labeled in blue.

| Tree Set | Count | Med. Size (Min / Max) | N50 | Human Content | | | Human Total | Med. MPL | Med. Species |
|-------------|-------|--------------------------|-----|---------------|------|------|----------------|-------------|-----------------|
| | | | | 0 | 1 | 2+ | | | |
| All | 18607 | 15 (2 / 400) | 139 | 0.50 | 0.30 | 0.20 | 19995 | 0.55 | 8 |
| (≤ 15) | 9378 | 3 (2 / 15) | 5 | 0.92 | 0.08 | 0.00 | 809 | 0.04 | 2 |
| (> 15) | 9229 | 54 (16 / 400) | 146 | 0.07 | 0.53 | 0.40 | 19186 | 1.04 | 47 |

Table 4.2: Summary of the set of Ensembl Compara root trees. The 'Human Content' columns represent the fraction of trees which contain the indicated number of human genes, and 'Human Total' is the total number of human genes contained within the tree set. 'Med. Species' is the median species count across all trees. Med. - median, MPL - mean path length

4.7 Analysis of the set of root Compara gene trees

Table 4.2 presents a summary of the set of root Compara gene trees and the subsets of trees with greater or fewer than 15 sequences.

A first observation was that, somewhat surprisingly, nearly half of all Compara gene trees contained few sequences: 9,378 out of 18,607 root trees constituted fewer than 15 sequences. Given the protein-based clustering performed by the Compara pipeline, one might expect many of these small trees to represent portions of larger fast-evolving gene trees whose high sequence divergences made the BLAST search step inaccurate or caused clustering via the *hcluster_sg* algorithm to be ineffective. Alternatively, these small gene trees might have resulted from extensive lineage-specific gene duplications or from pseudogenes being mis-annotated as genes, in either case causing tight clusters of very closely-related transcripts that would have been identified by *hcluster_sg* as independent clusters of homologs. Some evidence for the latter scenario could be found in the median species counts and mean path lengths of the smaller (fewer than 15 sequences) versus larger (greater than 15 sequences) trees, shown in the second and third rows of Table 4.2. The set of small root trees had a median species count of 2, compared to 47 for the large trees, indicating that the smaller trees encompassed sequences from a very small taxonomic range. Furthermore, the median mean path length (MPL) for small trees was 0.04 compared to 1.04 for the large subset, revealing a much smaller level of sequence divergence within each tree. These two summary statistics together provided evidence that the smallest gene trees in the Compara database consist of highly species-specific, closely-related proteins; it was likely that many of these small trees represented artifactual gene annotations and would be most appropriately removed prior to any downstream analysis.

Despite the existence of many small gene trees in the Compara database, they comprised only a small fraction of all protein-coding sequences. Only 809 human genes, or 4% of the

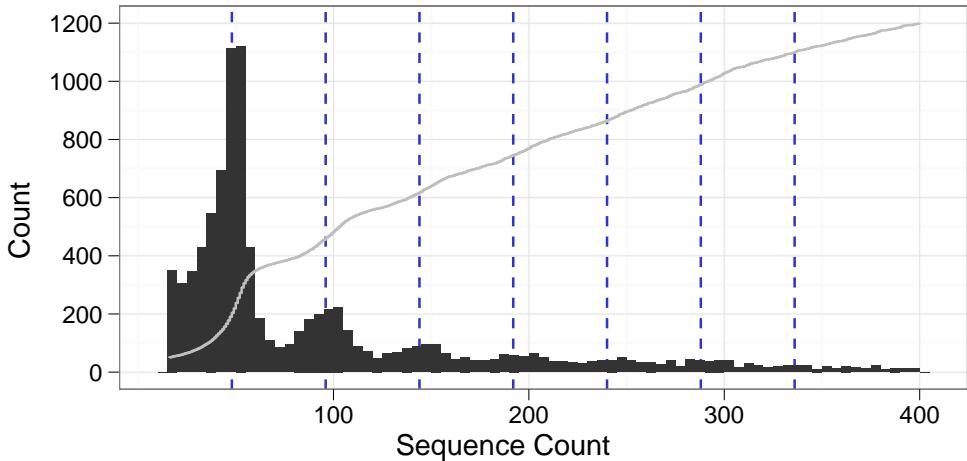


Figure 4.3: Sequence counts for the set of root Compara trees. Black bars show a histogram of sequence counts in bins of width 5, and the gray line shows the cumulative fraction of sequences contained within trees of that size or smaller. For clarity, 9,378 trees with 15 or fewer sequences are not shown. Dashed blue lines are drawn at integral multiples of 48, the number of vertebrate species within Ensembl.

total human gene set—a gene set which I expected be well-annotated and to contain few false positive genes due to the high level of manual curation and the large amount of continued scrutiny—was contained within the subset of small trees. This indicated that whatever process was causing the Compara pipeline to yield a high number of small gene trees did not have much of an impact on the overall placement of human protein-coding genes within the database of root Compara gene trees.

As mentioned in Section 4.5, a closer examination of the distribution of tree sizes in the set of root Compara trees presented a clear view of the over-clustering of mammalian LOTs. The black bars in Figure 4.3 show the distribution of sequence counts for all trees with more than 15 sequences, with vertical dashed lines overlaid at multiples of 48, the number of vertebrate species in Ensembl release 63. The highest peak of the histogram is at or slightly above 48 sequences, with the tree counts quickly diminishing at larger sizes. Weaker, but still discernable, peaks exist at larger tree sizes, with the location of these wave-shaped peaks corresponding closely to the second, third and fourth multiples of 48. The pattern of recurring peaks is indistinguishable at sizes above 200, but there is still a long tail of large trees extending out to a maximum size of 400 sequences. Overall, the distribution of tree sizes provided good support for the situation described in Section 4.5, with the Compara pipeline often clustering together two or more mammalian LOTs related

by a more distant duplication event.

I found it interesting to characterize the set of root Compara gene trees by the proportion of all protein sequences which are covered by trees of a given size or smaller. This value is plotted in Figure 4.3 as a gray line. First, this plot showed that trees with fewer than 15 sequences (which were excluded from the plot but included in the calculation of the cumulative distribution shown) represented a trifling fraction of the sequences within the Compara database; this observation was similar, but not identical, to the statistic noted above, that only 4% of human genes were contained within small trees. Second, the steady upward slope of the cumulative curve contrasted with the declining height of the histogram bars at higher sequence counts. This was a result of the increasing number of sequences encompassed by each of the larger trees: although relatively few trees contained more than 300 sequences, together they comprised around 10% of all protein-coding genes in Compara. Two points along this cumulative plot were of particular interest. First, one could use the height of the cumulative line at the beginning of the second hump in the histogram (at around 75 sequences) to identify the fraction of vertebrate proteins which exist in the Compara database as members of a duplicated mammalian LOT. The line showed that around 30% of proteins are covered by trees of 75 sequences or fewer, meaning that 70% of vertebrate proteins are contained within large gene trees containing sequence-based evidence of ancient paralogy. Second, using the cumulative line in the reverse direction could identify the tree size above which 50% of all sequences were clustered. This value represents the size of tree that an “average” protein might be clustered in, and in some ways is a more accurate characterization of a set of gene trees than the median tree size. A similar calculation is often performed to characterize the size distribution of contiguous sequence blocks (contigs) within a genome assembly. This statistic, referred to as the N50 length, is the contig length for which 50% of bases are contained in contigs of that size or larger [Miller *et al.*, 2010]. For the Compara gene trees, the N50 tree size was 139, slightly less than three times the number of vertebrate species (Table 4.3). Thus, the average protein in the set of Compara gene trees belonged to a tree with 139 sequences, corresponding roughly to three mammalian LOTs related by ancient duplication events, likely resulting from the two rounds of vertebrate genome duplication.

Another way to characterize the distribution of Compara gene trees was across the taxonomic space. Given the wide range of quality in genome assemblies and annotation sets contained within Ensembl, a question of particular interest was whether levels of gene presence and absence were consistent across different species groups and different levels of

assembly quality. To investigate this question in the context of the root Compara gene trees, data were collected by counting the number of sequences from each species contained within each gene tree. Results of this analysis are presented in Figure 4.4, showing the number of trees containing 0, 1, 2, or more than 3 genes from each of the 53 species within Ensembl. Comparing the ranges of values for each copy count (labeled 0, 1, 2 and 3+), I found that the largest number of trees contained zero copies from any one species (between 8,000–11,000 within vertebrate species), a smaller but large number of trees contained one copy from a species (between 4,000-6,000) and several thousand trees contained two, three or more copies (between 1,000-1,500 for 2 copies and 1,500-2,000 for 3+). Note that these numbers were tabulated independently between species; for example, the 10,000 zero-copy trees in human were counted independently of the 10,000 zero-copy trees in chimpanzee, so nothing could be said about how many trees contained zero copies in both human *and* chimpanzee. As noted in the analysis of tree sizes, the large number of zero-copy trees reflected the large number of small, species-specific trees within the set of Compara gene trees. Similarly, the large number of trees with many copies from each species reflected the trees with multiple mammalian LOTs clustered together.

Comparing numbers horizontally across the range of species in Figure 4.4 revealed that the zero-copy count tended to increase with increased evolutionary distance from human, while the 1, 2 and 3+ copy counts tended to decrease as the distance from human increased. Both trends were most striking at the distant end of the tree, where the five non-vertebrate species are shown. For the increase in zero-copy trees and the decrease in single-copy trees, the strength of the trend at the highest level of divergence can be partly explained by the long branch lengths connecting those species to each other and to the more densely sampled vertebrate clade: the distance-based clustering algorithm might reasonably be expected to produce more false negatives in longer branches for a number of reasons, including the behavior of the *hcluster_sg* algorithm, inaccurate BLAST E-values at large distances, and heterogeneity in evolutionary rates across lineages [Whelan, 2008]. However, the reduced numbers of 2 and 3+ copy count trees in non-vertebrate species were most likely due to the whole-genome duplications at the basal vertebrate lineage, resulting in the non-vertebrate species, which underwent no such genome duplications, being strongly depleted of multi-copy duplicates compared to their vertebrate relatives.

Knowing that many of the low-coverage genomes were annotated using projections based on the human set of transcripts, it was slightly concerning that human and its close primate relatives contained fewer zero-copy genes and more one-copy and two-copy genes

than any other group of vertebrates in the Compara set of trees. There was no *ab initio* biological reason to expect this to be the case, so I suspect that this pattern, which was fairly small in effect, was due to the widespread reliance on human annotation and protein experimental data in the annotation of non-human genomes. There was, however, one region of Figure 4.4 where human and primates did not contain the largest number of one-copy or 2-copy trees: the 3+ copy count for the fish species was greater than that for human or any primates, which was most likely a result of gene duplicates retained after the third round of genome duplication in the teleost ancestor [Jaillon *et al.*, 2004]. The signal resulting from the teleost genome duplication event was clearer in the Fish set of TC-defined subtrees, so I will defer further discussion of this effect to Section 4.8 where that set of trees is described.

Finally, the differences in copy counts between species with low- and high-coverage genome sequences showed the tendency of low-coverage genome sequences to be absent from gene trees or to show reduced copy counts. Low-coverage species contain more zero-copy, roughly the same number of one-copy, and noticeably fewer multi-copy genes than closely-related high-coverage species. These clear differences showed that gene absence in low-coverage genomes should not be taken as evidence for actual gene loss, and that gene duplications are systematically underrepresented in low-coverage genomes. The former point was emphasized in a recent critical analysis of the impact of low-coverage genomes on gene duplication inference in the Compara database [Milinkovitch *et al.*, 2010], but the latter point was largely ignored. Again, this signal was also stronger in the sets of TC-defined mammalian subtrees and will be revisited in the next section.

The preceding analysis of the set of root Compara gene trees, in which I characterized the distribution of trees with respect to size and across the taxonomic space, showed that despite the over-representation of small, species-specific trees, most sequences were contained in trees with biologically plausible sizes given the history of vertebrate whole-genome duplications. A gene tree-based equivalent of the N50 statistic was developed for summarizing the distribution of protein sequences within differently-sized gene trees, and two visual summaries of this distribution were introduced (in Figures 4.4 and 4.3), providing evidence for the clustering of paralogous mammalian sub-trees and for species-based and genome quality-based trends in the breakdown of gene copy counts within the root Compara gene trees.

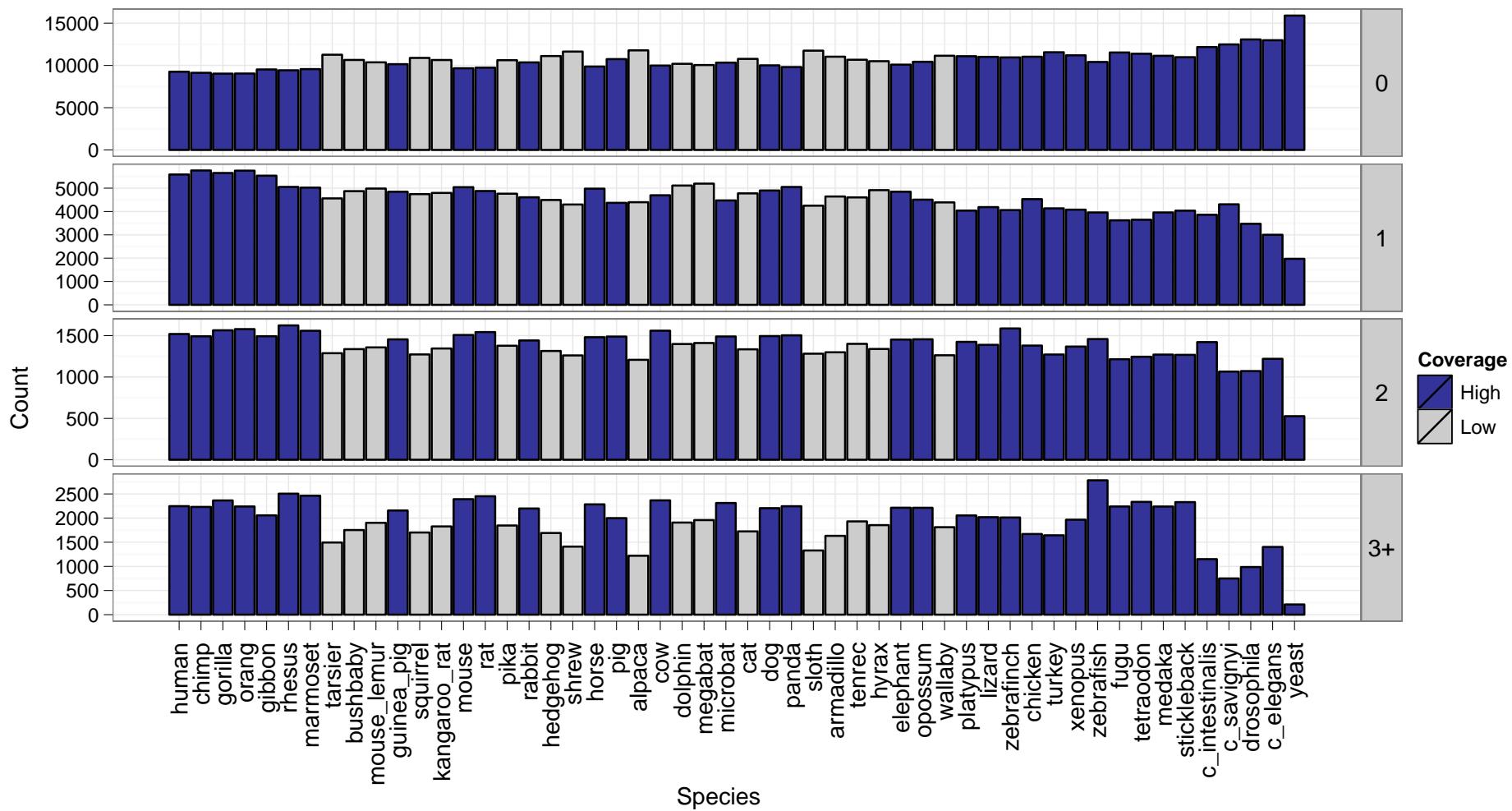


Figure 4.4: Taxonomic distribution of gene copy counts for the root Ensembl trees. The number of trees containing 0, 1, 2 or more than 3 sequences from each species is shown. Bars are colored blue and gray for species with high- and low-coverage genomes, respectively. Note that the y-axis scale is not the same for each panel.

4.8 Analysis of sets of subtrees defined by taxonomic coverage and orthology annotation

The sets of trees resulting from applying the subtree splitting scheme with various TCCs to the Compara gene trees are summarised in Table 4.3, with a summary of the root Compara gene trees and a summary of the set of seven-species amniote gene trees from the Optic database [Heger & Ponting, 2008] included at the bottom for comparison.

The Ensembl Roots and Drosophila Orthologs sets were two clear outliers among the subtree sets shown in Table 4.3, with much higher N50 values (139 and 125 vs. the next highest value of 56) and more trees with multiple human copies (0.20 and 0.43 vs. the next highest value of 0.14) than any other subtree set. In fact, these two subtree sets were very similar except for the excess of small species-limited trees in the Ensembl Roots set: the Drosophila Orthologs set contained fewer trees than the Ensembl Roots (9,210 vs. 18,607) and a larger average tree size (60 vs. 15), and the summary values closely resembled the set of Ensembl Roots with small trees removed (Table 4.2).

Within the Ingroups category, methods based on mammalian TCCs (Primates, Glires and Laurasiatheria) produced largely similar sets of trees, with the Primates set containing around 2,000 more trees and covering around 1,000 more human genes than the Glires and Laurasiatheria sets. There was no readily apparent reason for the higher number and human gene coverage of Primate trees, although it may speculatively be due to an excess of primate-specific gene trees that were not captured by non-primate TCCs. Further investigation of the trees unique to this set would be required to reveal the root cause of this minor discrepancy.

The Sauria set of subtrees was noticeably different from the mammal-based TCCs sets from the Ingroups category. The Sauria clade was represented by only four species in Ensembl and diverged from the mammalian ancestral population at an early point in the evolution of amniotes; it is plausible that the lower clade size and long branch length separating Sauria from the other vertebrate clades caused the moderately lower number of trees (13,046 vs. 15,764 for Laurasiatheria) and the increased proportion of trees containing multiple human genes (0.14 vs. 0.09 for Laurasiatheria).

The Fish clade TCC produced a strikingly different set of trees, resulting from the impact of the teleost-specific whole-genome duplication on the structure of gene trees in fish. Although the Fish subtree set yielded a N50 value of 49, which was no different from the N50 of the other Ingroups sets, Table 4.3 highlights three major differences in the Fish

set: it contained many more trees, a higher proportion of trees with zero human copies, and a lower total human gene count than the other Ingroups sets.

The reason for the drastically different Fish tree set was that the tree splitting procedure was designed to identify the largest non-overlapping set of subtrees that satisfied the given TCCs. Genes where one or both of the duplicate copies from the teleost-specific whole-genome duplication were lost would appear as one-to-one orthologs or deletions with respect to the other vertebrate lineages. Genes that were retained in duplicate form, however, would result in a gene tree with two teleost-specific subtrees, each containing a high TC value (i.e., near or equal to 1.0) for the Clupeocephala clade that contains the fish species within Ensembl. In this case, the splitting procedure would produce two small fish-specific subtrees, “ignoring” the surrounding set of mammalian orthologs because two smaller non-overlapping trees already exceeded the TC threshold of 0.6. If, however, one of the duplicate gene copies were lost, then the tree would resemble a typical singly-orthologous vertebrate gene tree, and the splitting procedure would select a subtree encompassing the entire vertebrate clade. It follows that the presence of small, teleost-specific gene trees in the Fish set is a signal of retained duplicate copies, and the size distribution of trees from the Fish set, shown in Figure 4.5, shows that several thousand trees fit the expected model. If we assume that all trees from the Fish subset which contain zero human copies, span 5 or fewer species, and contain 40 or fewer sequences are likely retained duplicate genes, a total of 6,980 retained duplicates are identified, yielding a retention rate of 17.5%, which is very much in line with a previously published estimate of 15% based on a comparison of tetraodon, fugu and zebrafish genes [Brunet *et al.*, 2006].

The sets of subtrees resulting from the Outgroup methods were of special interest, as the clades used to define these TCCs contained all or nearly all of the mammalian species whose orthologous genes I wished to study. The resulting sets of subtrees showed little variation, owing perhaps to the large sizes of the clades and their similar species composition. Each subtree set contained between 15 to 17 thousand trees, N50 values of around 49, and greater than 90% of trees containing exactly one human sequence. These measures provided strong evidence that the tree-splitting method was accurately isolating mammalian LOTs. Some slight differences between subtree sets were apparent, however, with the tree count decreasing, the proportion of trees with human duplications increasing, and the overall human gene coverage decreasing as the clade size used for the TC calculation increased. These trends could potentially be explained by the minimum required tree size increasing along with the clade size, as a result of the increased number of species required

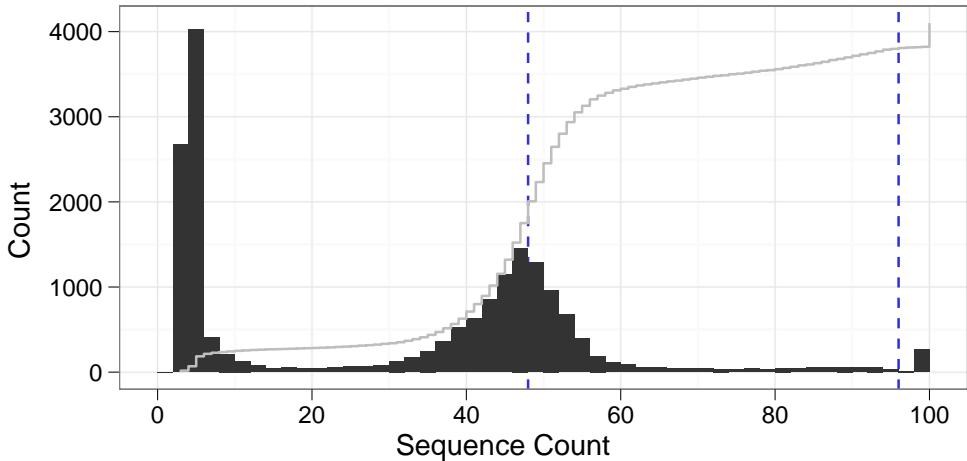


Figure 4.5: Sequence counts for the set of subtrees identified using the Fish clade taxonomic coverage constraint, showing an excess of small subtrees resulting from the teleost genome duplication. Black bars show a histogram of sequence counts in bins of width 2, a gray line shows the cumulative fraction of sequences contained within trees of that size or smaller, and dashed blue lines are drawn at integral multiples of 48, the number of vertebrate species within Ensembl. Trees with more than 100 sequences are included in the topmost bin.

to produce a TC value of 0.6. The minimum subtree size ranged from 21 for Eutheria to 32 for Fungi/Metazoa.

The Subgroups methods did not appear to produce subtrees of any higher quality or more biological interest than the Outgroups methods. The MammalsSubgroups set produced more trees than the Outgroups sets, but the N50 was slightly lower (46 vs. 49) and the proportion of zero-copy human trees was higher (0.18 vs. 0.01), suggesting that the additional trees in the MammalsSubgroups set were spurious subtrees containing limited species coverage. The addition of an outgroup requirement to the MammalSubgroupsPlusOutgroup method produced a tree set more closely resembling the Outgroup methods, but the human gene coverage was lower than that for any Outgroup method despite the overall higher tree count.

