# MIPhy: Identify and quantify rapidly evolving members of large gene families

David M Curran [Corresp., 1] , John S Gilleard [2] , James D Wasmuth [Corresp. 1]

[1] Department of Ecosystem and Public Health, Faculty of Veterinary Medicine, University of Calgary, Calgary, AB, Canada

[2] Department of Comparative Biology and Experimental Medicine, Faculty of Veterinary Medicine, University of Calgary, Calgary, AB, Canada

Corresponding Authors: David M Curran, James D Wasmuth
Email address: dmcurran@ucalgary.ca, jwasmuth@ucalgary.ca

After transitioning to a new environment, species often exhibit rapid phenotypic innovation. One of the fastest mechanisms for this is duplication followed by specialization of existing genes, which leaves a phylogenetic signature of lineage-specific expansions and contractions. These can be identified by analyzing the gene family across several species and identifying patterns of gene duplication and loss that do not correlate with the known relationships between those species. This signature, termed phylogenetic instability, has been previously linked to adaptations that change the way an organism samples and responds to its environment; conversely, low phylogenetic instability has been previously linked to proteins with endogenous functions.

Here, we present MIPhy, a method to identify and quantify phylogenetic instability by quantifying the incongruence of a gene's evolutionary history. The motivation behind MIPhy was to produce a tool to aid in interpreting phylogenetic trees. It can predict which members of a gene family are under adaptive evolution, working only from a gene tree and the relationship between the species under consideration.

We demonstrate the usefulness of MIPhy by accurately predicting which members of the mammalian cytochrome P450 gene superfamily metabolize xenobiotics and which metabolize endogenous compounds. Our predictions correlate very well with known substrate specificities of the human enzymes. We also analyze the *Caenorhabditis* collagen gene family and use MIPhy to predict genes that produce an observable phenotype when knocked down in *C. elegans*, and show that our predictions correlate well with existing knowledge. The software can be downloaded and installed from https://github.com/dave-the-scientist/miphy under a BSD 2-clause license. It is also available as an online web tool at http://miphy.wasmuthlab.org.

# MIPhy: identify and quantify rapidly evolving members of large gene families

David M. Curran[1, †], John S. Gilleard[2], James D. Wasmuth[1]

[1] Department of Ecosystem and Public Health, Faculty of Veterinary Medicine, University of Calgary, Calgary, Alberta, Canada

[2] Department of Comparative Biology and Experimental Medicine, Faculty of Veterinary Medicine, University of Calgary, Calgary, Alberta, Canada

[†] Current Address: Department of Molecular Medicine, the Hospital for Sick Children, Toronto, Ontario, Canada

Corresponding authors:

David Curran[1]

Email address: dmcurran@ucalgary.ca

James Wasmuth[1]

Email address: jwasmuth@ucalgary.ca

## Abstract

After transitioning to a new environment, species often exhibit rapid phenotypic innovation. One of the fastest mechanisms for this is duplication followed by specialization of existing genes, which leaves a phylogenetic signature of lineage-specific expansions and contractions. These can be identified by analyzing the gene family across several species and identifying patterns of gene duplication and loss that do not correlate with the known relationships between those species. This signature, termed phylogenetic instability, has been previously linked to adaptations that change the way an organism samples and responds to its environment; conversely, low phylogenetic instability has been previously linked to proteins with endogenous functions.

Here, we present MIPhy, a method to identify and quantify phylogenetic instability by quantifying the incongruence of a gene's evolutionary history. The motivation behind MIPhy was to produce a tool to aid in interpreting phylogenetic trees. It can predict which members of a gene family are under adaptive evolution, working only from a gene tree and the relationship between the species under consideration.

We demonstrate the usefulness of MIPhy by accurately predicting which members of the mammalian cytochrome P450 gene superfamily metabolize xenobiotics and which metabolize endogenous compounds. Our predictions correlate very well with known substrate specificities of the human enzymes. We also analyze the *Caenorhabditis* collagen gene family and use MIPhy to predict genes that produce an observable phenotype when knocked down in *C. elegans*, and show that our predictions correlate well with existing knowledge. The software can be downloaded and installed from https://github.com/dave-the-scientist/miphy under a BSD 2-clause license. It is also available as an online web tool at http://miphy.wasmuthlab.org.

## Introduction

In the absence of specific selective pressures, the phylogeny of a multi-species gene family will tend to agree with the underlying species tree. However, gene events such as gene duplication/loss, horizontal gene transfer (HGT), and incomplete lineage sorting (ILS) – where a polymorphic locus in an ancestral species results in incongruence with the species tree – may become fixed in a species due to evolutionary processes. These events can result in lineage-specific variations in gene family size and incongruence between the gene family phylogeny and the species tree, properties that have collectively been referred to as 'phylogenetic instability' (Thomas 2007). Attempting to work backwards and determine the sequence of events that led from the species tree to the observed gene family is a process called event-inference reconciliation.

It has been hypothesized that the change in environment during a speciation event may lead to higher levels of phylogenetic instability (Lynch and Conery 2000; Zhang 2003; Hurley, Hale, and Prince 2005), especially in genes involved in responding to molecules from the environment (xenobiotics). This has been observed in gene families involved in the immune response (de Bono, Madera, and Chothia 2004; Nei, Gu, and Sitnikova 1997; Su et al. 1999), chemosensory receptors (Niimura and Nei 2005; Thomas et al. 2005), detoxification (Thomas 2007), and host-pathogen interactions (Wasmuth et al. 2012).

Here, we propose using phylogenetic instability to predict the functional roles of the members of a gene family using a new tool, MIPhy (Minimizing Instability in Phylogenetics). Specifically, to identify which family members are under pressure to duplicate and contribute to altered or new functions, with the possibility of new phenotypes. Understanding the effects of these selective pressures is of more than purely theoretical importance; as one example the rapid evolution of

63    drug resistance remains one of the most significant challenges in managing both human

64    (Saunders and Lon 2016) and livestock parasites (Kaplan and Vidyashankar 2012), and the

65    mechanisms underlying these resistant phenotypes is often unknown. We show the usefulness of

66    MIPhy by validating it against two data sets: the cytochrome P450 (*cyp*) genes from ten species

67    of vertebrates, and the collagens from eight species of free-living nematodes. This tool can be

68    used to prioritize genes for further study, for example by predicting the origin of some species-

69    specific function, or identifying essential genes as new therapeutic targets in pathogens. The

70    process to detect phylogenetically unstable genes is two-fold. First, a tree of a large multi-

71    member gene family is split into meaningful clusters - termed minimum instability groups

72    (MIGs) - by incorporating an event-inference model of gene evolution. Second, each MIG is

73    independently scored for phylogenetic instability.

74

75    **Related work**

76    There are several existing algorithms for species/gene tree reconciliation, but none are able to

77    segregate a gene tree into meaningful clusters, quantify the stability of those gene clusters, or

78    score each gene in order to compare and rank the individual family members. CAFE 3 uses a

79    stochastic birth-death model of gene family evolution to infer the size of ancestral families (De

80    Bie et al. 2006; Han et al. 2013). It implements a sampling procedure to determine the statistical

81    significance of those gene families that differ from their expected values, and models the effects

82    of genome assembly and gene annotation errors to provide a more accurate estimate of its

83    evolutionary rates. CAFE 3 uses only the gene family counts without considering the

84    phylogenetic relationships within them, and so would be unable to distinguish inherited paralogs

85    from independently duplicated genes. Further, the algorithm calculates whether an entire gene

86    family is under adaptive evolution, while we are interested in the relative differences between

87    specific clusters of genes within a family. Because of this, it is more suited for large-scale

88    analyses of many gene families at once.

89    BadiRate is similar to CAFE 3, implementing several additional stochastic models of evolution,

90    and providing three statistical frameworks to calculate significance (Librado, Vieira, and Rozas

91    2012). While it allows for more detailed analyses of species traits, it still relies on gene count

92    data and so is unsuitable here for the same reasons as CAFE. It also requires a species tree with

93    meaningful branch lengths, the creation of which is in itself a challenging analysis.

94    NOTUNG (Chen, Durand, and Farach-Colton 2000; Vernot et al. 2008; Stolzer et al. 2012)

95    implements a parsimony-based reconciliation algorithm. It finds the sequence of gene events

96    (gene duplication, gene loss, HGT, and ILS) explaining the differences between the observed

97    gene tree and the underlying species relationships that minimizes a weighted sum. Uniquely

98    amongst other reconciliation methods, it allows for the species or gene tree to be non-binary; as

99    the true history of many species is unclear, polytomies can be useful to describe the current state

100    of knowledge. Important in this consideration is that NOTUNG explicitly models HGT and

101    assumes that ILS is a very rare event, only considering it at polytomies in the gene tree. A recent

102    paper has proposed a similar algorithm, with advances in identifying ILS and HGT (Y. Chan,

103    Ranwez, and Scornavacca 2017). Identifying HGT is a computationally intensive process and is

104    unlikely to play an important role in gene families from multicellular organisms, and we assume

105    that incongruence (as produced by ILS, adaptive evolution, or any other mechanism) is a

106    common enough event to allow throughout the tree (Carstens and Knowles 2007; Mirarab,

107    Bayzid, and Warnow 2016; Scally et al. 2012). RANGER-DTL is another reconciliation method,

108    and has been reported to be 1,000-1,000,000x faster than software like NOTUNG (Bansal, Alm,

109    and Kellis 2012). Unfortunately, this model proved unsuitable as it too does not allow for

110    incongruence events.