Finally, the ortholog annotation-derived subtrees provided for an interesting comparison between the three different ortholog sources and between the overlapping and non-overlapping sets of subtrees. The *Drosophila* ortholog set was highly contrasted with the vertebrate sets due to the two rounds of whole genome duplication, while there was minimal variation among the other ortholog sets. It is interesting to note that the protein-coding transcripts used by the Compara pipeline included 21,873 mouse protein-coding genes and

only 19,991 human genes, indicating either a larger number of true protein-coding genes in mouse or a higher tolerance for false positive gene predictions in the mouse genome compared to the human genome. Zebrafish, on the other hand, contained 24,540 genes; this number agreed well with the 17.5% proportion of retained duplicate genes that I estimated earlier in this section. Overall, 76% and 81% of mouse and zebrafish genes contained an apparent orthologous relationship with one human gene, which was only slightly lower than the 92% of Eutheria subtrees containing one human sequence.

| Subtree Method | | Med. Size | | | Human Content | | | Human | Med. | Med. |
|----------------|-----------------------------|-----------|---------------|-----|---------------|------|------|-------|------|---------|
| Category | Name | Count | (Min / Max) | N50 | 0 | 1 | 2+ | Total | MPL | Species |
| Ingroups | Primates | 17673 | 46 (6 / 388) | 48 | 0.02 | 0.93 | 0.05 | 19024 | 0.68 | 42 |
| | Glires | 15786 | 48 (8 / 391) | 49 | 0.02 | 0.90 | 0.08 | 17904 | 0.73 | 44 |
| | Laurasiatheria | 15764 | 48 (8 / 391) | 49 | 0.01 | 0.90 | 0.09 | 17952 | 0.73 | 44 |
| | Sauria | 13046 | 49 (3 / 391) | 51 | 0.06 | 0.80 | 0.14 | 14988 | 0.78 | 45 |
| | Fish | 18291 | 40 (3 / 391) | 49 | 0.43 | 0.52 | 0.06 | 12183 | 0.58 | 38 |
| Outgroups | Eutheria | 16477 | 47 (21 / 391) | 49 | 0.01 | 0.92 | 0.07 | 18343 | 0.71 | 43 |
| | Amniotes | 15899 | 48 (26 / 391) | 49 | 0.01 | 0.91 | 0.08 | 18094 | 0.73 | 44 |
| | Vertebrata | 15634 | 48 (29 / 391) | 49 | 0.01 | 0.91 | 0.08 | 17938 | 0.74 | 44 |
| | Fungi/Metazoa group | 14957 | 48 (32 / 391) | 50 | 0.01 | 0.90 | 0.09 | 17623 | 0.76 | 44 |
| Subgroups | MammalSubgroups | 21179 | 40 (4 / 159) | 46 | 0.18 | 0.79 | 0.03 | 18595 | 0.54 | 37 |
| | MammalSubgroupsPlusOutgroup | 17155 | 46 (5 / 159) | 48 | 0.05 | 0.90 | 0.05 | 17640 | 0.71 | 43 |
| Orthologs | Human Orthologs | 19991 | 49 (2 / 367) | 52 | 0.00 | 1.00 | 0.00 | 19991 | 1.07 | 44 |
| | Mouse Orthologs | 21873 | 50 (2 / 352) | 54 | 0.10 | 0.81 | 0.09 | 28256 | 1.01 | 43 |
| | Zebrafish Orthologs | 24540 | 51 (2 / 392) | 56 | 0.11 | 0.76 | 0.13 | 30063 | 1.14 | 46 |
| | Drosophila Orthologs | 9210 | 60 (2 / 399) | 125 | 0.08 | 0.49 | 0.43 | 17625 | 1.22 | 50 |
| Default Trees | Ensembl Trees | 18607 | 15 (2 / 400) | 139 | 0.50 | 0.30 | 0.20 | 19995 | 0.55 | 8 |
| | Optic Trees | 17372 | 9 (2 / 789) | 9 | 0.12 | 0.79 | 0.09 | 18477 | 0.00 | 8 |

Table 4.3: Summary of Ensembl subtrees identified using taxonomic criteria or Ensembl ortholog annotations. The set of Compara gene trees from Table 4.2 and the set of trees from the Optic database [Heger & Ponting, 2008] are included at the bottom for comparison. Cells in numeric columns are shaded according to their value relative to other rows, with low values in white and high values in blue. The 'Human Content' columns represent the fraction of trees which contain the indicated number of human genes. 'Med. Species' is the median species count across all trees. Med. – median, MPL – mean path length

In order to easily compare patterns between subtree sets and across the taxonomic space, I tabulated gene copy counts within each Ensembl species for each generated subtree set; the results of this analysis are shown in Figure 4.6. By way of reference, the values shown in each of the separate panels of Figure 4.4 appear in the bottom panel of Figure 4.6 as stacked bars of different colors. Although the various summary characteristics of each of the subtree methods have already been discussed at length, this more taxonomic-oriented view revealed some salient features of the patterns of gene deletion and duplication within the tree sets and showed the pervasive impact of genome-wide duplications on the evolution of vertebrate genes. The large fraction of species with multiple copies in the Drosophila Orthologs subtree set (Figure 4.6, blue and purple bars) showed the effect of two rounds of vertebrate genome evolution, producing a large number of trees with multiple Drosophila orthologs in vertebrate species; similarly, the elevated fraction of multi-copy fish trees in the Outgroups subtree sets (e.g., the blue bars in the Eutheria panel of Figure 4.6) showed the impact of the teleost-specific duplication event on the number of mammalian LOTs with multiple fish copies.

4.9 Gene duplication and loss in the set of Eutherian largely orthologous trees

The subtrees defined by the Eutheria TCC were chosen as the final set of gene trees for use in the downstream sitewise analysis. This set was chosen due to its slightly larger number of trees and better coverage of human genes compared to the other subtree sets from the Outgroups category (Table 4.3).

The distribution of tree sizes for the set of Eutherian subtrees is shown in Figure 4.8. The histogram of tree sizes showed a single peak at exactly the number of vertebrate species in Ensembl, with no trees containing fewer than 20 sequences and a few hundred trees with more than 100 sequences. The distribution of tree sizes was consistent with the set of 16,477 gene trees representing an accurate set of genome-wide mammalian LOTs, with variations in sequence counts resulting from sporadic gene duplication or loss events or unannotated genes in low-coverage genomes.

I also analyzed the detailed taxonomic distribution of gene duplications and losses implied by the set of Eutherian subtrees, as the relative prevalence of zero-copy and multi-copy trees in this set of LOTs might provide some indication of whether gene deletion or gene duplication has been more prevalent in the evolution of vertebrate genomes. Figure

?? shows the gene copy counts for the set of Eutherian subtrees across all Ensembl species, similar to what Figure ?? showed for the set of root Compara gene trees. The excess of zero-copy trees and deficit of duplications in low-coverage genomes was immediately apparent from Figure ??, confirming the trend observed in the set of root Compara trees. Similarly, the trend of increased zero-copy trees in non-mammalian and non-vertebrate species was stronger than in Figure ??, as was the excess of trees with multiple copies in fish genomes.

A quantitative comparison of the number of multi-copy versus zero-copy trees in each species showed that gene duplication has had a greater apparent impact on gene copy counts than gene deletion, at least within primates and most mammals with high-coverage genomes. For example, human contained 200 zero-copy trees and 1,140 combined two- or three-copy trees within the set of Eutherian subtrees, showing evidence for a greater prevalence of gene duplication than gene deletion in LOTs since the common Eutherian ancestor. On the other end of the spectrum in primates was gibbon, which had the most zero-copy trees (878) and the fewest multi-copy trees (846) of all the primates, showing roughly equal tendencies towards gene deletion and duplication across the set of Eutherian LOTs. This pattern was consistent across the mammalian tree, save for a few exceptions: n guinea pig showed roughly the same number of zero-copy and multi-copy trees (1012 vs. 1197, respectively), rabbit showed slightly more zero-copy trees (1569 vs. 1394), and pig showed a much higher number of zero-copy trees (3412 vs. 1953). Beginning with opossum, vertebrates more distantly related to the Eutherian common ancestor showed greater numbers of zero-copy trees, leveling off at ca. 4,000 zero-copy trees, and higher variation in the number of multi-copy genes; at the low end were chicken and turkey with 876 and 972 multi-copy genes, respectively, and at the high end (excluding the fish species) were zebrafinch and platypus with 1904 and 2307 multi-copy genes, respectively.

In the analysis of vertebrate genomes, one must be vigilantly aware of potential biases arising from the frequent reliance on homology with the well-annotated human and mouse genomes in the annotation of newly sequenced genomes. Such a bias could have plausibly led to anomalous results in the present analysis of gene copy counts, for example by reducing the chance of correctly identifying gene trees containing deletions in human or mouse, thus over-inflating the prevalence of duplications versus deletions. The level of consistency seen in the relative numbers of zero-copy and multi-copy trees across the range of mammals provided some evidence against such a bias, although it did not entirely rule out the possibility, as all mammalian genomes may have been similarly affected by the use of

human or mouse proteins in the Ensembl gene annotation pipeline.

An unexpected result of the comparison between mammalian species was the identification of certain genomes with uncharacteristically high numbers of zero-copy or multi-copy trees. The most striking example was pig, which contained nearly double the number of zero-copy trees of any other high-coverage mammalian genome and a noticeably elevated number of multi-copy trees, but rabbit and guinea pig also deviated from the normal patterns. Given the otherwise consistently low number of zero-copy trees throughout the range of Eutherian mammals, I would expect the number of zero-copy trees in these species to decrease as finished-quality genome sequences are produced and the gene annotation pipelines are further optimized to work with each species. In particular the anomalous nature of the pig gene trees may be related to the draft quality of the genome sequence and assembly; I would expect the number of zero-copy trees to be substantially reduced once a finished-quality genome sequence and annotation set is made available [[Archibald *et al.*, 2010](#)].

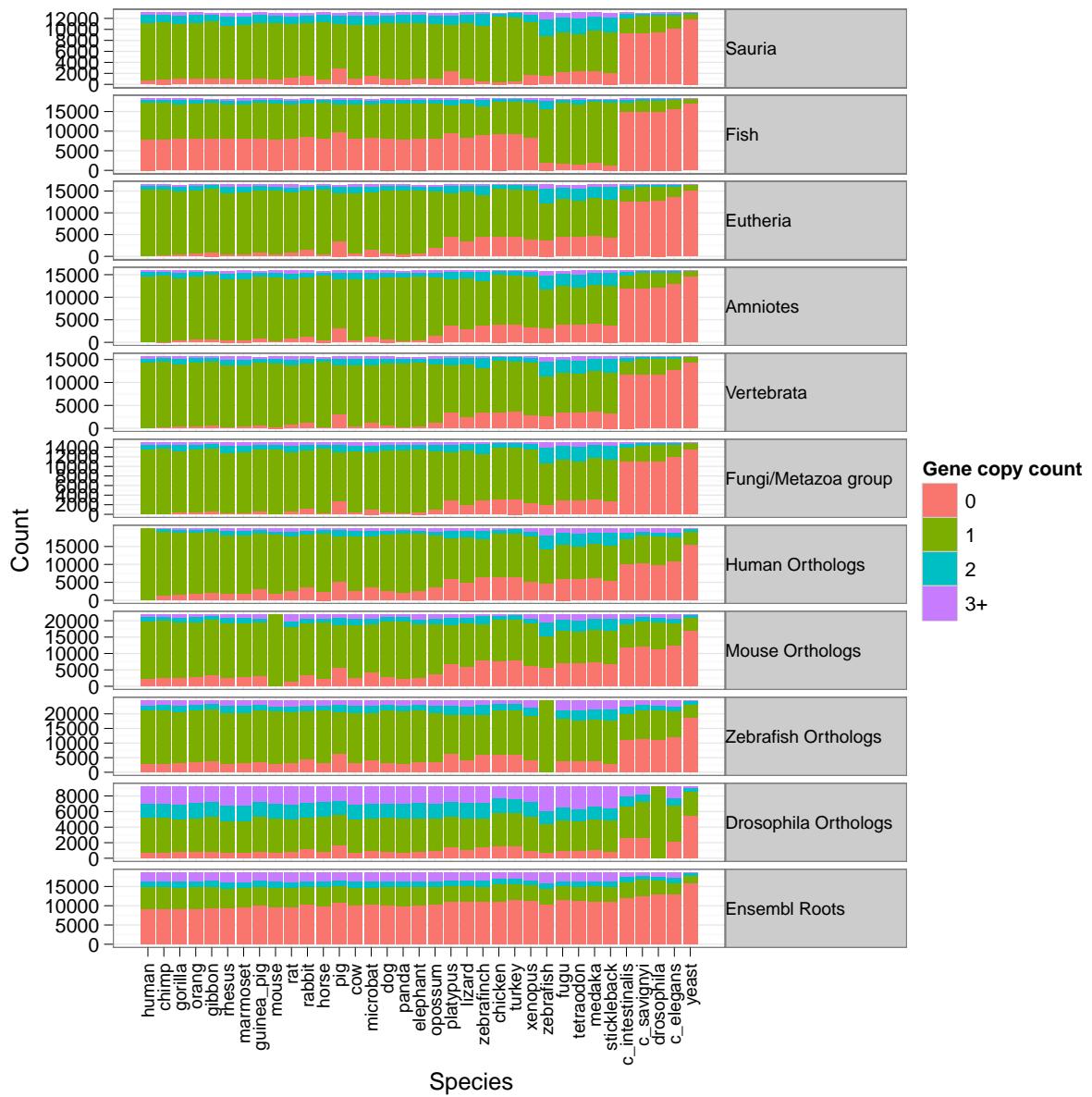


Figure 4.6: Taxonomic distribution of gene copy counts for different subtree methods. The numbers of trees containing 0 (red), 1 (green), 2 (blue) or more than 3 (purple) sequences from each species are shown as stacked colored bars. The Ingroup and Subgroups methods were omitted for clarity, as were species with low-coverage genomes. Note that the y-axis scale is not the same for each panel.

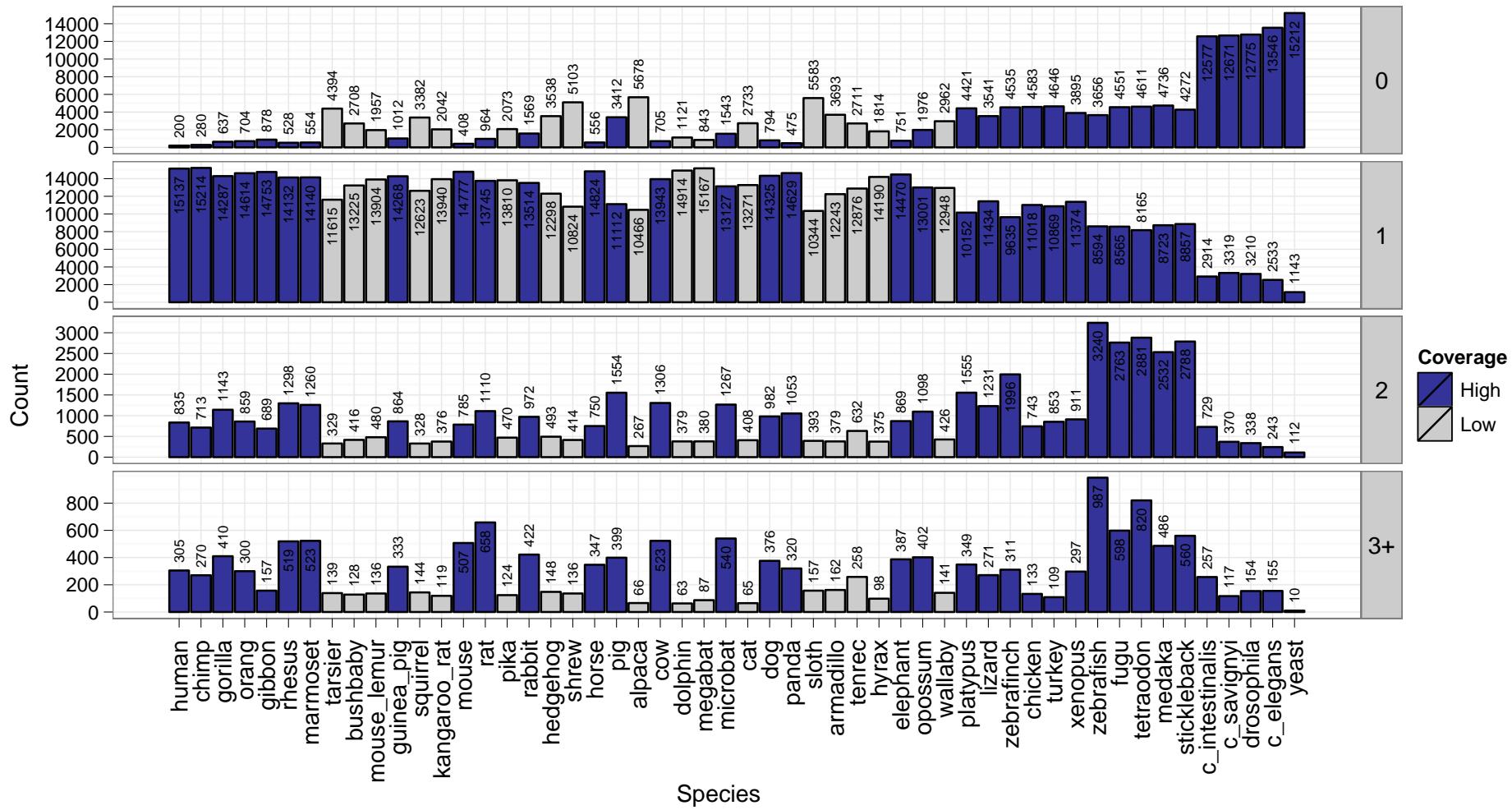


Figure 4.7: Taxonomic distribution of gene copy counts for the Eutheria subtrees defined by taxonomic coverage constraints. Each panel from top to bottom shows the number of trees containing 0, 1, 2 or more than 3 sequences from each species. Bars are colored blue and gray for species with high- and low-coverage genomes, respectively. Note that the y-axis scale is not the same for each panel.

4.10 Comparison to gene trees from the Optic database of amniote orthologs

Given the tendency of the root Compara gene trees to contain multiple mammalian LOTs, I wished to evaluate whether a different phylogenetic tree-based orthology pipeline produced a similar distribution of gene tree sizes. The Optic database, a project from the Ponting group which used an independently designed gene-building and comparative genomics pipeline combined with the TreeBeST software to infer duplication-resolved gene trees within a variety of vertebrate and invertebrate clades, was ideal for such a comparison [Heger & Ponting, 2008]. I downloaded the set of Optic gene trees inferred from a group of eight vertebrate genomes (human, mouse, dog, opossum, platypus, chicken, zebrafinch, and tetraodon), characterized the set of gene trees using a variety of summary statistics, which are included in the bottom row of Table ??, and plotted the distribution of gene tree sizes in Figure 4.9.

After accounting for differences in the number of sampled species, the set of Optic gene trees more closely resembled the set of Eutherian subtrees than the root Compara gene trees. Nearly 80% of the Optic gene trees contained exactly one human sequence and only 9% contained two or more human sequences; the large proportion of single-copy trees suggested that the Optic trees did not contain many “over-clustered” mammalian LOTs, and the 9% of trees with multiple human genes was close to the 7% seen in the set of Eutherian subtrees. Figures 4.8 and 4.9 clearly show the similar sequence count distributions between the Optic and Eutherian trees; each histogram has a distinct peak at a sequence count corresponding to the number of species included in the database, with no sign of the long tail of large gene trees seen for the distribution of root Compara gene trees in Figure ??.

4.11 Conclusions

This chapter was concerned with identifying a set of gene trees for subsequent sitewise evolutionary analysis. Three characteristics were desired in the ideal set of trees: maximal coverage of available mammalian genomes, minimal inclusion of paralogous relationships, and consistent taxonomic representation throughout the set of trees. I chose the Ensembl Compara database as a source of root gene trees due to its well-established methodology [Heger & Ponting, 2008; Vilella *et al.*, 2009] and its increased species sampling compared

to equivalent vertebrate databases such as Treefam and Optic.

Using the set of root Compara gene trees as a basis for further analysis, I characterized the distribution of gene trees in a variety of ways, including a tree-based analogue of the N50 statistic commonly used in evaluating genome assemblies. This analysis showed that the “over-clustering” of multiple Eutherian LOTs within single Compara gene trees had a major impact on the composition of the gene tree set. I then developed a simple scheme for isolating LOTs from within larger gene trees using flexible TCCs and applied the system to the set of Compara trees to generate several sets of genome-wide TC-defined subtrees.

These sets of TC-defined subtrees provided a number of insights into the orthology and paralogy relationships within and between mammalian and vertebrate genomes, including a quantification of the proportion of duplicate genes that were retained after the teleost whole-genome duplication that matched the prediction resulting from a detailed analysis of fish genomes [Brunet *et al.*, 2006]. The comparison between subtree sets resulting from different TCCs showed that a simple threshold based on TC in the Eutheria clade produced a set of subtrees that contained a high percentage of human genes and satisfied the three desired characteristics for subsequent sitewise analysis.

An analysis of the number of trees with zero, one, two or more sequences for a given species revealed patterns of variation in gene copy counts across vertebrate genomes. Importantly, low-coverage genomes showed a large increase in zero-copy trees (e.g. trees with no sequences from the low-coverage genome) and a notable decrease in multi-copy trees. Other, more subtle trends were identified, including individual genomes, such as pig, with apparently anomalous numbers of zero- or multi-copy trees, and more widespread trends, such as the relatively few multi-copy trees in avian genomes and the higher prevalence of multi-copy trees versus zero-copy trees within mammals.

The set of Eutherian LOTs identified in this chapter represented, to my knowledge, the largest and most accurate set of mammalian orthologous trees identified to date. Although the set of Eutherian subtrees were isolated from the Compara root trees and not independently derived, the separation of the “over-clustered” Compara trees into biologically realistic sets of Eutherian genes sharing mostly orthologous relationships was an important and fundamental part of preparing these gene trees for subsequent evolutionary analysis.

One might argue that a gene tree with fully-resolved duplication and speciation events, as provided by the Compara database, should be enough to allow for identification of orthologous versus paralogous subtrees. In theory this is true, but such an approach to identifying LOTs would have failed on two grounds: first, it would not distinguish between

ancient paralogs resulting from the two rounds of whole-genome duplication in vertebrates and more recent paralogs resulting from the gradual accrual of gene duplications over time, causing recent duplication events to be handled in the same way as ancient duplications absent the incorporation of taxonomic information. Second, in practice the resolution of gene deletions within the Compara gene trees was far from perfect, a point raised by Milinkovitch et al. [2010] in their criticism of the inclusion of low-coverage genomes in the Ensembl Compara orthology pipeline. A lack of data, caused either by missing genes from low-coverage genomes or by the limited amount of sequence data available for phylogenetic reconstruction, may often contribute to errors in the accurate identification of duplication and speciation nodes within the Compara gene trees; the use of TC information sidestepped this source of uncertainty by not relying on accurate resolution of duplication nodes.