111    There are also several probabilistic reconciliation methods available (Rasmussen and Kellis

112    2007; Rasmussen and Kellis 2011; Ma et al. 2008; Doyon et al. 2010; Doyon, Hamel, and

113    Chauve 2012). While these models make use of more sophisticated models of evolution, they are

114    far more computationally intensive and are only applicable to species for which speciation times

115    and/or ancestral population size estimates are available, which is not the case for most species.

116    PHYLDOG overcomes some of these limitations as it is able to estimate the most likely gene

117    trees, species tree, and evolutionary history of a large number of gene families at once (Boussau

118    et al. 2013). Though it does not explicitly model ILS, the authors state that the algorithms can

119    accommodate it as long as the signal is not too strong. This makes it unsuitable, as we expect

120    gene families involved in direct environmental interactions to have a strong ILS signal. Further,

121    this software is designed to combine the information from many gene families at once, and

122    requires extremely significant computational resources (Chaudhary et al. 2015).

123

124    **Validation**

125    A previous study conducted a detailed analysis of the vertebrate *cyp* gene family (Thomas 2007),

126    and found that enzymes with known xenobiotic substrates (about half of the gene family)

127    exhibited high phylogenetic instability, while those with known endogenous substrates were

128    strikingly phylogenetically stable, with clearly defined orthologous relationships. We validate the

129    accuracy of MIPhy by comparing its predictions to the results of that study. That work relied

130    upon the author's detailed knowledge of the gene family under study, and so was not quantified.

131    As the genomes of an increasing number of species are being made available, manual analysis of

132     large gene families from hundreds of species will become intractable. Further, it is desirable to

133     use an algorithm that is consistent and deterministic.

134     Nematode collagens are a large multi-gene family of structural proteins. The *C. elegans* genome

135     contains 181 collagen genes (The C. elegans Sequencing Consortium 1998), many of which

136     encode for proteins that form a major part of the nematode cuticle, which molts five times in the

137     nematode life-cycle and protects the worm from environmental insult. A combination of high

138     throughput and targeted gene knock-down studies have shown that 28 of these genes are

139     associated with an observable phenotype, ranging from morphological variants to lethality

140     (reviewed in (Page and Johnstone 2007)). Available genome sequences from other

141     *Caenorhabditis* species reveal both conservation and divergence of genes and their role in

142     biochemical pathways (Stein et al. 2003; Fierst et al. 2015; Gilabert et al. 2016). To validate

143     MIPhy's predictions for researchers aiming to prioritize genes for functional characterization, we

144     test whether MIGs with lower phylogenetic instability scores were more likely to contain *C.*

145     *elegans* genes associated with phenotypic changes when knocked-out.

146     While an individual can manually cluster a small tree without much trouble, the large size of

147     some gene families and the ever-expanding availability of sequence data mean that this will

148     quickly become intractable. There are several software packages used to automatically cluster a

149     phylogenetic tree, but because of the ill-defined nature of clustering problems in general, the

150     methods generally come to different conclusions on the same data sets. We are aware of no

151     method that is targeted towards multi-species gene families, which means that none make use of

152     problem-specific information such as an event-inference model of gene evolution. The clustering

153     algorithm described here combines the similarity between each gene with the most parsimonious

154     explanation of gene events, to predict the ancestry of each observed member of the gene family.

## Methods

**Running MIPhy on a large phylogeny**

The NCBI genome database (https://www.ncbi.nlm.nih.gov/assembly/organism/) was filtered for

all animal genomes that were at a 'Chromosome' or 'Complete' level of assembly on July 26,

2016, yielding 98 hits. When there were multiple genome assemblies for a single species, only

that with the highest number of annotated proteins was kept. Finally, the *Bos indicus, Capra*

*aegagrus, Mus spretus, Nasalis larvatus,* and *Nomascus leucogenys* genomes were discarded as

they were judged to contain too few protein sequences to have reliable annotations (all had fewer

than 1,500). All protein sequences for the remaining 58 species were concatenated into one file,

which was queried with the 628 vertebrate Cyp proteins from (Thomas 2007) using BlastP

(Camacho et al. 2009), and resulting in 5,498 hits with an E-value $< 10^{-10}$. We note that this is

not a particularly rigorous procedure; some of these sequences may not actually be Cyp proteins,

and we may have missed some true hits. However, the purpose of this procedure was to generate

a very large and representative phylogeny as a test case for MIPhy, not to comment on animal

Cyps themselves.

The sequences were aligned using Clustal Omega (Sievers et al. 2011) with the command:

clustalo -i INPUT_FILE.fa --threads 10 --log INPUT_FILE-clustalO.log -v --force --use-kimura

--iter 10 -o OUT_FILE

The columns of this alignment with <75% gaps were used to build a phylogenetic tree using

RAxML (Stamatakis 2014) with the command:

raxml -s INPUT_FILE.phylip -T 10 -# 5 -m PROTGAMMAWAG -j -p 12345 -n OUT_FILE

177 **Analysis of nematode collagen genes**

178 From Wormbase (Howe et al. 2016), there are 157 genes from *C. elegans* annotated with the

179 gene class 'col'. To these we added the 19 genes listed in (Page and Johnstone 2007). A further

180 five were found by searching for the repetitive Gly-X-Y amino acid motif and checking each

181 entry in WormBase. Phenotype data from gene knock-down studies is available from Wormbase.

182 The protein sequences of *C. angaria, C. brenneri, C. briggsae, C. japonica, C. remanei, C.*

183 *sinica* and *C. tropicalis* were downloaded from Wormbase (version WS259). We searched the

184 181 *C. elegans* collagens against these protein sets using BLASTP (Camacho et al. 2009) and

185 confirmed the presence of the characteristic and repetitive Gly-X-Y amino acid motif. In

186 instances of different isoforms, we selected the longest for subsequent analysis. In total, 1349

187 genes were collected from the eight species.

188 The diversity of the N- and C-terminus across the collagens, coupled to the variable number of

189 the Gly-X-Y motif, precludes a standard sequence alignment based approach. Therefore, we

190 constructed a distance matrix based on k-mer frequency, using the jD2Stat program (C. X. Chan

191 et al. 2014) with the command:

192 java -Xmx20g -jar jD2Stat_1.0.jar -n 1 -k 8

193 We used the neighbor program with default parameters from the phylip suite to reconstruct the

194 phylogenetic tree (Felsenstein 1989). The species phylogenetic relationships had been previously

195 determined using the ITS-2 genetic barcode (Félix, Braendle, and Cutter 2014). Note that *C. sp.*

196 *5* has since been renamed as *C. sinica* (Huang et al. 2014).

197    When statistically evaluating the instability scores between MIGs with and without observable

198    knock-down phenotypes in *C. elegans*, neither set was normally distributed (via the Shapiro-

199    Wilk test). We therefore used a one-tail Mann-Whitney U test to compare them.

200

**201    Parsimony clustering of the gene tree using a model of gene family evolution**

202    The algorithm described in this work uses a model of gene family evolution derived from the

203    core reconciliation methods of NOTUNG (Chen, Durand, and Farach-Colton 2000; Vernot et al.

204    2008; Stolzer et al. 2012), with some modifications such as allowing incongruence throughout

205    the tree. We do this as apparent incongruence may arise for many reasons: due to errors in

206    sequencing or gene-finding, incompletely resolved branches in tree-building software, horizontal

207    gene transfer, incomplete lineage sorting, or it may be due to selective pressures acting on one or

208    more species. Using our model, each internal node of the gene tree is classified as representing

209    one gene event: duplication, speciation, or incongruence. Gene loss is also considered a gene

210    event, and is quantified at duplication nodes. The algorithm is detailed in Article S1, but

211    summarized here.

212    MIPhy was designed to identify members of a gene family under adaptive evolution, and so must

213    also cluster the given gene tree into MIGs. This is necessary to isolate 'unstable' genes from

214    'stable' genes, and has the effect of assigning all genes from all species in one MIG the same

215    phylogenetic instability score. This score is a function of the model of gene family evolution, and

216    for a given MIG it quantifies all gene events at or below the most recent common ancestor of that

217    group:

218    $$score(g) = \theta_D \cdot D(g) + \theta_I \cdot I(g) + \theta_L \cdot L(g) + \theta_P \cdot P(g),$$

219     where $D(g)$, $I(g)$, and $L(g)$ are the total duplications, incongruence, and loss events within the

220     MIG, respectively; $P(g)$ is a measure of the "relative spread" of the MIG (how dissimilar the

221     sequences are – set to 0 for this phase); and the $\theta$ values are the strictly positive weights applied

222     to each event. Under this definition, the score can be interpreted as a measure of the

223     incongruence experienced by a cluster of genes throughout their evolutionary history.