The performance of the TC-based scheme described in this chapter showed that simple taxonomic criteria could be used to identify largely orthologous subtrees. The utility of a taxonomic-based approach to identifying meaningful gene tree subtrees was indirectly validated by comparison to gene trees from the Optic orthology database, which showed many similar characteristics to the set of Eutherian subtrees; the pipeline used to create the Optic database includes an explicit step for splitting trees based on taxonomic criteria [Heger & Ponting, 2008], suggesting that a similar technique was found to be important in identifying LOTs within the context of an independently designed orthology pipeline.

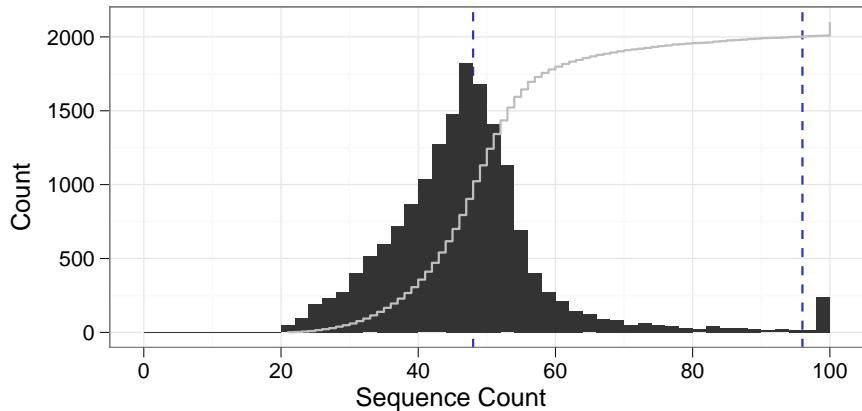


Figure 4.8: Sequence counts for the set of subtrees identified using the Eutheria clade taxonomic coverage constraint. Black bars show a histogram of sequence counts in bins of width 2, a gray line shows the cumulative fraction of sequences contained within trees of that size or smaller, and dashed blue lines are drawn at integral multiples of 48, the number of vertebrate species within Ensembl. Trees with more than 100 sequences are included in the topmost bin.

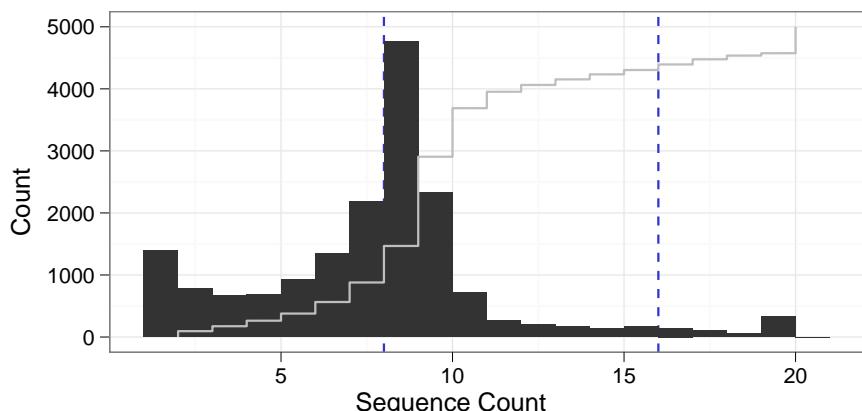


Figure 4.9: Sequence counts for gene trees in eight amniote species from the Optic orthology database [Heger & Ponting, 2008]. Black bars show a histogram of sequence counts, a gray line shows the cumulative fraction of sequences contained within trees of that size or smaller, and dashed blue lines are drawn at integral multiples of 8, the number of amniote species analyzed by Optic. Trees with more than 20 sequences are included in the topmost bin.

Chapter 5

Patterns of sitewise selection in mammalian genomes

5.1 Introduction

5.1.1 The Mammalian Genome Project

A major goal of mammalian comparative genomics has been to quantify, identify and understand the fraction of the human genome that is under evolutionary constraint. The first non-human mammalian genomes showed at least 5% of the human genome to be under purifying selection [Lindblad-Toh *et al.*, 2005; Mouse Genome Sequencing Consortium & Mouse Genome Analysis Group, 2002; Rat Genome Sequencing Project Consortium, 2004], but the small number of genomes available limited the extent to which regions of evolutionary constraint could be identified. The Mammalian Genome Project, a coordinated set of genome sequencing projects initiated in 2005 and organised by the Broad Institute of MIT and Harvard, was designed with the primary purpose of increasing the accuracy and confidence with which regions of the human genome that have evolved under evolutionary constraint in mammals could be identified [Margulies *et al.*, 2007]. In line with this goal, 20 mammalian species were chosen for sequencing in order to maximise the amount of evolutionary divergence available for comparative analysis when combined with the 9 already available sequenced genomes [Margulies *et al.*, 2005]. To save on sequencing costs most of the 20 additional species were only sequenced to a target twofold coverage, meaning each genomic base pair would be covered on average by two sequence reads and roughly 85% of genomic sequence would be covered by at least one read.

As the Mammalian Genome Project proceeded from its sequencing to analysis phase in late 2008, it became clear that the additional branch length afforded by the 29-species phylogeny would enable a number of improved evolutionary analyses beyond the identification of constrained non-coding regions. Among others, these included the evolutionary characterisation of gene promoters, identification of exapted non-coding elements, detection of evolutionary acceleration and deceleration in non-coding regions, and detection of purifying and positive selection in protein-coding genes. Given the prior involvement of the Goldman group in analysing the ENCODE comparative sequencing data [[ENCODE Project Consortium, 2007](#); [Margulies *et al.*, 2007](#)] and Massingham’s work on a method and software program for sitewise evolutionary analysis [[Massingham & Goldman, 2005](#)], the group became involved in the protein-coding evolutionary analysis for the Mammalian Genome Project. This chapter describes my work on the project, which began in late 2008; all of the work described below was performed by me in consultation with members of the Goldman group (Nick Goldman and Tim Massingham), EnsEMBL team (Albert Vilella, Javier Herrero, Ewan Birney) and organisers and members of the Mammalian Genome Project (Manolis Kellis, Kerstin Lindblad-Toh, Mike Lin, Katie Pollard), and the major results from this analysis have been published [[Lindblad-Toh *et al.*, 2011](#)].

5.1.2 The Sitewise Likelihood Ratio test

As described in Section [1.3.1](#), differential survival of non-synonymous and synonymous mutations based on the degeneracy of the genetic code can be used as a source of information on the continued importance of mutations at a given protein-coding site over evolutionary time: a lower rate of non-synonymous substitution compared to synonymous substitution is indicative of purifying selection, or natural selection in favor of maintaining protein structure and function; equal rates of non-synonymous and synonymous substitutions is indicative of neutral selection, or no differential survival of protein-altering mutations; a greater rate of non-synonymous than synonymous substitution is indicative of positive selection, or natural selection in favor of protein-altering mutations.

Early evolutionary analyses of protein sequences showed large variation in the rates of amino acid change both within and between proteins resulting from the myriad structures and functions embodied by different proteins and protein domains [[Kimura & Ohta, 1974](#)]. Continued work suggested that the overall evolutionary rate [Koonin & Wolf \[2006\]](#) and the pattern of localised selective pressures [Nielsen & Yang \[1998\]](#); [Yang & Nielsen \[1998\]](#) of a gene can reveal important insight into its role in the organism, establishing the study of

rates of non-synonymous and synonymous substitution in proteins as an effective method for using evolutionary information to investigate the functional characteristics of genes.

Maximum likelihood methods, introduced in Section 1.3.2, are commonly applied to biological sequence analysis due to their desirable statistical features... [[Finish out the paragraph with a quick recap of ML methods and their use in detecting selection]]

The Sitewise Likelihood Ratio (SLR) test is based on the mechanistic Goldman-Yang codon model of evolution, with an additional parameter for each site in the alignment representing the sitewise ω value. The inclusion of an additional parameter per alignment site makes the model extremely complex and difficult to optimise, but the dimensionality of the likelihood optimisation is reduced by making the assumption that the ω at each site does not contribute significantly to the overall likelihood, thus allowing for separate optimisation of the global parameters (including XYZ) across the whole alignment and the ω parameter at each site. In this way, SLR performs an approximate likelihood ratio test (LRT) for non-neutral evolution at each site...

5.1.3 Data quality concerns: alignment and sequencing error

The possibility that errors in the source alignments might cause false positives in the detection of sitewise positive selection was a major concern for this analysis. Although the SLR test and other sitewise maximum likelihood methods have been shown to be conservative for detecting positive selection even when the amount of data is low or the null model is violated [Anisimova *et al.*, 2002, 2003; Massingham & Goldman, 2005], most evolutionary analyses are based on the assumption that all sites within an alignment column are truly homologous. This assumption can be violated in a number of ways, some of which are described below.

Alignment error results from the difficulty of reconstructing the evolutionary history of sequences evolving with indels and can cause non-homologous codons to be placed in the same alignment column. In Chapters 2 and 3 I explored the tendency of multiple aligners to produce such errors, showing that PRANK_C alignments would be expected to introduce few falsely identified positively-selected sites resulting from alignment errors at mammalian-like divergence levels.

Errors resulting from the inclusion of incorrect genomic sequence in alignments was an additional concern. Twenty of the genomes under study were sequenced at low coverage and were not assembled into chromosomes or finished to completion, making the likelihood of miscalled bases, spurious insertions or deletions, or shuffled regions due to mis-assembly

relatively high [Green, 2007]. The potential effect of each of the aforementioned types of sequence errors on the detection of positive or purifying selection depends on the nature of the inference method, the type of sequencing error, and the branch length of the terminal lineage leading to the species containing the error.

As most codon-based inference methods assume independence between amino acid sites, I first consider the effect—in isolation—of a single spuriously-assigned homologous codon on the maximum likelihood estimation of ω . Two cases can be considered: a single sequence error causing one spurious substitution within a codon, and one or multiple sequence or assembly errors causing multiple spurious substitutions within a codon. In the case of a single spurious substitution, if we assume no large difference between the natural mutational process and the process that caused the erroneous mutation, then the effect would be to shift the estimated ω in the branch containing the error towards 1. The sequence error would be incorporated into the maximum likelihood optimisation as an additional neutral substitution, inflating the estimated substitution rate but not affecting the relative non-synonymous and synonymous rates. This effect may be biased towards higher or lower ω values if a significant difference exists between the neutral biological mutational process and the pseudo-mutational process causing the erroneous substitution. On the other hand, a codon with multiple erroneous bases may cause greater elevation of the inferred substitution rate and ω , due to the necessity of maximum likelihood methods to infer a multi-step path of single substitutions between the two codons on either side of a given evolutionary branch. The path estimated between two completely non-homologous codons depends on the estimated codon frequencies, the genetic code, and the nature of the process causing misalignment of nonhomologous codons; while a detailed investigation of the expected effect on inferred ω values is beyond the scope of my analysis, it is not unreasonable to expect a greater number of false positive PSSs resulting from codons with multiple erroneous bases than from codons with single errors.

Given the potentially greater impact of codons with multiple errors, the propensity of each of the common sequencing error types identified above (miscalled bases, spurious indels, and shuffled/repeated/collapsed regions due to mis-assembly) to cause single or multiple errors within codons could strongly affect its effect on the detection of positive selection. On its own, a miscalled base would obviously result in a single spurious substitution. However, low-quality bases tend not to be uniformly distributed among or within sequence reads, which makes for a larger probability of multiple errors within a codon resulting from miscalled bases. Spurious indels within coding regions may be even more

likely than miscalled bases to cause multiple errors within a codon due to the potential alignment and frameshift effects. Assembly errors, which result in larger-scale structural errors including missing, repeated, shuffled or inverted sequence regions, are most prone to produce codons with multiple erroneous errors due to the large amount of contiguous sequence data being misplaced.

I also note the impact of the inference method and terminal branch length on false positives resulting from sequence errors, which can be understood in terms of the information most directly affecting the inference of a positively selected site or a positively selected gene for a given detection method. Both the branch-site test and the sitewise tests (including SLR and PAML M8) are sensitive to substitutions at a subset of alignment sites, but the branch-site test is specifically sensitive to substitutions along the foreground branches of interest while the sitewise tests detect positive selection only throughout the entire tree. In the latter case, the effect of spurious synonymous and non-synonymous substitutions from sequence data depends on the ratio of the species' terminal branch length to the branch length of the entire tree: a longer terminal branch gives greater weight to the erroneous sequence data, making false positives more likely to result. In the former case of the branch-site test, the potential effect depends on the location and length of the foreground branches. If the terminal branch leading to the spurious sequence is within the foreground and the total foreground branch length is small, then false positives could easily result; if, however, the terminal branch is outside of the foreground then it should have little to no effect on the FPR of the branch-site test. Interestingly, this suggests that branch-site tests where the foreground only consists of internal branches may be less prone to false positives from sequencing error than tests that include terminal lineages in the foreground model.

To summarise, the expected effect of alignment errors on the sitewise detection of positive selection should be minimal when using a good aligner and analysing data within vertebrate divergence levels, but the number of false positives resulting from sequence errors depends on a number of factors including the frequency, spatial clustering, and phylogenetic branch length associated with sequencing-based errors when applied to detecting sitewise positive selection. In some cases even a large amount of sequencing error should not produce a strongly elevated FPR (e.g., when the total branch length is large, when analysing all mammals or vertebrates) but in other cases it could potentially bias results (e.g., when the branch length is small and/or many low-quality genomes are included, as in the major mammalian sub-clades).

Simulation studies similar to those I performed in Chapters 2 and 3 could improve our

understanding of the relative potential of different types of sequencing errors to introduce false positives in downstream analyses, but the absolute frequency and pattern of such errors would still difficult to predict without a reliable model for their generation. This is especially true for larger-scale errors from misassembly or misannotation, which are less easily modeled than base calling errors and could have potentially larger negative effects. For estimates of false positives resulting from these types of sequence errors, an empirical approach seems more appropriate.

Two empirical studies in mammals have provided convincing evidence that sequence, alignment and annotation errors can drastically increase the number of false positive PSGs in the branch-site test for positive selection.

Schneider et al. [2009] performed a genome-wide scan for positive selection in the terminal branches of 7 mammalian genomes using the branch-site test and analysed the fraction of PSGs within subsets of high- or low-quality genes according to three sequence and alignment quality metrics. They found that the fraction of PSGs was significantly higher for genes exhibiting lower quality sequence, annotation and alignment metric, with genes in the highest-quality and lowest-quality categories showing a 7.2-fold difference in the inferred fraction of PSGs [Schneider *et al.*, 2009]. This observation provided evidence of a correlation between the chosen quality metrics and the tendency of an alignment to exhibit positive selection. It did not necessarily imply causation, however, as the same result might have been observed—even in the absence of sequence error—if some biological properties of the true PSGs caused them to yield lower quality metrics than non-PSGs. Looking at the three metrics used in their study (sequencing coverage, gene annotation status, and alignment quality according to the heads-or-tails method), it is plausible that properties associated with elevated ω ratios and positive selection, such as recent gene duplication [Beisswanger & Stephan, 2008; Casola & Hahn, 2009; Studer *et al.*, 2008], high GC content [Ratnakumar *et al.*, 2010] or functional shifts [Storz *et al.*, 2008; Wang & Gu, 2001] might have had an error-independent effect resulting in a higher proportion of PSGs in low-scoring categories.

Mallick et al. [2009] took a different approach to the same problem by performing a careful resequencing and reassembly of the chimpanzee genome (the initial assembly of which had lower coverage and lower quality than the human genome) and re-analysing the evidence for positive selection along the chimpanzee lineage in 59 genes which had previously been identified as chimpanzee PSGs. The authors, who were motivated by a concern that previous reports of a larger proportion of PSGs in chimpanzee than in

human [Bakewell *et al.*, 2007] were the result of its lower-quality genome rather than a biologically significant difference in levels of adaptation, found that the vast majority of PSGs identified in two previous studies showed no evidence for positive selection when using their reassembled and higher-coverage version of the chimpanzee genome [Mallick *et al.*, 2009]. This suggested that the original 4x coverage chimpanzee assembly contained a number of sequencing errors leading to false inferences of positive selection. A detailed analysis of 302 codons with multiple spurious non-synonymous substitutions in the original assembly showed roughly comparable effects of sequence error (explaining 23% of codons), assembly error (14% of codons) and local alignment error (30% of codons).

Taken together, the results of Schneider et al. [2009] and Mallick et al. [2009] provide strong evidence in support of the hypothesis that errors in sequencing, assembly, annotation and alignment can result in strongly elevated inferred ω values when using sensitive tests for detecting positive selection. The detailed identification and quantification of error sources performed by Mallick et al. [2009] is especially useful for designing filters to apply to an analysis based largely on low-coverage genomes; in particular, their observation that clusters of chimpanzee-specific mutations were responsible for many false positives motivated the window-based filter I developed and applied here and in the analysis of primate genomes in Chapters ?? and ??.

[Possible figure summarizing the types of error and potential effects on the inference?]

5.2 Preparing mammalian alignments for sitewise analysis

The effects of sequencing, annotation and alignment error on the results of comparative evolutionary analyses can be severe, with a high potential for false positive results when using sensitive evolutionary models [Mallick *et al.*, 2009; Schneider *et al.*, 2009]. In order to minimize the potential for false positive results in this study, sequences were prepared for input to SLR with a series of filters and realignment steps designed to remove low-quality sequences prone to high error rates, realign sequences using an aligner with better performance for detecting sitewise positive selection, and mask out short alignment regions with dubious elevated rates of non-synonymous substitution.

5.2.1 Filtering out low-quality genome sequence

Due to the presence of several low-coverage genome assemblies in the set of available mammalian genomes and the elevated sequencing error rates in such assemblies [Hubbard *et al.*, 2007], I applied a conservative filter to the set of input sequences based on sequence quality scores where available.

Most automated genome assembly pipelines, such as the Arachne tool used to sequence many of the low-coverage mammalian genomes included in Ensembl [Jaffe *et al.*, 2003], output a set of Phred quality scores alongside the genome sequence, with one Phred score per base ranging from 0 to 99. A Phred score represents the probability, calculated by the sequencing and/or assembly program, that a given base call is incorrect. This probability is usually concisely expressed as the negative logarithm of the probability of an error multiplied by ten, or $Q = -10\log_{10}P$, where Q is the Phred score and P is the probability of an incorrect base call [Cock *et al.*, 2010].

Unfortunately, Ensembl does not store quality scores from its source genome assemblies, so Phred quality scores had to be manually downloaded for all genomes with readily available Phred-like quality scores. Most quality scores were provided as a single file in FASTA format with one string of numerical scores per assembled contig. Since the process of filtering a single mammalian coding alignment required collecting scores from many different quality score files for many disjoint genomic locations, a custom script was designed to process each quality score file to allow for quicker score retrieval and better memory performance.

A suitable score threshold for filtering coding regions was chosen based on a study by Hubisz et al. [2011], who performed a detailed analysis of Phred quality scores and observed error rates in low-coverage mammalian genome assemblies by comparing the low-coverage assemblies to matched regions of high-quality sequence from the ENCODE comparative genomics dataset [ENCODE Project Consortium, 2007]. The authors identified a strong correlation between Phred scores and error rates for scores below 25, indicating that the scores were accurate in this range. Error rates did not decrease significantly at scores above 25, however, suggesting that the value of using an extremely high Phred score threshold would be minimal. Furthermore, Hubisz et al. noted that 85% of bases in the low-coverage mammalian genomes contain very high Phred scores (>45) and only 4% have low scores (<20).

Based on these considerations, a threshold Phred score of 25 was chosen as a reasonable trade-off between the potential benefit of avoiding miscalled bases and the potential cost of

masking out correctly sequenced bases. For each coding sequence (CDS) with quality scores available, a “minimum score” approach was used to filter codons: all codons containing one or more nucleotides with a score below 25 were masked out with three ambiguous nucleotides, ’NNN’.

The expected proportion of filtered nucleotides could be calculated from the fraction of bases below the Phred score threshold of 25. According to the cumulative distribution of quality scores found in Hubisz et al. [2011], approximately 5% of bases in low-coverage mammalian genomes contain Phred scores below 25. The worst case scenario, in terms of high-quality bases being masked as a result of using the minimum score, would be if only one base per codon had a score below the threshold. Were that the case, an expected 15% of nucleotides would be filtered, since 3 bases would be masked for every low-quality base. However, the distribution of low-quality bases is likely highly clustered, due to the uneven distribution of repetitiveness and GC content as well as the tendency for uncertain base calls to occur towards the end of sequence reads (all of which are known to affect read coverage and assembly performance, e.g. Teytelman *et al.* [2009]). A more clustered distribution of low-quality bases would cause fewer high-quality bases to become masked by the minimum score approach, reaching the limit of 5% total filtered bases if they always occurred in codon triplets. Thus, anywhere from 5% to 15% of nucleotides from low-coverage genomes were expected to be filtered by this approach.

5.2.2 Removing recent paralogs

To complement the subtree splitting process, which split apart ancient duplications to avoid paralogous comparisons in the sitewise analysis, a second paralog filtering step was applied to remove more recent paralogs. Some of these paralogs might have resulted from gene duplications that occurred subsequent to the two rounds of whole-genome duplication in the vertebrate ancestor, but it was also expected that some proportion of apparent paralogs in the Ensembl gene trees would be the result of errors in gene annotation or in the Compara pipeline.

A particular cause for concern in the current analysis was the possibility that stretches of missing or unassembled sequence in low-coverage genomes might produce gaps of missing data or assembly breakpoints between exons of a single gene, causing it to become annotated as two separate genes. These shortened genes would be treated as independent proteins by the Compara pipeline, likely being placed at very similar positions in the gene tree due to each having been derived from the same single source gene. While this result

might not be detrimental to sitewise analysis in itself (as each shortened gene might be correctly aligned and provide useful information to the alignment), a number of factors, including the low quality of genomic sequence and assembly within these shortened genes, problems with aligning small fractions of a gene against complete sequences, and the potential for incorrect placement of fragmented sequences within the gene tree, made it desirable to remove these shortened genes.

In the case of true recent paralogs, their inclusion in the sitewise analysis might skew the dataset towards increased levels of relaxed constraint or adaptive evolution, as has been hypothesized and observed to occur in recently-duplicated genes [Lynch & Conery, 2000]. Most models of evolution after gene duplication predict that one sequence will retain the ancestral function and diverge less from the common ancestor than the other, so the least-diverged copy would be the natural one to keep as the 'canonical' paralog for evolutionary analysis.

Both gene length and sequence divergence were used to identify which gene among a set of recent paralogs was most suitable to retain for sitewise analysis. It was expected that gene length would help discriminate spuriously shortened genes from true genes, while sequence divergence would distinguish between more and less diverged paralogs. The mean sequence divergence, estimated using the JC69 nucleotide model and the stock Compara gene tree alignments, was calculated between each putative paralog and the rest of the gene tree, and the ratio of the length of each putative paralog to the mean sequence length was also stored. Within each group of putative paralogs, the single gene to keep was chosen by the following rules, applied in order: (1) if only one sequence had a length ratio above 0.5 and all others had a length ratio below 0.5, the longest sequence was kept; (2) if at least one sequence yielded a meaningful (i.e., non-zero and non-infinite) mean distance estimate, the sequence with the lowest distance was kept; (3) if no sequence yielded a meaningful distance estimate (or if all estimated distances were identical), the longest sequence was kept.