224     Every node in the gene tree is evaluated in a depth-first post-order traversal; if the node is a leaf

225     a new MIG is defined as containing only that node. At each non-leaf node the score function is

226     used to compare merging all of that node's descendants into a single MIG, versus leaving them

227     with their existing clustering pattern. This is the initial clustering phase, and it generates a

228     preliminary clustering pattern.

229

230     **Cluster refinement**

231     This initial clustering pattern arises from the most parsimonious history of gene events required

232     to reconcile $T_G$ with $T_S$. It indicates which groups of genes, under this model, for the given

233     weights and while disregarding all branch lengths in $T_G$, most probably evolved from a single

234     homologue in an ancestral species. This second phase of the algorithm refines these predictions

235     by incorporating branch length information, specifically the pairwise distance information

236     between the sequences. If a sequence in the gene tree is separated by an uncommonly large

237     phylogenetic distance from its closest MIG, there should be a cost associated with the decision to

238     include it in that MIG.

239     This is accomplished by the "relative spread" term $P(g)$ in the score function, which measures

240     the spread within a cluster. It is a measure of how "good" a cluster is compared to the others:

241 $$P(g) = \frac{\sigma(g)}{\overline{\sigma}} - 1,$$

242 where $\sigma(g)$ is the standard deviation of the points representing the sequences in the MIG rooted

243 by $g$, and $\overline{\sigma}$ is the median standard deviation of all MIGs (excluding singleton clusters). The

244 spread quantity is normalized around 0, so $P(g) = 1.0$ indicates that the spread of MIG $g$ is 100%

245 larger than the median spread, while $P(h) = -0.3$ indicates that the spread of MIG $h$ is 30%

246 smaller than $\overline{\sigma}$. Many clustering metrics, including this one, can only be calculated from data in a

247 coordinate space, and so we first transform the phylogenetic tree into a set of points using multi-

248 dimensional scaling (Torgerson 1952) (see Article S1 for implementation details). Standard

249 deviation is used as a measure of the pairwise branch lengths within a MIG, and because it is

250 widely used and easily understood, but clustering-specific methods like the Davies-Bouldin

251 index (Davies and Bouldin 1979) or silhouette (Rousseeuw 1987) could be easily substituted. As

252 in the initial clustering phase, each node $g$ in $T_G$ is again visited in turn. The clustering procedure

253 is repeated, this time using the full score function.

254

255

256 **Results**

257 **Program input, workflow, and interface**

258 This software requires two input files: the gene tree in Newick format, and an information file

259 that contains the species tree (topology only; no branch lengths) as well as the assignment of

260 each sequence to one species. MIPhy is agnostic to the method used to generate the tree, and can

261 be used to analyze those produced from nucleotides, amino acids, or any other features. The

262 cluster analysis algorithm is written in Python and a local daemon server is started along with an

263 HTML document to display the results. This page has interactive controls and communicates

264 directly with the Python server, allowing the user to reanalyze their data and see the effects of

265 modifying any of the parameters in real time.

266 The visualization page displays the gene tree clustered into MIGs, the current parameter values,

267 summary statistics, and a sortable list of the MIGs (Fig. 1). Selecting a specific sequence or MIG

268 will provide additional details. The page also contains a usage description, and provides options

269 to modify visual elements like font sizes, the tree size, and the colour of each element. The tree

270 and legend can be exported and saved as an SVG image file, or the clustering pattern and

271 instability scores from one or more species can be exported and saved as a CSV file.

272 MIPhy was used to analyze a dataset of annotated vertebrate Cyp proteins, which consists of 628

273 sequences from 10 species (Thomas 2007). The algorithm calculated the optimal clustering

274 pattern in less than 0.2 seconds on a 2.7 GHz laptop. Loading the results in a web browser

275 required ~5 seconds. Modifying parameter weights causes the clustering analysis to be rerun, and

276 redrawing the new results is sped up as only a subset of the page elements need to be modified or

277 recreated (<1 second). To determine how MIPhy will scale to cope with the ever-increasing

278 number of genome sequences, we analyzed a tree of 5,498 Cyp protein sequences from 58

279  animal species. MIPhy completed the initial clustering phase in 30 seconds, the optional cluster

280  refinement phase in 7 minutes, and loaded the results in a web browser in 1.5 minutes.

281

282  **Phylogenetic instability of human Cyp proteins**

283  MIPhy was run with default parameters on the Cyp phylogenetic tree from (Thomas 2007), and

284  the 59 scores from human sequences were extracted and graphed (Fig. 2). These scores fell into

285  two broad categories: 31 were unstable with scores in the interval [18.2, 97.5], and 28 were

286  stable with scores in [0.1, 10.9]. Of the stable sequences, 23 had low scores in [0.1, 5.7], and the

287  remaining 5 had intermediate scores in [7.8, 10.9].

288  Among the MIGs with intermediate scores, Cyp-11B1 (steroid 11β-hydroxylase) and Cyp-11B2

289  (aldosterone synthase) appear to have been recently duplicated in the terrestrial vertebrates, and

290  likely played a role in the ancient transition from sea to land (Colombo et al. 2006). Their

291  instability score is elevated because rats appear to have two additional genes in that cluster and

292  no homologs were found in chicken or frog. It is unclear whether they are actually lost in these

293  species or simply absent from the assemblies.

294

295  **Parameter impact**

296  The default MIPhy weight values are set at 1, 1, 0.5, and 1, for duplications, loss, incongruence,

297  and spread, respectively. These have performed well in our testing and analyses. The effects of

298  modifying these values are considered in terms of the clustering pattern – which indicates which

299  sequences are clustered together – and the cluster rankings – which indicates the instability score

300  of each MIG relative to the others. Increasing the weight for gene loss had very little effect; at

301    even triple its default value it only caused four small MIGs out of the 47 from the vertebrate Cyp

302    tree to be merged with their sister groups. Decreasing the gene duplication weight had much the

303    same effect, causing five MIGs to be merged when it was set to 1/3 of the default value.

304    Increasing the weights for duplication and loss together had no effect on the clustering pattern,

305    and very minimal effect on the cluster rankings. Decreasing both weights together had the same

306    effect as increasing the spread weight, which tended to break up larger MIGs. Decreasing the

307    spread weight to zero had minimal impact, only merging two singleton groups with their

308    neighbors. Decreasing the incongruence weight had no effect, and increasing it had little impact

309    until it became very high, at which point it tended to break up groups.

310

311    **Phylogenetic instability of *Caenorhabditis* collagens**

312    Across the eight species of *Caenorhabditis*, we found 1349 collagen genes (Table S1). The

313    characteristic Gly-X-Y repeat domain can vary greatly in length, presenting a problem for usual

314    alignment guided phylogenetics. To overcome this, we used a k-mer based distance matrix (C. X.

315    Chan et al. 2014). Default settings were used to cluster the protein phylogeny and subsequently

316    score each cluster's phylogenetic instability (Fig. 3). A total of 244 MIGs were generated, with

317    41 MIGs containing proteins from all eight species, 60 MIGs covering any seven of the species,

318    and 151 MIGs containing at least one protein from *C. elegans*. Twenty-five of the 151 MIGs that

319    contained a *C. elegans* protein encoded by a gene whose knock-down is associated with an

320    observable phenotype. The distribution of scores from these 25 MIGs was significantly smaller

321    than the remaining 126 MIGs (medians=2.02 and 3.22; U=991; p=0.002).

322

323 **Discussion**

324 Positive selection, pseudogenization, and the presence of tandem gene arrays are characteristic of

325 rapidly evolving genes, such as those involved in xenobiotic interactions (Thomas 2007). Even

326 though MIPhy's analysis does not incorporate any of this information, every human Cyp

327 sequence with these characteristics received a high instability score (Fig. 2). These predictions

328 appear to extend to the functional role of the enzymes as well, as MIPhy performed very well at

329 classifying the human Cyp proteins into those primarily acting on xenobiotic or endogenous

330 substrates. All enzymes with known endogenous functions had low scores, while all but two with

331 primarily xenobiotic substrates had high instability scores; these exceptions were Cyp-1A1 and

332 Cyp-1A2. While the latter is one of the most important human enzymes involved in xenobiotic

333 metabolism, it has been suggested that both also have important endogenous roles (Zhou et al.

334 2009; Kapitulnik and Gonzalez 1993), which may have shaped their evolutionary history in the

335 vertebrate species studied here.

336 The predictions can be extended to species for which detailed substrate specificity information is

337 limited. The sequences from terrestrial species in the MIG containing human Cyp-27A1 appear

338 stable, but those of the aquatic or amphibious species do not. This observation suggests that these

339 paralogs may play some role specific to aquatic environments. A similar observation can be

340 made about the cluster containing human Cyp-2W1. It has the second-highest instability score,

341 and of the 43 total sequences there is only one each from human, macaque, mouse, and cow.

342 There are 16 from frog, 10 from zebrafish, and 4 from pufferfish, which would suggest that these

343 paralogs may also have evolved to metabolize substrates specific to an aquatic environment, and

344 that this capacity was lost in terrestrial species.