Figure 5.1 shows the distribution of gene lengths (relative to the mean across the alignment) for all putative paralogs, kept paralogs, and removed paralogs. The overall distribution of relative lengths shows that most putative paralogs have lengths similar to the alignment mean (with a peak at or slightly above 1), but the shape of the distribution is highly asymmetric with a strong bias towards shorter lengths. The length distribution of the kept paralogs shows that the bulk of highly-shortened genes were successfully removed. If anything, the distribution of kept genes is slightly biased towards lengths greater than

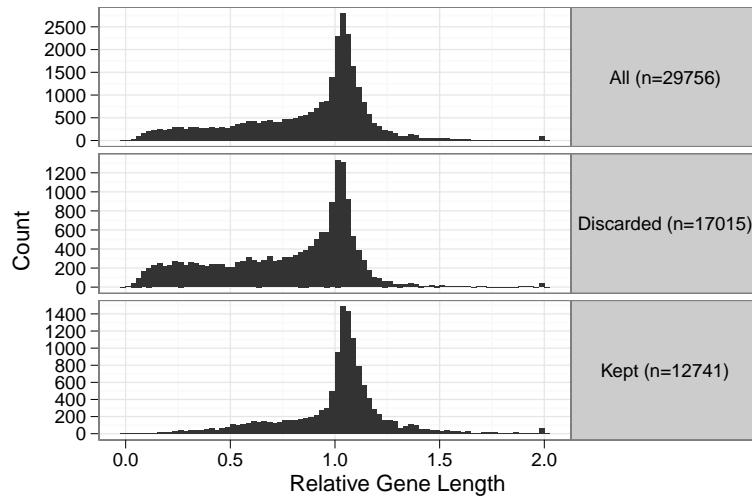


Figure 5.1: Gene lengths of all putative paralogs, normalized to the mean length all sequences in the enclosing tree of each paralog. Putatively paralogous genes (top panel) were either discarded (middle panel) or kept (bottom panel) according to rules based on their length and mean sequence divergence from other aligned sequences.

1, likely due to step 3 of the above process, where the longer gene was kept if putative paralogs yielded identical distance estimates.

An alternate view of the results of the paralog filtering process is shown in Figure 5.2, with the mean divergence of each discarded paralog compared that of the kept paralog, separated into panels according to the reason for discarding that gene. The spread of points above the diagonal in the first panel shows the difference in mean sequence divergence between the kept and discarded putative paralogs where divergence was used to choose between copies (TODO), and the middle panel represents paralogs whose mean divergences were identical (TODO). These two panels together show that although most recent paralogs within this set of gene trees contained indistinguishable levels of sequence divergence, around 40% showed a moderate difference that could be used for selecting the less-diverged copy. The rightmost panel shows the subset of apparent paralogs which were discarded due to their short gene length; a point worth noting here is that there was no bias towards higher or lower divergence levels in the discarded genes (the coefficient of the best-fit linear model for all non-negative values is TODO), suggesting (as expected) that many of the discarded short genes were in fact derived from a single orthologous gene.

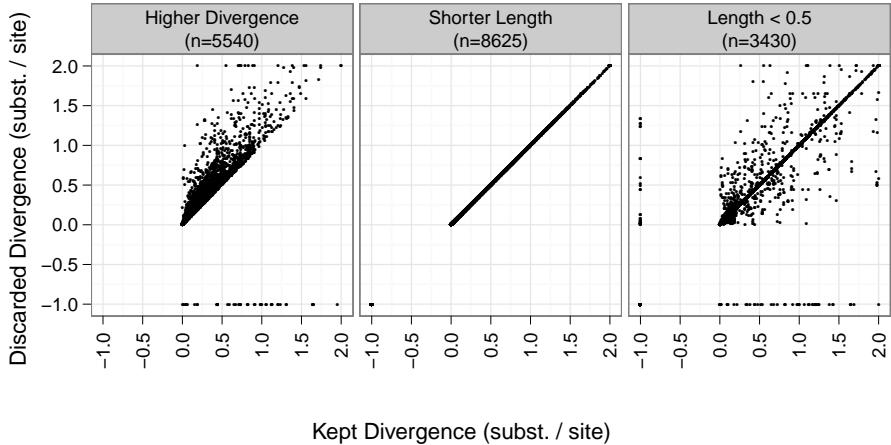


Figure 5.2: Sequence divergence of kept and discarded putative paralogs. Each point represents a gene which was discarded from the tree for one of three reasons: it had more sequence divergence than the kept gene (*Higher Divergence*; left panel), it had equal sequence divergence but shorter length than the kept gene (*Shorter Length*; middle panel), or it had a gene length (relative to the mean across all sequences) of less than 0.5 while the kept copy had a relative length greater than 0.5 (*Length < 0.5*; right panel). Divergence was measured as the mean pairwise divergence between the gene and all other sequences in the tree, and a value of -1 was assigned to genes for which no reliable divergence estimate could be attained due to a lack of sufficient data)

5.2.3 Realigning coding sequences

After filtering out codons with low quality scores and removing putative paralogs, sequences were aligned with PRANK [Löytynoja & Goldman, 2008] using its codon alignment model based on the empirical codon model [Kosiol *et al.*, 2007]. The simulation experiments described in Chapters 2 and 3 as well as numerous previously-reported empirical and simulation-based studies have shown PRANK’s codon-based alignments to be superior for avoiding false positives in the detection of sitewise positive selection, strongly supporting the choice of PRANK for this analysis.

5.2.4 Filtering out clusters of non-synonymous substitutions

A final filtering step was evaluated for possible application to PRANK-aligned sequences in order to ensure that stretches of aligned but nonhomologous sequence, resulting either from misalignment or from exon mis-annotation artifacts, were not causing elevated rates of non-synonymous substitutions within specific regions of genes. This filtering step motivated

by the expectation that errors resulting from either misalignment or exon mis-annotation would both lead to clustered regions of nonhomologous aligned nucleotides and, correspondingly, clustered regions of elevated non-synonymous substitution rates. Clustering of non-synonymous substitutions was expected because in both cases, the existence of one misaligned column is not independent of nearby columns: for mis-annotation the misalignment would span the length of the erroneously-annotated exon, while the global nature of the progressive pairwise alignments performed by PRANK (and all other alignment algorithms) causes any misalignment error to be strongly non-independent with respect to errors in nearby alignment columns.

These arguments provided some hope that clustered non-synonymous substitutions could be used as a signal to detect potential misalignment and mis-annotated regions, but the utility of such a signal for filtering alignments must depend on the strength of error-caused non-synonymous clusters relative to both the frequency of true non-synonymous substitutions and their tendency to cluster together along the length of the protein sequence. Sequences separated by longer branch lengths will clearly show higher densities of true non-synonymous substitutions, possibly drowning out the error-caused signal and reducing the ability to use substitution clusters as a discriminating factor. Furthermore, non-synonymous substitutions were shown to be significantly more clustered than expected by chance in a number of genomic analyses of mammals and insects [Bazykin *et al.*, 2004; Callahan *et al.*, 2011; ?], causing some concern that a filter based on detecting clusters of non-synonymous substitutions might attenuate the signal of true adaptive substitution that was one target of the present study.

[Write up the empirical investigation of non-synonymous substitution clustering.]

5.3 Genome-wide analysis of sitewise selective pressures in mammals

5.3.1 Mammalian species subsets for sitewise analysis

The SLR method was applied sequentially to several species subsets of each alignment of mammalian orthologs. For each subset, sequences corresponding to species within the subset were extracted from the alignment along with the corresponding subtree and input to SLR. If fewer than two sequences were available for a given subset, that subset was skipped and its absence from the dataset was recorded. Eight subsets in total were selected

for analysis; the species included in each subset and some phylogenetic measures of each subset are listed in Table 5.1.

Three subsets (Glires, Primates, and Laurasiatheria) were chosen because they represent the three mammalian superorders with the greatest taxonomic representation in Ensembl, providing an opportunity to compare the molecular evolutionary dynamics of three monophyletic mammalian groups containing varying levels of divergence, diverse biological characteristics, and a number of high-quality reference genomes. A fourth parallel mammalian subclade, named Atlantogenata and consisting of sloth, armadillo, tenrec, elephant and hyrax, was also included, but the monophyly of this group is still debated [Churakov *et al.*, 2009; Murphy *et al.*, 2007] and it contains only one high-coverage genome. As such, it was not considered a primary target for the mammalian superorder analysis.

Two larger species sets, Eutheria and Mammalia, were chosen for the purpose of measuring average sitewise selective pressures with high precision across a wider group of mammals. Using the Ensembl species tree as a guide, the estimated total synonymous branch lengths spanned by Ensembl species within Eutheria and Mammalia were 4.95 and 6.18, respectively. Simulations performed by Anisimova and Yang [2011; ?] and by myself in Chapter 2 predicted that the greater amount of branch length in the Eutherian and Mammalian trees—with two to three times the value of 1.71 for Laurasiatheria, the superorder with the largest total branch length—would result in significantly higher levels of power and accuracy for estimating sitewise ω and detecting sitewise positive selection. In this respect, Mammalia and Eutheria were more similar to each other than to any of the superorders.

However, the Mammalia and Eutheria subsets differed markedly in a different (and largely orthogonal) phylogenetic factor, the evolutionary depths of their last common ancestors. Whereas the ancestor of all Eutherian mammals lived ca. [125] mya, the Mammalian ancestor dates back to [320] mya. This suggested that a comparison between the sitewise results for the two groups might provide useful insight into the general effect of adding longer, deeper branches to a sitewise evolutionary analysis as well as some indirect evidence on the molecular evolutionary dynamics of our most distant mammalian relatives (the Eutheria and Mammalia groups only differ by the inclusion of wallaby, opossum and platypus in the Mammalia group).

Quantitatively, as measured by the MPL from the Ensembl species tree, the Eutheria subset ($MPL = 0.24$) is far more similar to either of the three superorders (MPL from 0.13 to 0.27) than to the Mammalia subset ($MPL = 0.54$). This is due to the striking

adaptive radiation of Eutherian mammals [Archibald, 1999; Bininda-Emonds *et al.*, 2007], which caused a quick succession of speciation events around the K-T boundary and gives a largely star-like structure to the eutherian evolutionary tree. Interestingly, according to the time-resolved mammalian tree from Bininda-Edmonds *et al.* [2007] the Diprotodontia order (containing wallaby and opossum, two outgroups to the Eutheria clade) experienced a radiation similar to, but less pronounced than, the Eutherian radiation; a comparison of the evolution of the deeply-rooted Diprotodontia clade to its sister Eutherian clade would be very enlightening, but the species representation of Diprotodontia (currently at one high-coverage and one low-coverage genome) is too limited to allow for a powerful analysis. Nevertheless, the inclusion of the three non-Eutherian species in the Mammalian species group was expected to provide an additional data point for aiding in an understanding of the complex relationship between branch length, power and biological variability in the analysis of sitewise selective pressures.

Finally, to further investigate the combined impact of evolutionary depth, biological variability and branch length on the results of large-scale sitewise analyses, two “sparse” subsets were created to act as controls relative to two existing species subsets. The species in the Sparse Glires group were chosen to approximate the total branch length of the Primate clade with species from the Glires clade, while the Sparse Mammals subset was constructed by selecting one species (preferably with a high-coverage genome) from each major mammalian branch, greatly reducing the total branch length covered but maintaining a similar evolutionary depth and distribution of branches within the species tree. The branch lengths in Table 5.1 show that the Sparse Glires group was only somewhat successful in its goal of approximating the Primates branch length (with total branch lengths of 0.99 and 0.68, respectively) while the Sparse Mammals group achieved a threefold lower total branch length compared to the full Mammalia group while maintaining a nearly identical MPL.

5.3.2 Evaluation of the bulk distributions and the design of a filtering approach

Sitewise data were collected from SLR and stored in a database for storage and further analysis. The Mammalia subset, containing the most branch length of all the datasets and representing the entire set of aligned sequences, and the Primate subset, containing the lowest overall branch length, were used to perform quality-control checks on the sitewise data. The point of these checks were to evaluate whether any additional filtering of the

| Name | (Species Count) | Species List | Ensembl | | Gene Median | |
|-----------------|-----------------|--|---------|------|-------------|------|
| | | | N_E | MPL | Total | MPL |
| Primates | (10) | Bushbaby, Chimpanzee, Gibbon, Gorilla, Human, Macaque, Marmoset, Mouse Lemur, Orangutan, Tarsier | 20000 | 0.13 | 0.68 | 0.16 |
| Glires | (7) | Guinea Pig, Kangaroo rat, Mouse, Pika, Rabbit, Rat, Squirrel | 230000 | 0.27 | 1.44 | 0.40 |
| Laurasiatheria | (12) | Alpaca, Cat, Cow, Dog, Dolphin, Hedgehog, Horse, Megabat, Microbat, Panda, Pig, Shrew | 34410 | 0.19 | 1.71 | 0.26 |
| Atlantogenata | (5) | Armadillo, Elephant, Hyrax, Sloth, Tenrec | 30000 | 0.20 | 0.83 | 0.26 |
| Eutheria | (35) | Alpaca, Armadillo, Bushbaby, Cat, Chimpanzee, Cow, Dog, Dolphin, Elephant, Gibbon, Gorilla, Guinea Pig, Hedgehog, Horse, Human, Hyrax, Kangaroo rat, Macaque, Marmoset, Megabat, Microbat, Mouse, Mouse Lemur, Opossum, Orangutan, Panda, Pig, Pika, Platypus, Rabbit, Rat, Shrew, Sloth, Squirrel, Tarsier, Tenrec, Tree Shrew, Wallaby | 110000 | 0.24 | 4.95 | 0.35 |
| Mammalia | (38) | Alpaca, Armadillo, Bushbaby, Cat, Chimpanzee, Cow, Dog, Dolphin, Elephant, Gibbon, Gorilla, Guinea Pig, Hedgehog, Horse, Human, Hyrax, Kangaroo rat, Macaque, Marmoset, Megabat, Microbat, Mouse, Mouse Lemur, Opossum, Orangutan, Panda, Pig, Pika, Platypus, Rabbit, Rat, Shrew, Sloth, Squirrel, Tarsier, Tenrec, Tree Shrew, Wallaby | 120000 | 0.54 | 6.18 | 0.66 |
| Sparse Glires | (5) | Guinea Pig, Kangaroo rat, Mouse, Rat, Squirrel | 230000 | 0.25 | 0.99 | 0.36 |
| Sparse Mammalia | (7) | Armadillo, Dog, Elephant, Human, Mouse, Platypus, Wallaby | 120000 | 0.51 | 2.18 | 0.60 |
| | | | | | | 2.80 |

Table 5.1: Species subsets used for sitewise analysis. Values under the “Ensembl” heading were calculated from subsets of the species tree used for evolutionary analyses by the Ensembl Compara pipeline, while values under the “Gene Median” heading were calculated as median values across the 15,XYZ gene trees analyzed (with branch lengths optimized by SLR). MPL – mean path length, Total – total branch length.

sitewise results was necessary before characterizing the global distribution of constraint in this and the other species subsets. Even if the sequence and alignment filters described above were successful at reducing the number of false positives due to incorrect input alignments, the behavior of SLR when applied to large datasets of heterogeneous alignments has not been well-studied, and a number of biases might have influenced the results. A particular point of concern was that columns with different patterns of gap and non-gap sequences, especially those with few non-gap sequences, might yield different performance characteristics. Although the SLR method was sensibly designed to account for uncertainty in the estimation of ω and detection of positive selection, one might reasonably expect less-desirable statistical properties from sites with 2 non-gap sequences compared to sites with 20.

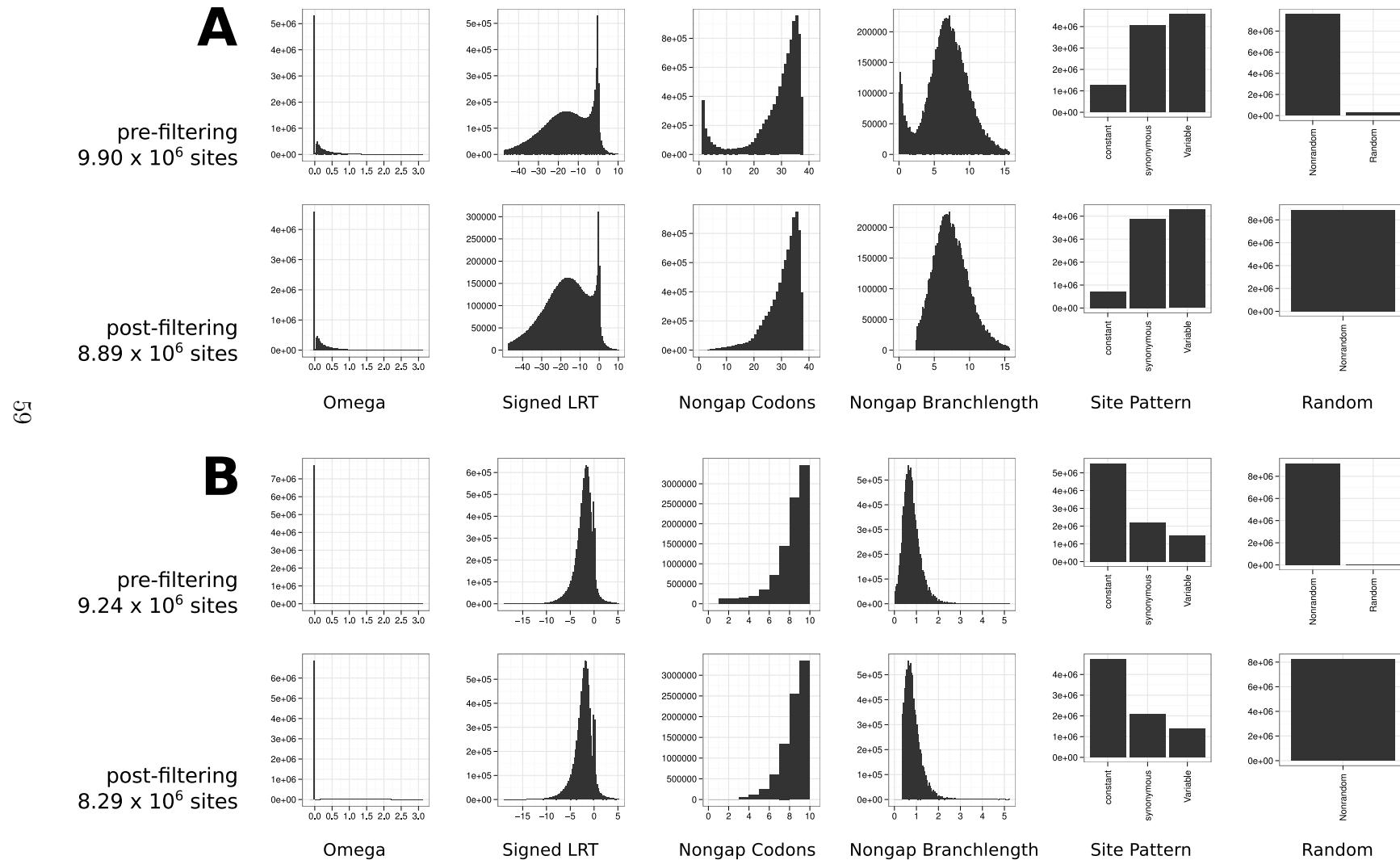


Figure 5.3: Distributions of sitewise values for Mammalia (A) and (B) Primates, before (top row) and after (bottom row) removing sites based on the filtering scheme (see text). Note: the y-axis scale varies between rows, and the x-axis scale varies between (A) and (B) for the Signed LRT, Nongap Codons and Nongap Branch Length values.

Figure 5.3 shows the distributions of six sitewise values: two continuous values output by SLR (Omega and Signed LRT), two categorical values from SLR (Site Pattern and Random), and two values calculated from the codon alignment (Nongap Codons and Non-gap Branch Length). The Nongap Codons value measures the number of non-gap codons in the given alignment column, while the Nongap Branch Length represents the total branch length connecting all non-gap sequences (using the gene tree with branch lengths optimized by SLR).

A prominent feature of the distribution of ω values for the unfiltered Mammalian data, shown in the top panel of Figure 5.3A, was the large number of sites with a maximum-likelihood estimated ω of zero. Further inspection of the data revealed that all zero- ω sites contained either synonymous or constant site patterns, and all sites with constant patterns (and nearly all sites with synonymous patterns) yielded a maximum likelihood ω estimate of zero. An estimate of zero for synonymous sites is intuitively appropriate, as the lack of any non-synonymous substitutions throughout the tree would seem to provide no evidence for a non-synonymous substitution rate of greater than zero. For constant sites the case is less clear, because no data regarding the rate of either synonymous or non-synonymous substitutions exists. However, given SLR's assumption of a constant synonymous substitution rate throughout each gene [Massingham & Goldman, 2005], the ω value which maximizes the likelihood of observing zero substitutions is zero, since that value minimizes the non-synonymous (and total) substitution rate.

It is interesting to note that a small proportion (ca. 0.2%) of synonymous sites resulted in maximum likelihood estimates greater than zero. Manual investigation of a number of these sites showed them all to include synonymous codons with multiple nucleotide differences (such as those coding for serine and arginine), for which SLR's mechanistic codon model—which does not allow for multiple simultaneous nucleotide changes—required the inference of multiple non-synonymous substitutions and, thus, a non-synonymous substitution rate of greater than zero. The existence of multiply-substituted codons in alignments has been previously reported [Averof *et al.*, 2000; Whelan & Goldman, 2004], and empirical results have supported the notion that codon models that allow for multiple simultaneous nucleotide changes better describe evolution than those that do not [Kosiol *et al.*, 2007]. However, the very low proportion of synonymous sites requiring nonzero non-synonymous substitution rates suggested that the impact of these effects on the current dataset was minimal; this is likely due to the relatively short branch lengths separating the nodes of the mammalian tree, making it less likely that codons with multiple substitutions (whether

Use the sit
and 'seq' t
do a more
tive analys

the result of simultaneous multiple nucleotide changes or successive single changes) would be observed [Kosiol *et al.*, 2007].

The distributions of Nongap Codons and Nongap Branch Length values in the top row of Figure 5.3A showed that many alignment columns contained only a few non-gap sequences. Both distributions were bimodal with a larger peak at the upper end of the range and a smaller peak at the lower end of the range. If the sites with low non-gap codon counts represented accurate evolutionary histories, the observed excess of mostly gapped sites might be taken as an indication that insertion events in terminal lineages or recent ancestral lineages were prominent enough in mammalian evolution to leave a noticeable signature of sites with low non-gap codon counts. This would be a very interesting observation, but unfortunately it is not likely a correct one. Given the many possible sources of error in the creation of Ensembl gene trees, however, a more likely scenario was that sites with low codon counts and low branch lengths came from stretches of sequence which only exist in a few species as a result of annotation or alignment error, with a higher probability of being nonhomologous and showing spurious signals of positive selection. This would make such sites prime candidates for filtering out prior to a large-scale analysis.

To test the hypothesis that sites with few non-gap sequences are less reliable than other sites, the Mammals and Primates data were split into deciles by nongap branch length and sites within each decile were summarized by the proportion of sites showing evidence for purifying and positive selection; the results of this analysis are presented in Table 5.2. The lowest decile appeared to be a clear outlier in the Mammalia dataset, with nearly 17% of sites having an estimated ω of greater than 1 and 2% of sites showing significant evidence for positive selection at a nominal 5% error rate. Deciles with greater nongap branch lengths showed lower proportions of sites with $\omega > 1$ and less evidence for positive selection, with a gradual increase in both values at progressively higher deciles. The gradual increase in evidence for positive selection with increasing nongap branch length could be explained by genes with higher overall dN/dS ratios (and perhaps more positive selection) producing, on average, higher estimated branch lengths due to the increased non-synonymous substitution rate. Overall, these patterns were consistent with the expectation that sites with few non-gap sequences were not consistent with the bulk of the dataset, and Table 5.2 showed that removing sites with the lowest 10% of nongap branch length would remove most of the apparently anomalous sites.