345    The collagens are a large gene family encoding for structural proteins. Most members have been

346    investigated in the past using gene knock-down assays in *C. elegans*, resulting in observed

347    phenotypic changes for approximately 15%. While not all the remaining genes have been

348    investigated, many have, suggesting wide-spread functional redundancy. Using MIPhy to cluster

349    and score the collagen gene phylogeny showed that we could prioritize genes for detailed

350    functional assays, and that we could predict those with observable knock-down phenotypes.

351    Further, this demonstrated that MIPhy is agnostic to the methods or underlying characters used

352    to construct a gene tree, and so is applicable to a wide range of data.

353    The predictions from these analyses would be complimentary to a between-genes positive

354    selection analysis, which is the most commonly used measure of adaptive evolution. While a

355    codon-based positive selection test measures the patterns of sequence variation, phylogenetic

356    instability combines the relative sequence variation between species (from the cluster spread and

357    incongruence events) with the most likely history of duplications and losses. However, MIPhy

358    does have its limitations. It is very sensitive to the given gene tree, and does not currently

359    incorporate any measures of uncertainty such as bootstrapping. There are also exceptions to the

360    assumption that phylogenetic instability is a hallmark of adaptive evolution; the most well-

361    known may be the beta-globin genes that form part of hemoglobin. These genes exhibit sequence

362    polymorphism within and between human populations, lineage-specific expansions and

363    contractions in gene cluster size, and yet continue to play a very vital endogenous role (Hill and

364    Wainscoat 1986; Opazo, Hoffmann, and Storz 2008).

365    An additional use of MIPhy is in the naming of genes, specifically towards generating

366    hierarchical naming conventions using an evolutionary framework. Because a sequence identity

367    threshold was used when annotating Cyp proteins, one may reasonably assume that Cyp-3A4 and

368     Cyp-3A5 have related functions, as they are likely closely related. Conversely, no such

369     assumptions may be made about many other gene families, whose members have often been

370     annotated in order of discovery. This can pose a problem with the discovery of novel genes. If

371     two species possess the example genes *pqr-21* and *pqr-22*, and one of them additionally

372     possesses a paralog to *pqr-21*, this paralog will be named with the next available number;

373     perhaps *pqr-42*. This single tiered naming system does not accommodate any way to suggest that

374     *pqr-21* and *pqr-42* are related to each other. We propose that a phylogenetic analysis like MIPhy

375     could be used to cluster such a gene family into sub-families, and that these clusters could be

376     used to inform a multi-tiered naming system that is better able to accommodate newly discovered

377     gene members. This is an issue that is going to arise more often as increasing numbers of species

378     are being sequenced.

379

380

## Conclusions

This work presents, to our knowledge, the first algorithm for simultaneous reconciliation and clustering of large gene families. MIPhy's instability score has proven to be a valuable tool in identifying the members of gene families that exhibit characteristics of adaptive evolution, predicting collagens that play an important functional role in *C. elegans*, and agrees very well with the known substrate specificity of human Cyp enzymes. It is a useful tool to gain an understanding of the evolution of large gene families, and to generate hypotheses about the potential functional roles of both the stable and unstable sequences.

## Acknowledgements

391 

392 We thank everyone that has tested the software during its development. We also wish to thank

393 Dr. Dannie Durand for her work on reconciliation algorithms, and for discussions on the nuances

394 and applications of this work. Finally, we are grateful to Dr. James Thomas for providing us with

395 his past data so that MIPhy could be validated on his published work.

396 **References**

397 Bansal, Mukul S., Eric J. Alm, and Manolis Kellis. 2012. "Efficient Algorithms for the

398       Reconciliation Problem with Gene Duplication, Horizontal Transfer and Loss."

399       *Bioinformatics* 28 (12): i283–91. doi:10.1093/bioinformatics/bts225.

400 Bie, Tijl De, Nello Cristianini, Jeffery P. Demuth, and Matthew W. Hahn. 2006. "CAFE: A

401       Computational Tool for the Study of Gene Family Evolution." *Bioinformatics* 22 (10):

402       1269–71. doi:10.1093/bioinformatics/btl097.

403 Bono, Bernard de, Martin Madera, and Cyrus Chothia. 2004. "VH Gene Segments in the Mouse

404       and Human Genomes." *Journal of Molecular Biology* 342 (1): 131–43.

405       doi:10.1016/j.jmb.2004.06.055.