The breakdown of Primates data in Table 5.2 showed a trend similar the Mammalia dataset, although the distinction between the lowest decile and the rest of the dataset was

| | BL | Nongap BL | | | Nongap Codons | | | ω_{ML} , % | | |
|----------|------|-----------|-------|-------|---------------|-------|-----|-------------------|-------|------|
| | | Quantile | 25% | 50% | 75% | 25ff% | 50% | 75% | < 1 | > 1 |
| Mammalia | 0.10 | 0.31 | 0.74 | 1.46 | 2 | 3 | 6 | 81.87 | 18.13 | 2.09 |
| | 0.20 | 3.28 | 3.78 | 4.17 | 19 | 30 | 35 | 95.01 | 4.99 | 0.52 |
| | 0.30 | 4.77 | 5.02 | 5.24 | 27 | 33 | 36 | 96.90 | 3.10 | 0.35 |
| | 0.40 | 5.67 | 5.86 | 6.04 | 28 | 33 | 36 | 96.94 | 3.06 | 0.35 |
| | 0.50 | 6.38 | 6.55 | 6.72 | 29 | 33 | 36 | 96.75 | 3.24 | 0.39 |
| | 0.60 | 7.06 | 7.22 | 7.40 | 29 | 33 | 36 | 96.41 | 3.59 | 0.44 |
| | 0.70 | 7.77 | 7.95 | 8.15 | 29 | 33 | 36 | 96.04 | 3.96 | 0.50 |
| | 0.80 | 8.58 | 8.79 | 9.03 | 29 | 33 | 35 | 95.43 | 4.57 | 0.61 |
| | 0.90 | 9.58 | 9.88 | 10.24 | 29 | 33 | 35 | 94.57 | 5.43 | 0.79 |
| | 1.00 | 11.22 | 12.00 | 13.28 | 29 | 32 | 35 | 92.95 | 7.04 | 1.14 |
| Primates | 0.10 | 0.17 | 0.25 | 0.30 | 4 | 6 | 8 | 94.42 | 5.58 | 0.61 |
| | 0.20 | 0.38 | 0.41 | 0.44 | 8 | 9 | 10 | 94.39 | 5.61 | 0.32 |
| | 0.30 | 0.49 | 0.52 | 0.54 | 8 | 9 | 10 | 93.64 | 6.36 | 0.30 |
| | 0.40 | 0.59 | 0.61 | 0.63 | 8 | 9 | 10 | 93.09 | 6.91 | 0.33 |
| | 0.50 | 0.67 | 0.69 | 0.71 | 8 | 9 | 10 | 92.39 | 7.61 | 0.35 |
| | 0.60 | 0.76 | 0.78 | 0.80 | 8 | 9 | 10 | 91.29 | 8.71 | 0.46 |
| | 0.70 | 0.85 | 0.87 | 0.90 | 8 | 9 | 10 | 90.68 | 9.32 | 0.50 |
| | 0.80 | 0.97 | 1.00 | 1.04 | 8 | 9 | 10 | 89.10 | 10.90 | 0.66 |
| | 0.90 | 1.13 | 1.19 | 1.25 | 9 | 9 | 10 | 87.13 | 12.87 | 0.86 |
| | 1.00 | 1.44 | 1.61 | 1.95 | 8 | 9 | 10 | 84.64 | 15.36 | 1.24 |

Table 5.2: Proportions of sites with evidence for purifying and positive selection in the Mammalia and Primates datasets broken down by nongap branch length. Sites were separated into 10 equally-sized bins of nongap branch length and summarized by the 25th, 50th and 75th percentiles of nongap branch length and nongap codons, the percentage of sites with ω estimated below or above 1, and the percentage of sites with significant evidence of positive selection at a nominal 5% FPR.

less clear. The F_{pos} in the lowest decile was only slightly higher than in the next-highest decile, and $F_{<1}$ was lower than in all other bins. The gradual increase of $F_{>1}$ and F_{pos} in higher branch length deciles was similar to that seen in the Mammalia dataset, however. Despite weaker evidence in the Primates data for the erroneous nature of sites with few non-gap sequences, it still appeared that filtering sites in the bottom decile would improve the overall quality and consistency of the data.

Turning back to the bulk distributions in Figure 5.3, the rightmost panel depicts the small set of sites designated as “random” by SLR. These sites were flagged by SLR as having a site pattern not significantly different from random [Massingham & Goldman, 2005], and they were also targeted for removal before analysis of the global distribution.

The final filtering protocol applied to each sitewise dataset included three steps. First, all sites within the bottom 10% of nongap branch length values were removed; second, sites

flagged by SLR as “random” were removed; third, all sites with fewer than four non-gap sequences were removed.

The most prominent effect of the filter on the bulk distributions in Figures 5.3A and 5.3B was, as expected, the removal of the excess of sites with low non-gap branch lengths and non-gap codon counts. The distribution of ω estimates and Signed LRT statistics were largely unchanged, indicating that the overall characteristics of each dataset were not significantly altered by the filter. The lack of large-scale change was a somewhat reassuring result, given that the filter only removed roughly 10% of sites from each dataset.

A more detailed comparison of various summary statistics for the filtered and unfiltered datasets showed that filtering had a noticeable impact on three quantities of interest: it reduced the proportion of constant sites, lowered the mean ω , and slightly decreased the percentage of positively-selected sites. Tables 5.3 and 5.4 contain summary statistics and calculations performed on the filtered and unfiltered Mammalia and Primates data. Most of the data contained in these tables will be more fully described in the next subsection, but a comparison of the filtered and unfiltered rows for Primates and Mammalia provided a means by which to quantitatively assess the effect of filtering on particular aspects of the dataset. First, the percentage of constant sites was reduced in the post-filtering data, moving from 60.04% to 57.62% in Primates and from 12.62% to 8.17% in Mammalia. This was expected, as the sites removed by filtering were enriched in constant site patterns due to their lower non-gap branch lengths. Second, the mean ω value was slightly reduced in Primates (from 0.32 to 0.27) and greatly reduced in Mammalia (from 0.49 to 0.20), likely due to the removal of sites containing a small number of nonhomologous codons which might have produced abnormally high sitewise maximum likelihood ω estimates. Third, the proportion of positively-selected sites (shown for a range of significance thresholds in Table 5.4) was moderately reduced in both Primates (from 0.56% to 0.53% at $P_{\chi^2_1} < 0.05$) and Mammalia (from 0.72% to 0.56%). These three effects of filtering, each showing a shift indicating more useful data (e.g., a lower percentage of constant sites) or less evidence for positive selection (e.g., lower mean ω and proportion of positively-selected sites) in the post-filtering datasets, together provided evidence supporting the inclusion of such a filtering step prior to the analysis and comparison of sitewise estimates in different species sets. Although most quantities of interest were not noticeably changed, those that were affected by filtering shifted to more conservative values, which was taken to be a positive sign given the persistent concern regarding the presence of false positives in detecting positive selection.

| Name | Sites | Site Pattern, % | | | Codons | Nongap BL | | | ω_{ML} | | ω_{ML} Above / Below, % | | | |
|-----------------|----------|-----------------|-------|-------|--------|-----------|------|------|---------------|------|--------------------------------|-------|------|-------|
| | | Const. | Syn. | Nsyn. | | Med. | Mean | SD | Mean | SD | < 0.5 | < 1 | > 1 | > 1.5 |
| Primates (raw) | 9.24e+06 | 60.04 | 24.08 | 15.88 | 9 | 0.73 | 0.85 | 1.36 | 0.19 | 0.50 | 86.21 | 91.08 | 8.92 | 5.79 |
| Primates | 8.24e+06 | 57.62 | 25.49 | 16.89 | 9 | 0.78 | 0.92 | 1.24 | 0.20 | 0.51 | 85.45 | 90.86 | 9.14 | 5.72 |
| Glires | 7.70e+06 | 33.99 | 41.95 | 24.06 | 7 | 1.81 | 2.00 | 1.24 | 0.15 | 0.35 | 89.33 | 96.01 | 3.99 | 1.76 |
| Laurasiatheria | 8.16e+06 | 32.52 | 38.93 | 28.54 | 10 | 2.03 | 2.24 | 1.32 | 0.19 | 0.41 | 86.46 | 94.08 | 5.92 | 2.95 |
| Atlantogenata | 5.15e+06 | 55.95 | 29.75 | 14.30 | 5 | 0.98 | 1.09 | 0.62 | 0.13 | 0.38 | 89.03 | 94.91 | 5.09 | 2.55 |
| Eutheria | 8.79e+06 | 11.54 | 44.49 | 43.97 | 30 | 5.83 | 6.17 | 2.76 | 0.19 | 0.36 | 87.01 | 95.27 | 4.73 | 1.91 |
| Mammalia (raw) | 9.90e+06 | 12.62 | 40.96 | 46.43 | 32 | 6.89 | 6.88 | 3.74 | 0.21 | 0.40 | 86.26 | 94.29 | 5.71 | 2.95 |
| Mammalia | 8.89e+06 | 8.17 | 43.56 | 48.26 | 33 | 7.23 | 7.55 | 3.32 | 0.19 | 0.35 | 87.52 | 95.74 | 4.26 | 1.63 |
| Sparse Glires | 6.72e+06 | 44.22 | 38.40 | 17.38 | 5 | 1.28 | 1.43 | 1.14 | 0.12 | 0.34 | 90.07 | 96.35 | 3.65 | 1.69 |
| Sparse Mammalia | 7.52e+06 | 26.32 | 44.99 | 28.69 | 6 | 2.55 | 2.82 | 1.81 | 0.14 | 0.32 | 90.69 | 96.85 | 3.15 | 1.33 |

Table 5.3: Summary statistics and maximum likelihood ω estimates for sitewise data in eight species groups. Rows labeled “Primates (raw)” and “Mammalia (raw)” were computed based on unfiltered data and are included for reference. Columns under the “ ω_{ML} Above / Below, %” heading measure the cumulative percentage of sites with ω_{ML} below or above the indicated value. Med.—median, Const.—constant, Syn.—synonymous, Nsyn.—non-synonymous, ML—maximum likelihood.

| Name | Positively Selected Sites (%) | | | | $P_{\chi_1^2} < 0.1, \%$ | | | $P_{\chi_1^2} < 0.05, \%$ | | |
|-----------------|-------------------------------|-----------------------|-----------------------|----------------|--------------------------|-------|-------|---------------------------|-------|-------|
| | $P_{\chi_1^2} < 0.1$ | $P_{\chi_1^2} < 0.05$ | $P_{\chi_1^2} < 0.01$ | FDR<0.05 | Pos | Neg | Neut. | Pos | Neg | Neut. |
| Primates (raw) | 89734 (0.97) | 52005 (0.56) | 15973 (0.17) | 357 (0.00) | 0.97 | 29.86 | 69.16 | 0.56 | 14.31 | 85.13 |
| Primates | 75629 (0.92) | 43382 (0.53) | 12985 (0.16) | 158 (0.00) | 0.92 | 33.17 | 65.92 | 0.53 | 15.97 | 83.50 |
| Glires | 17388 (0.23) | 8735 (0.11) | 1953 (0.03) | 0 (0.00) | 0.23 | 76.73 | 23.04 | 0.11 | 65.58 | 34.31 |
| Laurasiatheria | 57494 (0.70) | 34693 (0.43) | 11589 (0.14) | 471 (0.01) | 0.70 | 74.49 | 24.81 | 0.43 | 64.65 | 34.92 |
| Atlantogenata | 11987 (0.23) | 5766 (0.11) | 1166 (0.02) | 0 (0.00) | 0.23 | 46.82 | 52.94 | 0.11 | 24.24 | 75.64 |
| Eutheria | 80406 (0.91) | 55063 (0.63) | 24804 (0.28) | 9378 (0.11) | 0.91 | 85.09 | 14.00 | 0.63 | 82.10 | 17.27 |
| Mammalia (raw) | 107595 (1.09) | 71167 (0.72) | 28894 (0.29) | 7822 (0.08) | 1.09 | 80.15 | 18.77 | 0.72 | 76.98 | 22.30 |
| Mammalia | 72278 (0.81) | 49311 (0.55) | 22123 (0.25) | 8335 (0.09) | 0.81 | 86.51 | 12.68 | 0.55 | 83.87 | 15.57 |
| Sparse Glires | 10582 (0.16) | 4928 (0.07) | 933 (0.01) | 0 (0.00) | 0.16 | 69.33 | 30.51 | 0.07 | 48.67 | 51.26 |
| Sparse Mammalia | 12737 (0.17) | 6347 (0.08) | 1396 (0.02) | 0 (0.00) | 0.17 | 79.85 | 19.98 | 0.08 | 72.43 | 27.49 |

Table 5.4: Proportions of sites subject to positive, purifying and neutral selection at various LRT_{SLR} thresholds. The Benjamini-Hochberg method [Benjamini & Hochberg, 1995] was used to identify the LRT_{SLR} threshold at which FDR<0.05. For columns under the headings “ $P_{\chi_1^2} < 0.1, \%$ ” and “ $P_{\chi_1^2} < 0.05, \%$ ”, Pos. and Neg. are the percentage of sites with significant evidence for positive and negative selection, respectively, and Neut. is the percentage of “neutral” sites not showing significant evidence for non-neutral selection.

5.3.3 The global distribution of sitewise selective pressures in mammals

Each set of sitewise data was filtered as described above. The resulting global distributions of site patterns, sitewise ω estimates, and 95% confidence intervals are shown in Figure 5.4; the left panel in each row plots the number of sites with constant, synonymous, and non-synonymous patterns. All sites with $\omega_{ML} = 0$ had constant or synonymous patterns, and all sites with $\omega_{ML} > 0$ had non-synonymous patterns; the distributions of ω_{ML} for these non-synonymous sites are shown as histograms on the right panel in each row.

5.3.3.1 Site patterns and ω_{ML} values reveal the prevalence of purifying selection in mammalian proteins

The site pattern counts in Figure 5.4 showed that the branch length of each species group had a strong effect on the overall composition of the sitewise data. Species groups covering little branch length, such as Primates and Atlantogenata, contained a majority of constant sites, while groups comprising lots of branch length, such as Eutheria and Mammalia, contained few constant sites and roughly equal proportions of synonymous and non-synonymous sites. Comparing the Glires and Mammalia data with their corresponding “sparse” datasets confirmed that this trend was largely due to branch length as opposed to biological factors: the Sparse Glires data yielded a smaller proportion of non-synonymous sites and a greater proportion of constant sites than the Glires data (17.41% versus 24.08% for non-synonymous sites, 44.21% versus 33.98% for constant sites; numbers from Table 5.3), and the pattern for Sparse Mammalia and Mammalia was qualitatively the same.

Turning to the distribution of these ω_{ML} estimates, shown in Figure 5.4 as a series of histograms representing the ω_{ML} density (for nonzero values only) and a series of solid lines representing the cumulative density (for all values), it is clear that the vast majority of protein-coding sites have evolved under purifying selection in mammals. The Mammalia group, which contained a very small proportion of potentially uninformative constant sites (8.17%), showed a maximum density of nonzero ω_{ML} estimates at $\omega \approx 0.1$, and the vast majority of sites showed some evidence of purifying selection, with ω_{ML} estimates below 1. The height of the ω_{ML} cumulative distribution at $\omega = 1$ corresponds to the proportion of such sites; the exact value, included in Table 5.3 under the “ < 1 ” column, is 95.74%. The nonzero ω_{ML} values were more evenly spread in the other species groups, with Glires showing a maximum nonzero ω_{ML} density at around $\omega \approx 0.25$ and Primates at $\omega \approx 0.7$.

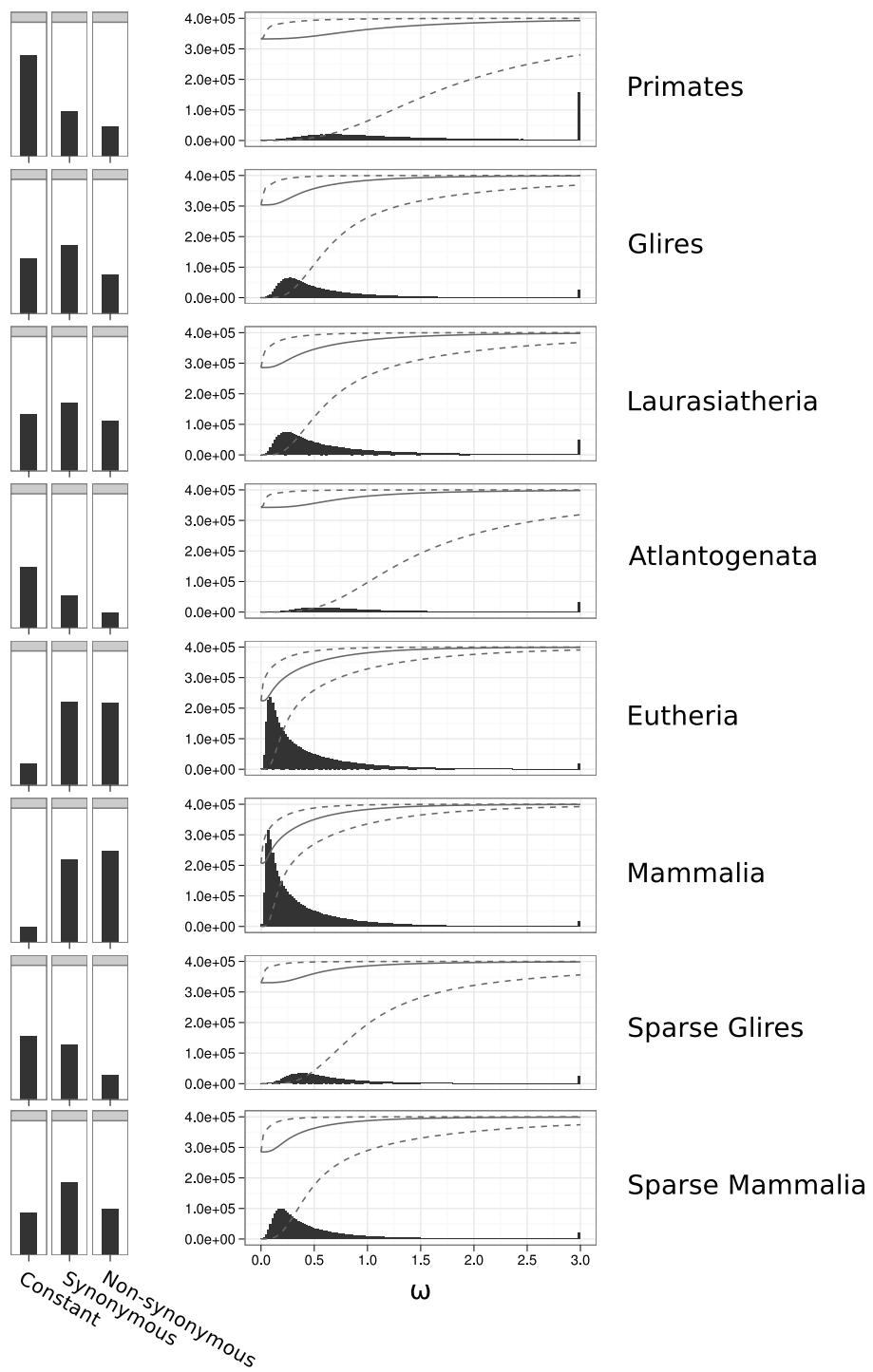


Figure 5.4: Global distributions of site patterns and ω estimates for eight species groups. Left, bars represent the number of sites showing constant, synonymous, and non-synonymous patterns. Note, the y-axis is held constant between rows. Right, bars represent a histogram of maximum likelihood ω estimates where $\omega > 0$. Sites where $\omega > 3$ are counted in the bin at $\omega = 3$. A solid line is drawn showing the cumulative proportion of sites with ω below the current value, and dashed lines are drawn above and below the solid line, showing the cumulative proportion of sites with the lower or upper range, respectively, of their 95% confidence interval below the current value.

This upwards shift in nonzero ω_{ML} estimates relative to Mammalia was likely due to the greater proportion of constant and synonymous sites in lower-branch length datasets: sites which were truly evolving with $\omega > 0$, but where no non-synonymous or synonymous substitutions were observed, would have their ω_{ML} estimate “pushed” towards zero, causing an increase in constant sites and a concomitant upwards shift in the distribution of the remaining nonzero ω_{ML} values.

5.3.3.2 Sitewise confidence intervals and LRT statistics identify sites with significant evidence for purifying and positive selection

An important component of SLR’s output is the set of statistics providing information about the confidence with which purifying or positive selection was detected. These values include the lower and upper bounds of $CI_{95\%}$, the 95% confidence interval for each ω_{ML} estimate, and the LRT statistic, which corresponds to the strength of evidence for purifying or positive selection. Following Massingham [2005], I used a signed version of the LRT statistic (hereafter LRT_{SLR}), formed by negating the LRT statistic for sites where $\omega_{ML} < 1$, as a way to sort sites according to their evidence, or lack thereof, for purifying and positive selection. Thus, sites with $LRT_{SLR} < 0$ showed at least some evidence for purifying selection and sites with $LRT_{SLR} > 0$ showed at least some evidence for positive selection. It should be noted that the LRT_{SLR} is a measure of the strength of evidence for purifying or positive selection, but not necessarily the actual strength of that selection. For example, an alignment covering a very large branch length might yield a strongly negative LRT_{SLR} for a site with ω_{ML} only moderately below 1, because the difference between dN and dS was highly statistically significant; on the other hand, a strongly-purifying site in an alignment covering less branch length might produce a much less-negative LRT_{SLR} , even with ω_{ML} near zero.

Figure 5.5A shows the empirical relationship between LRT_{SLR} , ω_{ML} and the $CI_{95\%}$ width for sites from the Mammali dataset. The left panel, comparing the LRT_{SLR} to nonzero ω_{ML} estimates, shows that the two values are highly correlated, with the greatest number of low ω_{ML} estimates occurring at sites with strongly negative LRT_{SLRs} . Correspondingly, the middle panel shows an even stronger relationship between the LRT_{SLR} magnitude and the $CI_{95\%}$ width, with the tighest windows at sites with very strong evidence for purifying selection. The rightmost panel compares the ω_{ML} of each site with the width of its $CI_{95\%}$, revealing a more linear and diffuse positive relationship between ω_{ML} and the size of the $CI_{95\%}$. The equivalent plots for Primates, shown in Figure 5.5B, reveal

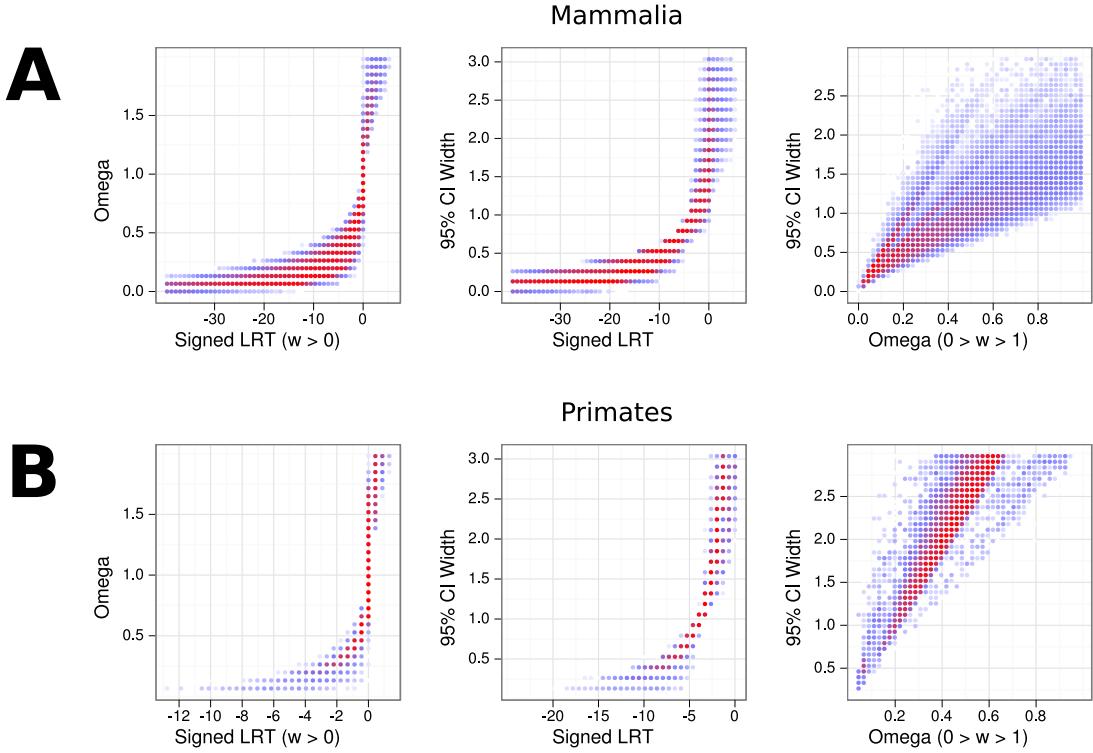


Figure 5.5: The relationship between LRT_{SLR} , ω_{ML} , and $CI_{95\%}$ width in (a) Mammalia and (b) Primates datasets. Each point represents the binned density of sites; no points are drawn where no density exists, while blue and red points are drawn at areas of low and high density, respectively. The left panel shows sites where $\omega_{ML} < 0$, the middle panel shows all sites, and the right panel shows sites where $0 < \omega_{ML} < 1$. Note the change in x-axis scale between plots in (a) and (b), reflecting the paucity of sites in Primates with strong evidence ($LRT_{SLR} <-12$) for purifying selection.

similar patterns, but with more weight towards less-negative LRT_{SLR} values, higher ω_{ML} , and larger $CI_{95\%}$. These differences highlight the impact of branch length on the amount of confidence with which ω can be estimated, showing that the low branch length of the Primates clade rarely yields ω_{ML} estimates with $CI_{95\%}$ intervals smaller than 1, while the bulk of sites from the Mammalia dataset have a relatively small $CI_{95\%}$. Thus, the ω_{ML} distribution from datasets with low branch lengths (e.g., Figure 5.4) should be interpreted with caution, and any comparison between different sites or datasets must be sensitive to the amount of statistical confidence placed on each estimate.

The statistical information at each site could be used to identify sites evolving under purifying or positive selection with statistical confidence. Sites with a CI_{upper} , the upper bound of the $CI_{95\%}$ interval, below $\omega = 1$ showed evidence of purifying selection with

an expected 5% FPR, and sites with a CI_{lower} above $\omega = 1$ showed evidence of positive selection with an expected 5% FPR. In both cases, the 5% FPR was expected under SLR's null model of neutral evolution. There was a strong relationship between CI_{upper} and the χ_1^2 approximation to the LRT_{SLR} distribution, whereby the set of sites with $\text{CI}_{upper} < 1$ was exactly equivalent to the set of sites with LRT_{SLR} below the negative χ_1^2 95% critical value. Similarly, the sites with $\text{CI}_{lower} > 1$ were those with LRT_{SLR} above the χ_1^2 95% critical value. Because of this equality, I will refer to LRT_{SLR} values instead of $\text{CI}_{95\%}$ intervals when discussing sites with significant evidence for purifying or positive selection. In some cases, however, use of the full $\text{CI}_{95\%}$ will be preferable, as the LRT_{SLR} critical values only correspond to one end of the $\text{CI}_{95\%}$, depending on whether the site shows greater evidence for purifying or positive selection.

Table 5.4 summarizes the results from using the χ_1^2 approximation to the LRT_{SLR} distribution to identify sites subject to purifying selection and positive selection at various FPR thresholds. The left group of columns show the number and proportion of sites with evidence for positive selection at nominal 10%, 5%, and 1% FPR thresholds, respectively, as well as an expected 5% FDR calculated using the Benjamini Hochberg method for FDR control [Benjamini & Hochberg, 1995]. The two groups of columns on the right show the result of breaking sites into three groups (positive, negative, and neutral) based on the result of a χ_1^2 test at a given FPR threshold.

The positive selection results demonstrated that between 0.2% to 1% of sites could be identified as under positive selection in mammals at nominal FPR thresholds, but different species groups yielded strikingly different estimates of the proportion of positively-selected sites. At a 5% FPR threshold, Primates, Laurasiatheria, Eutheria, and Mammalia produced roughly comparable proportions of positively-selected sites, ranging from 0.43% to 0.72%. The proportions of positively-selected sites in these groups were higher using a 10% FPR threshold (ranging from 0.70% to 0.97% of sites) and lower using a 1% FPR threshold (ranging from 0.14% to 0.28%). When the FDR was controlled using the Benjamini-Hochberg method, however, only the Eutheria and Mammalia groups yielded a substantial number of positively-selected sites; the Primates and Laurasiatheria data were likely limited in their power to yield positively-selected sites after FDR control due to their lower total branch lengths. Interestingly, the Glires, Atlantogenata, Sparse Glires and Sparse Mammalia datasets produced much lower proportions of positively-selected sites across all FPR thresholds. At $\text{FDR} < 0.05$, all four groups yielded zero significant PSCs.

In Mammalia, the breakdown of sites into positive, negative and neutral categories at

both FPR thresholds produced a pattern similar to the ω_{ML} distribution, with overwhelming purifying constraint (83.87% of sites at 5% FPR), a small proportion of neutrally-evolving sites (15.57%), and a small fraction of positively-selected sites (0.55%). As expected given the use of a fixed LRT_{SLR} threshold to identify purifying sites, the fraction of sites confidently identified as under purifying selection showed a strong dependency on the branch length of the species set, with a much higher power in Mammalia than in Primates (83.87% vs. 15.97%).

Even for the Mammalia dataset, which encapsulated roughly 7.5 synonymous substitutions per site on average, one might reasonably expect that the power to confidently identify sites under purifying selection, though higher than in Primates, was still less than 100%. If this is the case, then the proportion of confidently identified purifying sites must be an underestimate of the true proportion of negatively-selected sites (and by symmetry, absent any methodological bias, the same should be true for positively-selected sites). As a result, the fractions of positively- and negatively-selected sites in Table 5.4 should not be taken as best estimates of the actual proportions of such sites, but more appropriately as lower bounds on those proportions. In fact, In the next two sub-sections, I will separately consider the issue of using the sets of sitewise estimates in different species groups to estimate the proportion of sites subject to purifying and positive selection in mammals.

5.3.3.3 Estimating the proportion of negatively-selected sites

For sites under purifying selection,

In , the value of $\approx 95\%$ based on $\omega_{ML} < 1$ (Table 5.3) is likely closer to the true number; despite the caveats involved in ignoring the uncertainty involved in ω_{ML} estimates, the 95% number was surprisingly consistent across different species groups and branch lengths, ranging from $\approx 91\%$ in Primates to $\approx 97\%$ in Sparse Mammalia.

5.3.3.4 Estimating the proportion of positively-selected sites

The pattern of the prevalence of positive selection across species groups was surprising. First,

, showing no sign of the expected correlation with branch length. Theory predicts, and many studies have confirmed [Anisimova *et al.*, 2001, 2002; Massingham & Goldman, 2005], that the power of LRT-based tests for non-neutral selection should increase with branch length, as the discrimination between non-synonymous and synonymous substitution rates becomes more confident with more fixed substitutions. This was certainly the case for

identifying purifying selection, but there was no obvious correlation between higher branch lengths and higher numbers of confidently identified PSCs.

5.3.3.5 Correlations between branch length, effective population size and site-wise summary statistics

To quantify this observation, Spearman's rank correlation coefficients between the median branch length of each species set and each of several summary statistics were calculated (Table 5.5). Although the number of samples was small with only eight species groups, these correlations should be able to provide some indication as to which aspects of the sitewise data might be easily attributed to the effects of branch length and which aspects suggested an alternative (e.g., biological or artefactual) cause for the differences between species groups. The results were quite striking: the site classifications (constant, synonymous and non-synonymous) and the proportion of negatively-selected sites at $P_{\chi^2} < 0.05$ were strongly correlated with branch length, while the other factors, including mean ω_{ML} , the proportion of sites with $\omega_{ML} < 1$ and $\omega_{ML} < 0.5$, and the proportion of positively-selected sites, were weakly correlated or largely uncorrelated with branch length.

The results in Table 5.5 can be interpreted in a number of ways. First, they emphasized the unambiguous correlation between the branch length in a tree and the power to detect sitewise purifying selection. However, they also showed that various measures based on sitewise estimates contained variation between species groups that was not well-explained by branch length.

Refer to C
Pollard 20
showed that
mates and
vores have
BDS (bias
gent subst
while glire

5.3.4 Modeling the global distribution of sitewise selective pressures

| Value | Branch Length | | N_E | |
|---------------------|---------------|---------|-------|---------|
| | Rho | P-value | Rho | P-value |
| Median BL | - | - | 0.37 | 0.36 |
| Population Size | 0.37 | 0.36 | - | - |
| Mean ω_{ML} | 0.14 | 0.75 | -0.54 | 0.17 |
| Constant | -1.00 | 0.00 | -0.37 | 0.36 |
| Synonymous | 0.88 | 0.01 | 0.48 | 0.23 |
| Nsynonymous | 0.98 | 0.00 | 0.35 | 0.40 |
| $\omega_{ML} < 1$ | 0.38 | 0.36 | 0.87 | 0.01 |
| $\omega_{ML} < 0.5$ | 0.12 | 0.79 | 0.73 | 0.04 |
| PSC _{5%} | 0.07 | 0.88 | -0.70 | 0.05 |
| NSC _{5%} | 0.98 | 0.00 | 0.48 | 0.23 |

Table 5.5: Correlations between median non-gap branch length, N_E , and various summary statistics of the sitewise data across eight species sets. The magnitude (Rho) and significance (P-value) of Spearman's rank correlations between variables were calculated using N_E values from Table 5.1 and branch lengths and summary statistics from Tables 5.3 and 5.4. All summary statistics except for mean ω_{ML} were measured as a fraction of total sites. More highly significant correlations are shaded in darker blue. N_E – estimated effective population size, BL – non-gap branch length, PSC_{5%} – positively-selected codons at $P_{\chi_1^2} < 0.95$.

ω

| Species Set | Data Type | Log-normal | | Gamma | | Exponential | | Beta | | Weibull | |
|-----------------|---------------|----------------|-------|----------------|-------|----------------|-------|----------------|-------|----------------|-------|
| | | $\bar{\omega}$ | % > 1 |
| Primates | ω_{ML} | 0.11 | 1.53 | 0.25 | 7.31 | 0.25 | 1.82 | 0.25 | 0.00 | 0.14 | 2.96 |
| Glires | | 0.15 | 2.07 | 0.16 | 3.35 | 0.16 | 0.16 | 0.20 | 0.00 | 0.12 | 2.56 |
| Laurasiatheria | | 0.33 | 3.52 | 0.20 | 5.34 | 0.20 | 0.74 | 0.23 | 0.00 | 0.19 | 4.08 |
| Atlantogenata | | 0.04 | 0.59 | 0.15 | 3.33 | 0.15 | 0.12 | 0.20 | 0.00 | 0.07 | 1.30 |
| Eutheria | | 0.72 | 6.10 | 0.20 | 4.62 | 0.20 | 0.64 | 0.24 | 0.00 | 0.22 | 5.18 |
| Mammalia | | 0.77 | 6.57 | 0.19 | 4.26 | 0.19 | 0.57 | 0.23 | 0.00 | 0.22 | 5.16 |
| Sparse Glires | | 0.06 | 0.87 | 0.13 | 2.64 | 0.13 | 0.06 | 0.18 | 0.00 | 0.08 | 1.46 |
| Sparse Mammalia | | 0.20 | 2.64 | 0.15 | 2.85 | 0.15 | 0.11 | 0.19 | 0.00 | 0.13 | 2.74 |
| Primates | $CI_{95\%}$ | 0.40 | 4.09 | 0.42 | 4.56 | 0.37 | 6.83 | 0.42 | 0.00 | 0.41 | 5.55 |
| Glires | | 0.21 | 0.14 | 0.21 | 0.17 | 0.18 | 0.44 | 0.21 | 0.00 | 0.20 | 0.27 |
| Laurasiatheria | | 0.22 | 0.95 | 0.23 | 0.45 | 0.22 | 1.05 | 0.23 | 0.00 | 0.23 | 0.64 |
| Atlantogenata | | 0.30 | 0.58 | 0.32 | 0.55 | 0.26 | 2.14 | 0.32 | 0.00 | 0.30 | 1.26 |
| Eutheria | | 0.19 | 2.28 | 0.18 | 0.79 | 0.18 | 0.42 | 0.18 | 0.00 | 0.18 | 1.32 |
| Mammalia | | 0.18 | 2.35 | 0.18 | 0.74 | 0.17 | 0.31 | 0.17 | 0.00 | 0.17 | 1.27 |
| Sparse Glires | | 0.23 | 0.14 | 0.24 | 0.24 | 0.20 | 0.67 | 0.25 | 0.00 | 0.23 | 0.43 |
| Sparse Mammalia | | 0.16 | 0.10 | 0.17 | 0.09 | 0.15 | 0.13 | 0.17 | 0.00 | 0.16 | 0.11 |

Table 5.6: Mean ω and the percentage of sites with $\omega > 1$ based on maximum-likelihood fits of parametric distributions to sitewise estimates. For each species set and each one of five distribution types (log-normal, gamma, exponential, beta, and weibull) 100 replicate datasets of 1 million sites were sampled with replacement from the genome-wide dataset and the maximum likelihood distribution parameters were numerically optimized (see text for details). Columns show, for each distribution, the median value across 100 replicates of mean ω ($\bar{\omega}$) and the probability mass with $\omega > 1$, expressed as a percentage (%> 1). The top eight rows show the results based on fitting parameters to sitewise ω_{ML} estimates, and the bottom eight rows show the results based on fitting parameters to sitewise $CI_{95\%}$ estimates. Note the greater consistency of $\bar{\omega}$ and %> 1 across distribution types for the $CI_{95\%}$ -based fits.

We used the `fitdistr` function of the MASS package for R to fit five distributions (gamma, lognormal, beta, Weibull, and exponential) to the vertebrate dN/dS values and subsequently calculated Akaike's Information Criterion (AIC) for each fit. For all optimizations, a constant value of 0.001 was added to sites where $dN/dS = 0$ in order to satisfy the optimizers requirement that the probability functions have a defined value for all input data. Similarly, sites with $dN/dS \leq 1$ were excluded from the analysis for the beta optimization. All distributions were also separately fit to the subset of sites with $dN/dS > 1$; the AIC values from these optimizations were used to compare the fit of the beta distribution to the others.

The `fitdistr` produced the following optimized parameters for each function: gamma (`shape=0.271, rate=1.203`), lognormal (`meanlog=-4.079, sdlog=2.863`), beta (`shape1=0.257, shape2=1.431`), Weibull (`shape=0.3882, scale=0.07151`) exponential (`rate=4.441`). The beta distribution yielded the lowest AIC when compared to the fit of other distributions to the subset of sites where $dN/dS > 1$ (-2.33e7 versus the next best equivalent AIC of -2.11e7 for the lognormal). Of the distributions which were fit to the whole dataset, the lognormal distribution yielded the lowest AIC (-2.11e7), followed by gamma (-2.03e7) and exponential (-6.45e6).

5.3.5 Simulations to evaluate the fit of empirical data to a simulated poisson process at varying branch lengths and ω ratios

The large difference in branch lengths covered by the species sets under investigation led to some uncertainty in whether differences in the observed summary statistics, such as the mean ω_{ML} or the proportion of sites with $\omega_{ML} < 1$, were due to biological differences between species groups (i.e., differences in the efficacy of purifying selection resulting from population size differences) or to branch length effects. For example, two

I performed a small simulation study to further investigate the interactions between branch length, effective population size, and the summary statistics presented in Table 5.3,

| Species Set | Distribution | AIC | dAIC | Parameter A | Parameter B |
|-----------------|--------------|------------|----------|-------------|-------------|
| Primates | Beta | 161168.45 | 0.00 | 1.78 | 2.50 |
| | Lognormal | 217743.70 | 56575.25 | -1.14 | 0.65 |
| | Gamma | 231082.15 | 13338.45 | 2.12 | 5.06 |
| | Weibull | 242149.95 | 11067.80 | 1.34 | 0.45 |
| | Exponential | 261904.40 | 19754.45 | 2.68 | |
| Glires | Lognormal | 312788.30 | 0.00 | -1.76 | 0.59 |
| | Beta | 332776.25 | 19987.95 | 1.57 | 5.75 |
| | Gamma | 341872.95 | 9096.70 | 1.70 | 8.02 |
| | Weibull | 364256.80 | 22383.85 | 1.15 | 0.21 |
| | Exponential | 385083.15 | 20826.35 | 5.42 | |
| Laurasiatheria | Lognormal | 514030.30 | 0.00 | -1.80 | 0.77 |
| | Beta | 528667.45 | 14637.15 | 1.25 | 4.27 |
| | Gamma | 558537.00 | 29869.55 | 1.49 | 6.52 |
| | Weibull | 570535.20 | 11998.20 | 1.12 | 0.23 |
| | Exponential | 575734.00 | 5198.80 | 4.56 | |
| Atlantogenata | Beta | 109348.25 | 0.00 | 2.37 | 4.93 |
| | Lognormal | 119888.05 | 10539.80 | -1.34 | 0.53 |
| | Gamma | 126416.00 | 6527.95 | 2.75 | 8.72 |
| | Weibull | 143678.60 | 17262.60 | 1.33 | 0.33 |
| | Exponential | 175474.60 | 31796.00 | 3.84 | |
| Eutheria | Lognormal | 1180885.50 | 0.00 | -2.40 | 1.20 |
| | Weibull | 1270446.50 | 89561.00 | 0.80 | 0.16 |
| | Gamma | 1295309.50 | 24863.00 | 0.80 | 4.35 |
| | Beta | 1307444.50 | 12135.00 | 0.69 | 3.16 |
| | Exponential | 1308964.50 | 1520.00 | 5.46 | |
| Mammalia | Lognormal | 1310778.00 | 0.00 | -2.51 | 1.26 |
| | Weibull | 1403090.50 | 92312.50 | 0.78 | 0.15 |
| | Gamma | 1433878.50 | 30788.00 | 0.75 | 4.30 |
| | Beta | 1456833.00 | 22954.50 | 0.65 | 3.11 |
| | Exponential | 1459667.50 | 2834.50 | 5.77 | |
| Sparse Glires | Beta | 161165.45 | 0.00 | 2.04 | 6.12 |
| | Lognormal | 161264.00 | 98.55 | -1.60 | 0.53 |
| | Gamma | 176568.30 | 15304.30 | 1.99 | 8.23 |
| | Weibull | 198910.90 | 22342.60 | 1.20 | 0.24 |
| | Exponential | 225609.40 | 26698.50 | 5.00 | |
| Sparse Mammalia | Lognormal | 403243.50 | 0.00 | -2.09 | 0.68 |
| | Gamma | 447431.20 | 44187.70 | 1.31 | 7.74 |
| | Beta | 450936.70 | 3505.50 | 1.18 | 5.70 |
| | Weibull | 462813.80 | 11877.10 | 1.06 | 0.16 |
| | Exponential | 471324.15 | 8510.35 | 6.67 | |

Table 5.7: AIC values and parameters for maximum-likelihood fits of parametric distributions to sitewise CI_{95%} estimates. Distributions were fit to 100 replicate datasets for each species set as in Table 5.6. For each species set and distribution type, the median Akaike information criterion (AIC) and parameter estimates are shown. Distributions are sorted according to their median AIC value (where a lower AIC corresponds to a better fit to the data), and the difference in AIC to the next-best fitting distribution is displayed (dAIC). The lognormal and beta distributions are ranked first or second in the mammalian superorder subgroups, while lognormal, weibull and gamma distributions fit the Eutheria and Mammalia datasets best. The named parameters corresponding to parameters A and B for each distribution are as follows: lognormal (A=meanlog, B=sdlog), weibull (A=shape, B=scale), gamma (A=shape, B=rate [where rate=1/scale]), beta (A=shape1 [α], B=shape2 [β]), exponential (A=rate [λ]).

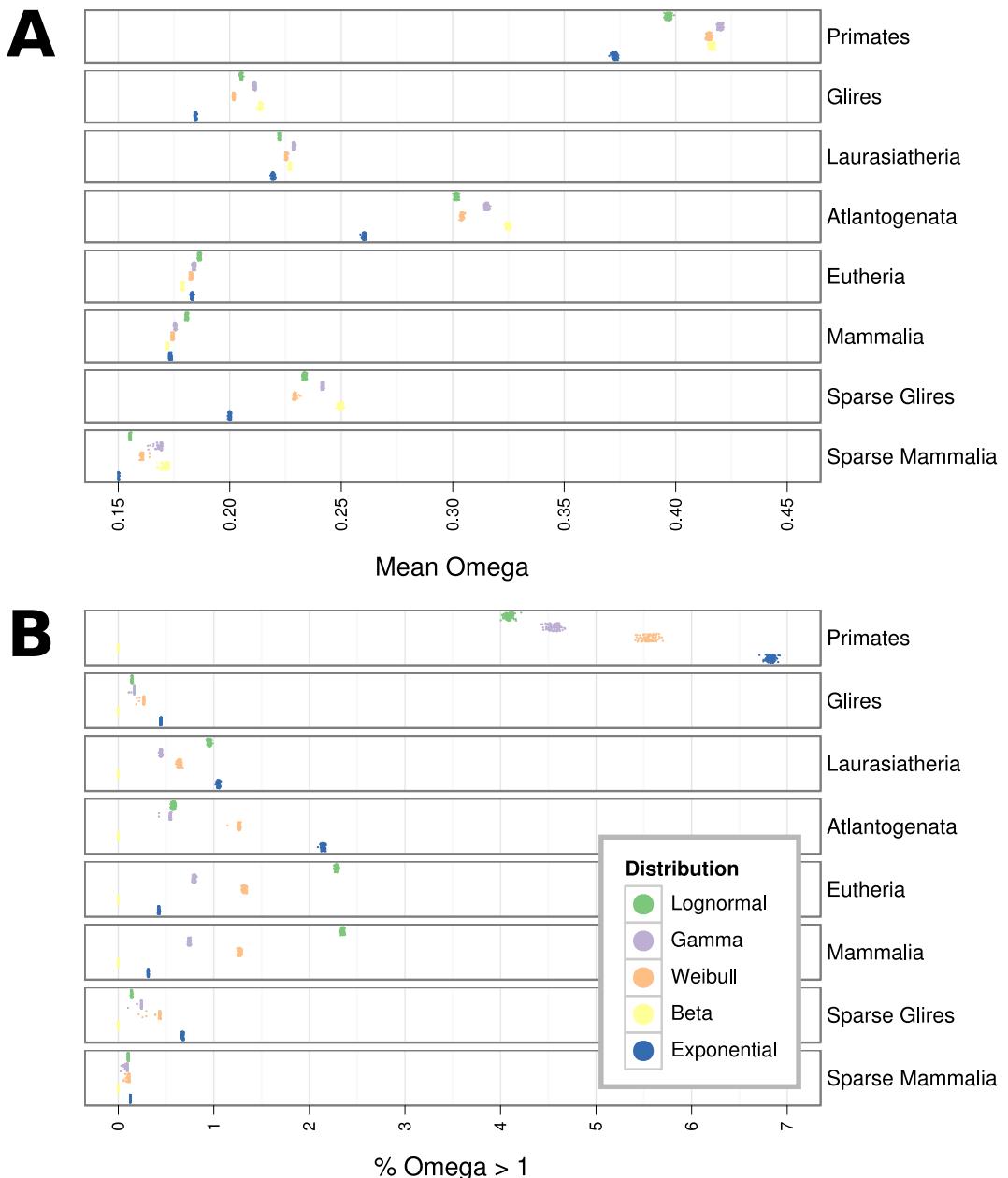


Figure 5.6: The mean value (A, Mean Omega) and the percentage of probability density with $\omega > 1$ (B, Omega> 1) from maximum likelihood fits of parametric distributions to sitewise CI_{95%} estimates. Each point represents the maximum likelihood fit of one distribution to 1 million points sampled from the genome-wide dataset with replacement; 100 such fits were generated for each distribution and species set. Note that the beta distribution is only defined on the interval (0,1), so the percentage of sites with $\omega > 1$ was always zero.

5.3.6 Simulations to evaluate the power to detect positive selection and estimate selective pressures

Previous simulations on the power and accuracy of maximum-likelihood methods for detecting sitewise positive selection have provided strong evidence for increased power with increased branch length and number of taxa [Anisimova and Yang 2002 PMID:12032251, Massingham and Goldman 2005 PMID:15654091]. Other major effects observed have been a reduction in power when branch lengths are very short, due to the scarcity of data in the form of observed substitutions, and a reduction in accuracy when branches are very long and the ancestral reconstruction procedure becomes inaccurate due to saturation of substitutions at synonymous sites [Anisimova and Yang 2002 PMID:12032251]. These results have provided general guidance to empirical analyses, but the potential effects of tree shape, divergence level, distribution of dN/dS levels, and misalignment error make it difficult to extrapolate expected power levels or error rates from generic simulations to specific empirical analyses and real-world datasets. In order to estimate the power of the SLR method when applied to the present set of species and to quantify the power gained from the additional 20 mammalian genomes, we ran a series of simulation experiments with parameters tuned specifically to the analysis of mammalian gene families. Based on the previous studies described above we hypothesized that the addition of 20 mammalian genomes to the available dataset would significantly improve SLRs sensitivity for detecting positive selection at a reasonable error rate, with some portion of that improvement coming from a reduction in alignment error due to the shorter average length of branches in the phylogenetic tree being aligned.

We used the Indelible program [Fletcher and Yang 2010 PMID:19423664] to simulate 100 replicate codon alignments with a root sequence length of 500 codons for each of three trees: all 29 Eutherian mammals, the 9 mammals with high-coverage genomes, and the four-species Human-Mouse-Rat-Dog quartet used in a number of previous comparative analyses. The dN/dS value at each simulated site was drawn from a discretized lognormal distribution with $\log(\text{mean})=-1.864$ and $\log(\text{sd})=1.201$ with the maximum dN/dS capped at 3. This distribution yielded a mean dN/dS of 0.277 and 6% of sites with dN/dS > 1 , which is consistent with our estimates from the global mammalian distribution. We ran each set of simulations twice: once with an insertion and deletion (indel) rate of zero, and once with an indel rate of 0.05 indel events per substitution event. The length of each indel event was drawn from a discretized power-law distribution with a parameter of 1.8, a maximum insertion or deletion length of 40 codons, and equal insertion and deletion lengths

and probabilities. Each simulated alignment was aligned using PRANKs codon model of evolution and analyzed with SLR. The resulting SLR score at each human sequence position was compared to the true dN/dS at the equivalent site in the true alignment and used to calculate the power and accuracy of the detection of positive selection for each set of simulated alignments. A summary of the results is provided in Supplementary Table S17.12.

The simulation results without indels show a dramatic increase in the ability to detect sitewise selective pressures and sitewise positive selection in larger mammalian trees. We used ROC curves based on the true and inferred dN/dS values to calculate a number of statistics summarizing the performance of sitewise inference under each tree tested. The Spearmans rank correlation coefficient between inferred and true dN/dS was 0.749, 0.849, and 0.942 in the 4-taxon, 9-taxon, and 29-taxon trees, respectively, representing a 10% increase in the accuracy of inferred maximum-likelihood dN/dS values as a result of the added 20 mammalian species. The number of true positives recovered when controlling for a false discovery rate (FDR) below 0.1 was 50, 782, and 3760 for the three trees; for FDR ≥ 0.05 , the numbers of true positives were 10, 429, and 2990. This represents a 4.8-fold increase for FDR ≥ 0.1 and a 6.9-fold increase for FDR ≥ 0.05 resulting from the additional 20 species in the tree. It should be noted that the two previous power estimates re

5.3.7 Evaluation of the effect of GC content, recombination rate, and codon usage on sitewise dN/dS estimates and the detection of positive selection

Multiple lines of evidence have lent support to the hypothesis that GC-biased gene conversion (BGC) has been a major force in the evolution of mammalian genomes [Galtier et al 2001 PMID:11693127, Galtier 2003 PMID:XYZ, Dreszer et al 2007 PMID:17785536]. Both empirical and theoretical results have shown that BGC can significantly affect patterns of observed substitutions in both selectively neutral and functionally constrained sites [Galtier et al 2009 PMID:19027980, Berglund et al 2009 PMID:19175294]. Recently, Ratnakumar et al. [2010, PMID:20643747] re-analyzed the dataset of positively-selected genes from Kosiol et al. [2008, PMID:18670650] for signatures of BGC and found that up to 20% of cases of identified elevated dN/dS ratios could be due to BGC rather than adaptive evolution. However, the strongest signals of BGC were found only in genes showing signals of positive selection along short branches in the phylogenetic tree using so-called

branch-site models of evolution; when the authors looked for similar BGC signatures in genes with evidence for positive selection at specific sites throughout the mammalian tree (e.g., genes with significant LRTs for PAMLs sites model) they found no evidence for a strong BGC influence [Ratnakumar et al. 2010 PMID:20643747].

The above evidence suggests that although BGC has the potential to produce misleading signals of branch-specific positive selection near recombination hotspots, the positively-selected sites we detected should not be strongly influenced by the non-adaptive effects of BGC since the dN/dS level detected by SLR is estimated from across the entire input phylogeny [Massingham 2005 PMID:15654091]. This is consistent with the observation that recombination hotspots (where most recombination in humans and other mammals occurs [Myers et al. 2005 PMID:16224025]) tend not to be maintained over long evolutionary periods, although larger-scale recombination rates are likely more conserved [Winckler et al 2005 PMID:15705809]. Still, due to the potential confounding implications of BGC on the interpretation of signals of positive selection, we found it worthwhile to empirically test for any BGC effect on our data.

The BGC model predicts a recombination-associated drive towards the fixation of GC alleles at heterozygous sites, resulting in an expected correlation between AT to GC (or weak-to-strong, W-S) mutational bias and recombination rate [Galtier and Duret 2007, PMID:17418442]. This bias can lead to elevated dN/dS estimates in coding regions, particularly in GC-rich regions where W-S mutations are more likely to result in nonsynonymous changes [Berglund et al 2009]. Ratnakumar and colleagues identified three ways of distinguishing potential BGC effects from true signals of positive selection in protein-coding regions: (a) positive selection is not expected on its own to result in a strong W-S bias, (b) a BGC-associated W-S biased mutation pattern should extend to noncoding sites flanking the affected coding region, and (c) BGC is associated with recombination hotspots and regions of high recombination rates (and most strongly with male-specific rates) while there is no empirical evidence linking positive selection with higher recombination rates in mammals, although natural selection should theoretically be more efficient in regions of high recombination [Ratnakumar et al. 2010, PMID:20643747]. We could not use (a) or (b) to detect possible BGC influence since we did not calculate inferred ancestral mutations for either the coding or flanking noncoding regions of the mammalian gene families studied here. Instead, we turned to point (c) and tested for a correlation between signals of positive selection and an increase in recombination rates, especially the male-specific rate and in regions of high GC content. The predictions of the BGC hypothesis suggest

that if our sitewise data do contain a strong BGC influence, then the positively-selected sites we detected would be expected to be associated with regions of high male-specific recombination.

We combined the sitewise codon data with male, female, and sex-averaged recombination rates derived from the deCODE map (using rates averaged over genomic bins of 1Mb downloaded from the UCSC human genome browser hg19 release) and human GC content calculated in 10-kb windows and analyzed sites within various quantiles of GC content, mean recombination rate, and sitewise statistics. Supplementary Table S17.13 contains summaries for each subset. The LRT statistic section shows that sites with higher LRT statistics (which corresponds to weaker purifying selection when the value is below zero and stronger positive selection when the value is above zero) show decreasing recombination rates; this trend holds true even for the highest quantile (mean signed_lrt between 3.648 and 108.850), which is composed entirely of sites with evidence for positive selection. In other words, the bulk of positively-selected sites are in regions of lower than average male recombination rates – exactly opposite what would be expected in the face of strong BGC effects. The Male Recombination quantiles show a similar trend, with the mean dN/dS, mean signed LRT and the proportion of sites identified as positively-selected (pos_f) all consistently decreasing as the recombination rate increases. The GC content quantiles showed a slightly different pattern. Although the mean LRT decreased and male recombination increased monotonically with increasing GC content, the mean dN/dS and fraction of positive sites started low, increased to a maximum in the middle range of GC content, and decreased again in regions of high GC content. Thus, although the GC content quantiles were similar to the male recombination quantiles in their higher range (with similar mean dN/dS, mean LRT, and pos_f values), they differed slightly in their lower range (with lower dN/dS and pos_f for low GC quantiles). Although the exact reason for such a pattern is unclear, it is consistent with the existence of altered or constrained selective or mutational dynamics at the extreme ends of the genomic distribution of GC content. As GC content has been shown to correlate with myriad structural and evolutionary features of mammalian genomes [Xia et al. 2009 PMID:19521505], the existence of other (possibly unrelated) confounding influences such as CpG mutability or isochore structure is likely.

Theoretical and empirical evidence pointed towards an increased sensitivity of dN/dS estimates to BGC influence in regions of high GC content, so we separated out the top 10% of sites by GC content and analyzed them according to quantiles of male recombination rate (Supplementary Table S17.13, High GC, Male Recombination. The middle four

recombination quantiles showed a similar pattern to that observed for all GC contents, with mean LRT decreasing with increased male recombination and mean dN/dS and pos_f decreasing or hovering around values slightly lower than those observed across all GC contents (e.g., mean dN/dS in the 1-25% bin is 0.207 for the top 10% GC sites, but 0.249 for the same recombination bin across all sites). The highest recombination bin of the top 10% GC sites showed a strikingly different pattern, however, with mean dN/dS=0.348, mean LRT=-11.262, and pos_f=0.0338. These values suggest a strong shift towards higher dN/dS values and more positively-selected sites. This jump in values in the highest recombination bin is not seen in the highest male recombination bin across all GC contents (mean dN/dS=0.207, mean LRT=-16.616, pos_f=0.00932) or for the highest female recombination bin for the top 10% of GC sites (mean dN/dS=0.164, mean LRT=-18.475, pos_f=0.00449). Although the small number of sites in the bin of interest compared to other bins suggests possible stochastic artifacts, the shift is dramatic, directly opposite to the trends observed for the female recombination rates and for male recombination rates in regions of lower GC content, and is in agreement with the BGC prediction of elevated dN/dS estimates in regions of high GC content and male-specific recombination rates. This evidence raises the interesting possibility that BGC may have a detectable, if rather minor, impact on sitewise dN/dS estimates across the mammalian phylogeny. It is highly unlikely, however, that any such effect – which in our analysis was only detectable in 0.05% of sites with the most extreme GC content and recombination rates – has contaminated our codon-specific estimates with more than a negligible amount of noise resulting from the neutral but biased process of BGC.

5.4 Conclusions

[...]

Chapter 6

The use of sitewise selective pressures to characterise the evolution of genes and domains in mammals

6.1 Introduction

A clear definition of what constitutes a positively selected gene (PSG) is elusive and depends heavily on the model the data. Neutral theory, however, does provide a precise definition: the dN/dS ratio gives a natural threshold for distinguishing purifying from positive selection at $dN/dS = 1$, above which point positive selection can be inferred. This leads logically to the most stringent definition of a PSG as a gene whose overall mean dN/dS ratio is greater than 1. However, as noted by Yang [5], this definition of a PSG is extremely conservative and is likely to be met by very few genes. To examine this hypothesis in the context of mammalian proteins, we averaged the site-wise dN/dS estimates from each gene tree tested to obtain an average dN/dS value, mean- dN/dS . At a threshold of mean- $dN/dS = 1$, 126 out of 15,451 gene trees analyzed showed evidence for positive selection (0.82%). At a threshold of mean- $dN/dS = 0.5$, 1,438 genes were PSGs (9.3%).

More powerful scans for positive selection relax the requirement that the overall dN/dS of a gene is above 1 under the assumption that the force of positive selection is unlikely to act along an entire protein sequence and throughout the entire evolutionary history of a gene tree. This led to models that allow for variation in the dN/dS ratio along the sequence or along the phylogeny. If the data shows statistically significant evidence for positive selection under a given model, then it is considered a PSG thus, the definition of

a PSG is tightly linked to the evolutionary model being applied and the evolutionary clade or lineage being analyzed. In the case of the SLR method (which allows for variation of dN/dS along the sequence but not along the phylogeny), the most appropriate definition of a PSG is a gene where one or more sites show statistical evidence for positive selection after correcting for multiple tests. With this definition, we find 2,865 positively selected gene trees in mammals (18.5%). We may also formulate a more conservative definition in an attempt to minimize the number of false positives resulting from a background of reduced purifying constraint, considering genes with mean-dN/dS of less than 0.5 and at least 2 positively-selected codons to be PSGs. This definition yields 739 positively selected gene trees (4.8%).

Although the mean-dN/dS > 1 definition of a PSG is unreasonably conservative (and, upon manual inspection, picks out many proteins with especially severe alignment or gene annotation errors), the numbers of PSGs resulting from both codon-based definitions applied to the current data (4.8% at the low end and 18.5% at the high end) are similar to the range of other recently published comparative scans for positive selection, which identified between 3.3% [6] and 16.6% [7] of genes as PSGs in mammals or primates.

6.2 Comparison of sitewise results to previously described sets of positively selected genes

6.3 Using sitewise selective pressures to characterise the evolution of genes

6.3.1 Identifying genes subject to positive selection

6.3.2 Identifying genes subject to strong or weak purifying selection

6.4 Using sitewise selective pressures to characterise the evolution of protein domains

[To cite: Moses & Durbin [2009] — They use residue probabilities to generate “preferred residue” predictions, and look at human polymorphism to assign whether a SNP is moving

to or from a preferred residue.]

In order to evaluate the prevalence of positive selection in mammals for various domain structures, we used the Pfam domain mappings from the Ensembl database (release version 54) to annotate the site-wise dN/dS values with domain assignments. We mapped Pfam protein domain annotations from all sequences in a gene tree onto the alignment, keeping only features with a hit score greater than 20 and alignment sites with greater than 4 columns and an inferred dN/dS value of less than 50. We then removed any domains with fewer than one thousand annotated sites, to avoid errors resulting from small sample sizes.

The domain annotations were collated in a variety of ways: (1) taking the mean omega value across all annotated sites [omega column in the table below], (2) counting the number of PSCs within all annotated sites [psc], and (3) counting the number of PSCs per annotated site [psc_corr]. The len column represents the number of alignment sites containing the given Pfam annotation. When the list of domains is sorted by the total number of annotated PSCs we see the more familiar positively-selected domains at the top of the list (immunoglobulin, collagen), followed by more novel positively-selected domains such as the 7 transmembrane receptor, protein kinase domain, and ion transport protein.

6.4.1 Identifying protein domains subject to positive selection

6.4.2 Identifying protein domains subject to strong or weak purifying selection

6.5 Identifying genes under unusual selective pressures in mammalian superorders

Chapter 7

Evolution of protein-coding genes in gorilla and the African apes

7.1 Introduction

- 7.1.1 The gorilla and other primate genome projects
- 7.1.2 Incomplete lineage sorting
- 7.1.3 Effective population sizes of extant and ancestral primate populations
- 7.1.4 Measuring shifts in selective pressures using branch-specific likelihood ratio tests
- 7.1.5 Data quality concerns: sequencing, assembly and alignment error

7.2 Constructing codon alignments of one-to-one orthologous genes in six primate species

- 7.2.1 Identification of genes with one-to-one homology
- 7.2.2 Collection of homologous DNA sequences from genome- or transcript-based multiple alignments
- 7.2.3 Filtering sequence regions with low sequence quality⁸⁶
- 7.2.4 Filtering sequence regions with high substitution counts
- 7.2.5 Filtering sequence regions with evidence of incomplete lineage sorting

Chapter 8

Gorilla part 2

- 8.1 Analysis of incomplete lineage sorting in the African great apes within and nearby protein-coding genes**
- 8.2 Analysis of dN/dS levels in six primate genomes**
 - 8.2.1 Genome-wide dN/dS in six primates and their ancestors**
 - 8.2.2 Genome-wide dN/dS in regions of differing sitewise constraint**
 - 8.2.3 Analysis of the impact of sequence and alignment filtering on primate dN/dS estimates**
- 8.3 Conclusions and future work**

Bibliography

- ALBERS, C., CVEJIC, A., FAVIER, R., BOUWMANS, E., ALESSI, M., BERTONE, P., JORDAN, G., KETTLEBOROUGH, R., KIDDLE, G., KOSTADIMA, M., READ, R., SIPOS, B., SIVAPALARATNAM, S., SMETHURST, P., STEPHENS, J., VOSS, K., NURDEN, A., RENDON, A., NURDEN, P. & OUWEHAND, W. (2011). Exome sequencing identifies NBEAL2 as the causative gene for gray platelet syndrome. *Nat Genet*, **43**, 735–7. [18](#)
- ALTSCHUL, S.F., MADDEN, T.L., SCHÄFFER, A.A., ZHANG, J., ZHANG, Z., MILLER, W. & LIPMAN, D.J. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research*, **25**, 3389–3402. [9](#)
- ANISIMOVA, M., BIELAWSKI, J. & YANG, Z. (2001). Accuracy and power of the likelihood ratio test in detecting adaptive molecular evolution. *Mol Biol Evol*, **18**, 1585–92. [70](#)
- ANISIMOVA, M., BIELAWSKI, J. & YANG, Z. (2002). Accuracy and power of bayes prediction of amino acid sites under positive selection. *Mol Biol Evol*, **19**, 950–8. [45](#), [70](#)
- ANISIMOVA, M., NIELSEN, R. & YANG, Z. (2003). Effect of recombination on the accuracy of the likelihood method for detecting positive selection at amino acid sites. *Genetics*, **164**, 1229–36. [45](#)
- ARCHIBALD, A., BOLUND, L., CHURCHER, C., FREDHOLM, M., GROENEN, M., HARIZIUS, B., LEE, K., MILAN, D., ROGERS, J., ROTHSCHILD, M., UENISHI, H., WANG, J., SCHOOK, L. & SWINE GENOME SEQUENCING CONSORTIUM (2010). Pig genome sequence–analysis and publication strategy. *BMC Genomics*, **11**, 438. [36](#)
- ARCHIBALD, J.D. (1999). Divergence times of eutherian mammals. *Science*, **285**, 2031. [57](#)

BIBLIOGRAPHY

- AVEROF, M., ROKAS, A., WOLFE, K. & SHARP, P. (2000). Evidence for a high frequency of simultaneous double-nucleotide substitutions. *Science*, **287**, 1283–6. [60](#)
- BAKEWELL, M., SHI, P. & ZHANG, J. (2007). More genes underwent positive selection in chimpanzee evolution than in human evolution. *Proc Natl Acad Sci U S A*, **104**, 7489–94. [49](#)
- BAZYKIN, G., KONDRASHOV, F., OGURTSOV, A., SUNYAEV, S. & KONDRASHOV, A. (2004). Positive selection at sites of multiple amino acid replacements since rat-mouse divergence. *Nature*, **429**, 558–62. [55](#)
- BEISSWANGER, S. & STEPHAN, W. (2008). Evidence that strong positive selection drives neofunctionalization in the tandemly duplicated polyhomeotic genes in drosophila. *Proc Natl Acad Sci U S A*, **105**, 5447–52. [48](#)
- BENJAMINI, Y. & HOCHBERG, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, **57**, 289–300. [64](#), [69](#)
- BININDA-EMONDS, O.R.P., CARDILLO, M., JONES, K.E., MACPHEE, R.D.E., BECK, R.M.D., GRENYER, R., PRICE, S.A., VOS, R.A., GITTLEMAN, J.L. & PURVIS, A. (2007). The delayed rise of present-day mammals. *Nature*, **446**, 507–512. [19](#), [57](#)
- BIRNEY, E., ANDREWS, D., CACCAMO, M., CHEN, Y., CLARKE, L., COATES, G., COX, T., CUNNINGHAM, F., CURWEN, V., CUTTS, T., DOWN, T., DURBIN, R., FERNANDEZ-SUAREZ, X., FLICEK, P., GRÄF, S., HAMMOND, M., HERRERO, J., HOWE, K., IYER, V., JEKOSCH, K., KÄHÄRI, A., KASPRZYK, A., KEEFE, D., KOKOCINSKI, F., KULESHA, E., LONDON, D., LONGDEN, I., MELSOPP, C., MEIDL, P., OVERDUIN, B., PARKER, A., PROCTOR, G., PRILIC, A., RAE, M., RIOS, D., REDMOND, S., SCHUSTER, M., SEALY, I., SEARLE, S., SEVERIN, J., SLATER, G., SMEDLEY, D., SMITH, J., STABENAU, A., STALKER, J., TREVANION, S., URETA-VIDAL, A., VOGEL, J., WHITE, S., WOODWARD, C. & HUBBARD, T. (2006). Ensembl 2006. *Nucleic Acids Res*, **34**, D556–61. [14](#)
- BRUNET, F., ROEST CROLLIUS, H., PARIS, M., AURY, J., GIBERT, P., JAILLON, O., LAUDET, V. & ROBINSON-RECHAVI, M. (2006). Gene loss and evolutionary rates following whole-genome duplication in teleost fishes. *Mol Biol Evol*, **23**, 1808–16. [30](#), [40](#)

BIBLIOGRAPHY

- CALLAHAN, B., NEHER, R., BACHTROG, D., ANDOLFATTO, P. & SHRAIMAN, B. (2011). Correlated evolution of nearby residues in drosophilid proteins. *PLoS Genet*, **7**, e1001315. [55](#)
- CASOLA, C. & HAHN, M. (2009). Gene conversion among paralogs results in moderate false detection of positive selection using likelihood methods. *J Mol Evol*, **68**, 679–87. [48](#)
- CEPAS, H., BUENO, A., DOPAZO, J. & GABALDÓN, T. (2007). PhylomeDB: a database for genome-wide collections of gene phylogenies. *Nucl. Acids Res.*, **36**, gkm899. [10](#)
- CHEN, F., MACKEY, A., VERMUNT, J. & ROOS, D. (2007). Assessing performance of orthology detection strategies applied to eukaryotic genomes. *PLoS One*, **2**, e383. [9](#)
- CHURAKOV, G., KRIEGS, J., BAERTSCH, R., ZEMANN, A., BROSIUS, J. & SCHMITZ, J. (2009). Mosaic retroposon insertion patterns in placental mammals. *Genome Res*, **19**, 868–75. [56](#)
- COCK, P., FIELDS, C., GOTO, N., HEUER, M. & RICE, P. (2010). The sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Res*, **38**, 1767–71. [50](#)
- TOCITE** (2011). Citation will be inserted at a later point in time. [56](#)
- CSUROS, M., ROGOZIN, I. & KOONIN, E. (2011). A detailed history of intron-rich eukaryotic ancestors inferred from a global survey of 100 complete genomes. *PLoS Comput Biol*, **7**, e1002150. [12](#)
- DATTA, R., MEACHAM, C., SAMAD, B., NEYER, C. & SJÖLANDER, K. (2009). Berkeley PHOG: PhyloFacts orthology group prediction web server. *Nucleic Acids Res*, **37**, W84–9. [10](#)
- DEHAL, P. & BOORE, J. (2005). Two rounds of whole genome duplication in the ancestral vertebrate. *PLoS Biol*, **3**, e314. [18](#)
- DEMUTH, J., DE BIE, T., STAJICH, J., CRISTIANINI, N. & HAHN, M. (2006). The evolution of mammalian gene families. *PLoS One*, **1**, e85. [19](#)
- EDDY, S. (2009). A new generation of homology search tools based on probabilistic inference. *Genome Inform*, **23**, 205–11. [9](#)

BIBLIOGRAPHY

- EDGAR, R. (2004). MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*, **32**, 1792–7. [14](#)
- ENCODE PROJECT CONSORTIUM (2007). Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature*, **447**, 799–816. [44](#), [50](#)
- FLICEK, P., AMODE, M., BARRELL, D., BEAL, K., BRENT, S., CHEN, Y., CLAPHAM, P., COATES, G., FAIRLEY, S., FITZGERALD, S., GORDON, L., HENDRIX, M., HOURLIER, T., JOHNSON, N., KÄHÄRI, A., KEEFE, D., KEENAN, S., KINSELLA, R., KOKOCINSKI, F., KULESHA, E., LARSSON, P., LONGDEN, I., McLAREN, W., OVERDUIN, B., PRITCHARD, B., RIAT, H., RIOS, D., RITCHIE, G., RUFFIER, M., SCHUSTER, M., SOBRAL, D., SPUDICH, G., TANG, Y., TREVANION, S., VANDROVCOVA, J., VILELLA, A., WHITE, S., WILDER, S., ZADISSA, A., ZAMORA, J., AKEN, B., BIRNEY, E., CUNNINGHAM, F., DUNHAM, I., DURBIN, R., FERNÁNDEZ-SUAREZ, X., HERRERO, J., HUBBARD, T., PARKER, A., PROCTOR, G., VOGEL, J. & SEARLE, S. (2011). Ensembl 2011. *Nucleic Acids Res*, **39**, D800–6. [12](#)
- GREEN, P. (2007). 2x genomes—does depth matter? *Genome Res*, **17**, 1547–9. [46](#)
- HAHN, M., HAN, M. & HAN, S. (2007). Gene family evolution across 12 drosophila genomes. *PLoS Genet*, **3**, e197. [10](#)
- HEGER, A. & PONTING, C. (2008). OPTIC: orthologous and paralogous transcripts in clades. *Nucleic Acids Res*, **36**, D267–70. [10](#), [14](#), [29](#), [33](#), [39](#), [41](#), [42](#)
- HUBBARD, T., AKEN, B., BEAL, K., BALLESTER, B., CACCAMO, M., CHEN, Y., CLARKE, L., COATES, G., CUNNINGHAM, F., CUTTS, T., DOWN, T., DYER, S., FITZGERALD, S., FERNANDEZ-BANET, J., GRAF, S., HAIDER, S., HAMMOND, M., HERRERO, J., HOLLAND, R., HOWE, K., HOWE, K., JOHNSON, N., KAHARI, A., KEEFE, D., KOKOCINSKI, F., KULESHA, E., LAWSON, D., LONGDEN, I., MELSOOPP, C., MEGY, K., MEIDL, P., OUVERDIN, B., PARKER, A., PRLIC, A., RICE, S., RIOS, D., SCHUSTER, M., SEALY, I., SEVERIN, J., SLATER, G., SMEDLEY, D., SPUDICH, G., TREVANION, S., VILELLA, A., VOGEL, J., WHITE, S., WOOD, M., COX, T., CURWEN, V., DURBIN, R., FERNANDEZ-SUAREZ, X., FLICEK, P., KASPRZYK, A., PROCTOR, G., SEARLE, S., SMITH, J., URETA-VIDAL, A. & BIRNEY, E. (2007). Ensembl 2007. *Nucleic Acids Res*, **35**, D610–7. [11](#), [50](#)

BIBLIOGRAPHY

- HUBISZ, M., LIN, M., KELLIS, M. & SIEPEL, A. (2011). Error and error mitigation in low-coverage genome assemblies. *PLoS One*, **6**, e17034. [12](#), [50](#), [51](#)
- JAFFE, D., BUTLER, J., GNERRE, S., MAUCELI, E., LINDBLAD-TOH, K., MESIROV, J., ZODY, M. & LANDER, E. (2003). Whole-genome sequence assembly for mammalian genomes: Arachne 2. *Genome Res*, **13**, 91–6. [50](#)
- JAILLON, O., AURY, J., BRUNET, F., PETIT, J., STANGE-THOMANN, N., MAUCELI, E., BOUNEAU, L., FISCHER, C., OZOUF-COSTAZ, C., BERNOT, A., NICAUD, S., JAFFE, D., FISHER, S., LUTFALLA, G., DOSSAT, C., SEGURENS, B., DASILVA, C., SALANOBAT, M., LEVY, M., BOUDET, N., CASTELLANO, S., ANTHOUARD, V., JUBIN, C., CASTELLI, V., KATINKA, M., VACHERIE, B., BIÉMONT, C., SKALLI, Z., CATTOLICO, L., POULAIN, J., DE BERARDINIS, V., CRUAUD, C., DUPRAT, S., BROTTIER, P., COUTANCEAU, J., GOUZY, J., PARRA, G., LARDIER, G., CHAPPLE, C., MCKERNAN, K., MC EWAN, P., BOSAK, S., KELLIS, M., VOLFF, J., GUIGÓ, R., ZODY, M., MESIROV, J., LINDBLAD-TOH, K., BIRREN, B., NUSBAUM, C., KAHN, D., ROBINSON-RECHAVI, M., LAUDET, V., SCHACHTER, V., QUÉTIER, F., SAURIN, W., SCARPELLI, C., WINCKER, P., LANDER, E., WEISSENBACH, J. & ROEST CROLIUS, H. (2004). Genome duplication in the teleost fish *tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. *Nature*, **431**, 946–57. [27](#)
- JUN, J., MANDOIU, I. & NELSON, C. (2009). Identification of mammalian orthologs using local synteny. *BMC Genomics*, **10**, 630. [9](#)
- KATOH, K., KUMA, K., TOH, H. & MIYATA, T. (2005). MAFFT version 5: improvement in accuracy of multiple sequence alignment. *Nucleic Acids Res*, **33**, 511–8. [14](#)
- KELLIS, M., BIRREN, B. & LANDER, E. (2004). Proof and evolutionary analysis of ancient genome duplication in the yeast *saccharomyces cerevisiae*. *Nature*, **428**, 617–24. [14](#)
- KIMURA, M. & OHTA, T. (1974). On some principles governing molecular evolution. *Proceedings of the National Academy of Sciences*, **71**, 2848–2852. [44](#)
- KOONIN, E. & WOLF, Y. (2006). Evolutionary systems biology: links between gene evolution and function. *Curr Opin Biotechnol*, **17**, 481–7. [13](#), [44](#)
- KOONIN, E., FEDOROVA, N., JACKSON, J., JACOBS, A., KRYLOV, D., MAKAROVA, K., MAZUMDER, R., MEKHEDOV, S., NIKOLSKAYA, A., RAO, B., ROGOZIN, I.,

BIBLIOGRAPHY

- SMIRNOV, S., SOROKIN, A., SVERDLOV, A., VASUDEVAN, S., WOLF, Y., YIN, J. & NATALE, D. (2004). A comprehensive evolutionary classification of proteins encoded in complete eukaryotic genomes. *Genome Biol*, **5**, R7. [16](#)
- KOONIN, E.V., MAKAROVA, K.S. & ARAVIND, L. (2001). HORIZONTAL GENE TRANSFER IN PROKARYOTES: Quantification and classification1. *Annual Review of Microbiology*, **55**, 709–742. [9](#)
- KOSIOL, C., HOLMES, I. & GOLDMAN, N. (2007). An empirical codon model for protein sequence evolution. *Mol Biol Evol*, **24**, 1464–79. [54](#), [60](#), [61](#)
- LASSMANN, T., FRINGS, O. & SONNHAMMER, E. (2009). Kalign2: high-performance multiple alignment of protein and nucleotide sequences allowing external features. *Nucleic Acids Res*, **37**, 858–65. [14](#)
- LINDBLAD-TOH, K., WADE, C.M., MIKKELSEN, T.S., KARLSSON, E.K., JAFFE, D.B., KAMAL, M., CLAMP, M., CHANG, J.L., KULBOKAS, E.J., ZODY, M.C., MAUCELI, E., XIE, X., BREEN, M., WAYNE, R.K., OSTRANDER, E.A., PONTING, C.P., GALIBERT, F., SMITH, D.R., DEJONG, P.J., KIRKNESS, E., ALVAREZ, P., BIAGI, T., BROCKMAN, W., BUTLER, J., CHIN, C.W., COOK, A., CUFF, J., DALY, M.J., DECAPRIO, D., GNERRE, S., GRABHERR, M., KELLIS, M., KLEBER, M., BARDELEBEN, C., GOODSTADT, L., HEGER, A., HITTE, C., KIM, L., KOEPFLI, K.P., PARKER, H.G., POLLINGER, J.P., SEARLE, S.M.J., SUTTER, N.B., THOMAS, R., WEBBER, C., BALDWIN, J., BROAD SEQUENCING PLATFORM MEMBERS & LANDER, E.S. (2005). Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature*, **438**, 803–819. [43](#)
- LINDBLAD-TOH, K., GARBER, M., ZUK, O. & ,ET AL. (64 CO-AUTHORS) (2011). A high-resolution map of evolutionary constraint in the human genome based on 29 eutherian mammals (in press). *Nature*, **0**, 0. [44](#)
- LÖYTINOJA, A. & GOLDMAN, N. (2008). Phylogeny-aware gap placement prevents errors in sequence alignment and evolutionary analysis. *Science*, **320**, 1632–5. [54](#)
- LYNCH, M. & CONERY, J. (2000). The evolutionary fate and consequences of duplicate genes. *Science*, **290**, 1151–5. [9](#), [52](#)
- MACKENZIE, D., DEFERNÉZ, M., DUNN, W., BROWN, M., FULLER, L., DE HERRERA, S., GÜNTHER, A., JAMES, S., EAGLES, J., PHILO, M., GOODACRE, R. & ROBERTS,

BIBLIOGRAPHY

- I. (2008). Relatedness of medically important strains of *saccharomyces cerevisiae* as revealed by phylogenetics and metabolomics. *Yeast*, **25**, 501–12. [14](#)
- MALLICK, S., GNERRE, S., MULLER, P. & REICH, D. (2009). The difficulty of avoiding false positives in genome scans for natural selection. *Genome Res*, **19**, 922–33. [48](#), [49](#)
- MARGULIES, E., VINSON, J., NISC COMPARATIVE SEQUENCING PROGRAM, MILLER, W., JAFFE, D., LINDBLAD-TOH, K., CHANG, J., GREEN, E., LANDER, E., MULLIKIN, J. & CLAMP, M. (2005). An initial strategy for the systematic identification of functional elements in the human genome by low-redundancy comparative sequencing. *Proc Natl Acad Sci U S A*, **102**, 4795–800. [43](#)
- MARGULIES, E., COOPER, G., ASIMENOS, G., THOMAS, D., DEWEY, C., SIEPEL, A., BIRNEY, E., KEEFE, D., SCHWARTZ, A., HOU, M., TAYLOR, J., NIKOLAEV, S., MONTOYA-BURGOS, J., LÖYTYNOJA, A., WHELAN, S., PARDI, F., MASSINGHAM, T., BROWN, J., BICKEL, P., HOLMES, I., MULLIKIN, J., URETA-VIDAL, A., PATEN, B., STONE, E., ROSENBLoom, K., KENT, W., BOUFFARD, G., GUAN, X., HANSEN, N., IDOL, J., MADURO, V., MASKERI, B., McDOWELL, J., PARK, M., THOMAS, P., YOUNG, A., BLAKESLEY, R., MUZNY, D., SODERGREN, E., WHEELER, D., WORLEY, K., JIANG, H., WEINSTOCK, G., GIBBS, R., GRAVES, T., FULTON, R., MARDIS, E., WILSON, R., CLAMP, M., CUFF, J., GNERRE, S., JAFFE, D., CHANG, J., LINDBLAD-TOH, K., LANDER, E., HINRICHs, A., TRUMBOWER, H., CLAWSON, H., ZWEIG, A., KUHN, R., BARBER, G., HARTE, R., KAROLCHIK, D., FIELD, M., MOORE, R., MATTHEWSON, C., SCHEIN, J., MARRA, M., ANTONARAKIS, S., BATZOGLOU, S., GOLDMAN, N., HARDISON, R., HAUSSLER, D., MILLER, W., PACHTER, L., GREEN, E. & SIDOW, A. (2007). Analyses of deep mammalian sequence alignments and constraint predictions for 1% of the human genome. *Genome Res*, **17**, 760–74. [43](#), [44](#)
- MASSINGHAM, T. & GOLDMAN, N. (2005). Detecting amino acid sites under positive selection and purifying selection. *Genetics*, **169**, 1753–62. [44](#), [45](#), [60](#), [62](#), [67](#), [70](#)
- MILINKOVITCH, M., HELAERS, R., DEPIERREUX, E., TZIKA, A. & GABALDÓN, T. (2010). 2x genomes—depth does matter. *Genome Biol*, **11**, R16. [15](#), [27](#), [41](#)
- MILLER, J., KOREN, S. & SUTTON, G. (2010). Assembly algorithms for next-generation sequencing data. *Genomics*, **95**, 315–27. [25](#)

BIBLIOGRAPHY

- MIRONOV, A., FICKETT, J. & GELFAND, M. (1999). Frequent alternative splicing of human genes. *Genome Res*, **9**, 1288–93. [12](#)
- MOSES, A. & DURBIN, R. (2009). Inferring selection on amino acid preference in protein domains. *Mol Biol Evol*, **26**, 527–36. [83](#)
- MOUSE GENOME SEQUENCING CONSORTIUM & MOUSE GENOME ANALYSIS GROUP (2002). Initial sequencing and comparative analysis of the mouse genome. *Nature*, **420**, 520–62. [43](#)
- MULLER, J., SZKLARCZYK, D., JULIEN, P., LETUNIC, I., ROTH, A., KUHN, M., POWELL, S., VON MERING, C., DOERKS, T., JENSEN, L. & BORK, P. (2010). eggNOG v2.0: extending the evolutionary genealogy of genes with enhanced non-supervised orthologous groups, species and functional annotations. *Nucleic Acids Res*, **38**, D190–5. [10](#)
- MURPHY, W., PRINGLE, T., CRIDER, T., SPRINGER, M. & MILLER, W. (2007). Using genomic data to unravel the root of the placental mammal phylogeny. *Genome Res*, **17**, 413–21. [56](#)
- NIELSEN, R. & YANG, Z. (1998). Likelihood models for detecting positively selected amino acid sites and applications to the HIV-1 envelope gene. *Genetics*, **148**, 929–36. [44](#)
- NOTREDAME, C., HIGGINS, D. & HERINGA, J. (2000). T-Coffee: A novel method for fast and accurate multiple sequence alignment. *J Mol Biol*, **302**, 205–17. [14](#)
- OHNO, S. (1970). *Evolution by gene duplication..* London: George Allen & Unwin Ltd. Berlin, Heidelberg and New York: Springer-Verlag. [9](#)
- PARMLEY, J., CHAMARY, J. & HURST, L. (2006). Evidence for purifying selection against synonymous mutations in mammalian exonic splicing enhancers. *Mol Biol Evol*, **23**, 301–9. [12](#)
- POLLARD, K., HUBISZ, M., ROSENBLoom, K. & SIEPEL, A. (2010). Detection of non-neutral substitution rates on mammalian phylogenies. *Genome Res*, **20**, 110–21. [20](#)
- PRUITT, K., HARROW, J., HARTE, R., WALLIN, C., DIEKHANS, M., MAGLOTT, D., SEARLE, S., FARRELL, C., LOVELAND, J., RUEF, B., HART, E., SUNER, M.,

BIBLIOGRAPHY

- LANDRUM, M., AKEN, B., AYLING, S., BAERTSCH, R., FERNANDEZ-BANET, J., CHERRY, J., CURWEN, V., DICUCCIO, M., KELLIS, M., LEE, J., LIN, M., SCHUSTER, M., SHKEDA, A., AMID, C., BROWN, G., DUKHANINA, O., FRANKISH, A., HART, J., MAIDAK, B., MUDGE, J., MURPHY, M., MURPHY, T., RAJAN, J., RAJPUT, B., RIDICK, L., SNOW, C., STEWARD, C., WEBB, D., WEBER, J., WILMING, L., WU, W., BIRNEY, E., HAUSSLER, D., HUBBARD, T., OSTELL, J., DURBIN, R. & LIPMAN, D. (2009). The consensus coding sequence (ccds) project: Identifying a common protein-coding gene set for the human and mouse genomes. *Genome Res*, **19**, 1316–23. [13](#)
- RASMUSSEN, M. & KELLIS, M. (2007). Accurate gene-tree reconstruction by learning gene- and species-specific substitution rates across multiple complete genomes. *Genome Res*, **17**, 1932–42. [14](#)
- RAT GENOME SEQUENCING PROJECT CONSORTIUM (2004). Genome sequence of the brown norway rat yields insights into mammalian evolution. *Nature*, **428**, 493–521. [43](#)
- RATNAKUMAR, A., MOUSSET, S., GLÉMIN, S., BERGLUND, J., GALTIER, N., DURET, L. & WEBSTER, M. (2010). Detecting positive selection within genomes: the problem of biased gene conversion. *Philos Trans R Soc Lond B Biol Sci*, **365**, 2571–80. [48](#)
- REMM, M., STORM, C. & SONNHAMMER, E. (2001). Automatic clustering of orthologs and in-paralogs from pairwise species comparisons. *J Mol Biol*, **314**, 1041–52. [9](#)
- RUAN, J., LI, H., CHEN, Z., COGHLAN, A., COIN, L., GUO, Y., HÉRICHÉ, J., HU, Y., KRISTIANSEN, K., LI, R., LIU, T., MOSES, A., QIN, J., VANG, S., VILELLA, A., URETA-VIDAL, A., BOLUND, L., WANG, J. & DURBIN, R. (2008). TreeFam: 2008 update. *Nucleic Acids Res*, **36**, D735–40. [10](#), [14](#), [15](#)
- SCHNEIDER, A., SOUVOROV, A., SABATH, N., LANDAN, G., GONNET, G. & GRAUR, D. (2009). Estimates of positive darwinian selection are inflated by errors in sequencing, annotation, and alignment. *Genome Biol Evol*, **1**, 114–8. [48](#), [49](#)
- SIEPEL, A., BEJERANO, G., PEDERSEN, J., HINRICHES, A., HOU, M., ROSENBLoom, K., CLAWSON, H., SPIETH, J., HILLIER, L., RICHARDS, S., WEINSTOCK, G., WILSON, R., GIBBS, R., KENT, W., MILLER, W. & HAUSSLER, D. (2005). Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res*, **15**, 1034–50. [14](#)

BIBLIOGRAPHY

- SJÖLANDER, K., DATTA, R., SHEN, Y. & SHOFFNER, G. (2011). Ortholog identification in the presence of domain architecture rearrangement. *Brief Bioinform*, **12**, 413–22. 9
- STORZ, J., HOFFMANN, F., OPAZO, J. & MORIYAMA, H. (2008). Adaptive functional divergence among triplicated alpha-globin genes in rodents. *Genetics*, **178**, 1623–38. 48
- STUDER, R., PENEL, S., DURET, L. & ROBINSON-RECHAVI, M. (2008). Pervasive positive selection on duplicated and nonduplicated vertebrate protein coding genes. *Genome Res*, **18**, 1393–402. 48
- TEYTELMAN, L., OZAYDIN, B., ZILL, O., LEFRANÇOIS, P., SNYDER, M., RINE, J. & EISEN, M. (2009). Impact of chromatin structures on DNA processing for genomic analyses. *PLoS One*, **4**, e6700. 51
- VILELLA, A., SEVERIN, J., URETA-VIDAL, A., HENG, L., DURBIN, R. & BIRNEY, E. (2009). EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Res*, **19**, 327–35. 10, 11, 12, 14, 15, 39
- VILELLA, A., BIRNEY, E., FLICEK, P. & HERRERO, J. (2011). Considerations for the inclusion of 2x mammalian genomes in phylogenetic analyses. *Genome Biol*, **12**, 401. 15
- WALLACE, I., O'SULLIVAN, O., HIGGINS, D. & NOTREDAME, C. (2006). M-Coffee: combining multiple sequence alignment methods with T-Coffee. *Nucleic Acids Res*, **34**, 1692–9. 14
- WANG, Y. & GU, X. (2001). Functional divergence in the caspase gene family and altered functional constraints: statistical analysis and prediction. *Genetics*, **158**, 1311–20. 48
- WHELAN, S. (2008). Spatial and temporal heterogeneity in nucleotide sequence evolution. *Mol Biol Evol*, **25**, 1683–94. 26
- WHELAN, S. & GOLDMAN, N. (2004). Estimating the frequency of events that cause multiple-nucleotide changes. *Genetics*, **167**, 2027–43. 60
- YANG, Z. & NIELSEN, R. (1998). Synonymous and nonsynonymous rate variation in nuclear genes of mammals. *J Mol Evol*, **46**, 409–18. 44
- YUAN, Y., EULENSTEIN, O., VINGRON, M. & BORK, P. (1998). Towards detection of orthologues in sequence databases. *Bioinformatics*, **14**, 285–9. 9

BIBLIOGRAPHY

ZHU, J., HE, F., HU, S. & YU, J. (2008). On the nature of human housekeeping genes.
Trends Genet., **24**, 481–4. [13](#)