406 Boussau, Bastien, Gergely J. Szollosi, Laurent Duret, Manolo Gouy, Eric Tannier, and Vincent

407       Daubin. 2013. "Genome-Scale Coestimation of Species and Gene Trees." *Genome*

408       *Research* 23 (2): 323–30. doi:10.1101/gr.141978.112.

409 Camacho, C, G Coulouris, V Avagyan, N Ma, J Papadopoulos, K Bealer, and T L Madden.

410       2009. "BLAST plus: Architecture and Applications." *BMC Bioinformatics* 10 (421): 1.

411       doi:Artn 421\nDoi 10.1186/1471-2105-10-421.

412 Carstens, Bryan C, and L Lacey Knowles. 2007. "Estimating Species Phylogeny from Gene-Tree

413       Probabilities despite Incomplete Lineage Sorting: An Example from Melanoplus

414       Grasshoppers." *Systematic Biology* 56 (3): 400–411. doi:10.1080/10635150701405560.

415 Chan, Cheong Xin, Guillaume Bernard, Olivier Poirion, James M. Hogan, and Mark A. Ragan.

416       2014. "Inferring Phylogenies of Evolving Sequences without Multiple Sequence

417       Alignment." *Scientific Reports* 4: 6504. doi:10.1038/srep06504.

418    Chan, Yao-ban, Vincent Ranwez, and Céline Scornavacca. 2017. "Inferring Incomplete Lineage

419        Sorting, Duplications, Transfers and Losses with Reconciliations." *Journal of Theoretical*

420        *Biology* 432: 1–13. doi:10.1016/j.jtbi.2017.08.008.

421    Chaudhary, Ruchi, Bastien Boussau, J. Gordon Burleigh, and David Fernández-Baca. 2015.

422        "Assessing Approaches for Inferring Species Trees from Multi-Copy Genes." *Systematic*

423        *Biology* 64 (2): 325–39. doi:10.1093/sysbio/syu128.

424    Chen, K, D Durand, and M Farach-Colton. 2000. "NOTUNG: A Program for Dating Gene

425        Duplications and Optimizing Gene Family Trees." *J Comp Biol* 7 (3–4): 429–47.

426        doi:10.1089/106652700750050871.

427    Colombo, Lorenzo, L. Dalla Valle, C. Fiore, D. Armanini, and P. Belvedere. 2006. "Aldosterone

428        and the Conquest of Land." *Journal of Endocrinological Investigation*.

429    Davies, D L, and D W Bouldin. 1979. "A Cluster Separation Measure." *IEEE Transactions on*

430        *Pattern Analysis and Machine Intelligence* 1 (2): 224–27.

431        doi:10.1109/TPAMI.1979.4766909.

432    Doyon, Jean Philippe, Sylvie Hamel, and Cedric Chauve. 2012. "An Efficient Method for

433        Exploring the Space of Gene Tree/species Tree Reconciliations in a Probabilistic

434        Framework." *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 9 (1):

435        26–39. doi:10.1109/TCBB.2011.64.

436    Doyon, Jean Philippe, Celine Scornavacca, K. Yu Gorbunov, Gergely J. Szöllosi, Vincent

437        Ranwez, and Vincent Berry. 2010. "An Efficient Algorithm for Gene/species Trees

438        Parsimonious Reconciliation with Losses, Duplications and Transfers." In *Lecture Notes in*

439        *Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture*

440     *Notes in Bioinformatics)*, 6398 LNBI:93–108. doi:10.1007/978-3-642-16181-0_9.

441     Félix, Marie Anne, Christian Braendle, and Asher D. Cutter. 2014. "A Streamlined System for

442     Species Diagnosis in Caenorhabditis (Nematoda: Rhabditidae) with Name Designations for

443     15 Distinct Biological Species." *PLoS ONE* 9 (4). doi:10.1371/journal.pone.0094723.

444     Felsenstein, J. 1989. "PHYLIP—Phylogeny Inference Package (Version 3.2)." *Cladistics* 5 (2):

445     163–66. doi:10.1111/j.1096-0031.1989.tb00562.x.

446     Fierst, Janna L., John H. Willis, Cristel G. Thomas, Wei Wang, Rose M. Reynolds, Timothy E.

447     Ahearne, Asher D. Cutter, and Patrick C. Phillips. 2015. "Reproductive Mode and the

448     Evolution of Genome Size and Structure in Caenorhabditis Nematodes." *PLOS Genet* 11

449     (6): e1005323. doi:10.1371/journal.pgen.1005323.

450     Gilabert, Aude, David M. Curran, Simon C. Harvey, James D. Wasmuth, A Pires-daSilva, RJ

451     Sommer, RC Cassada, et al. 2016. "Expanding the View on the Evolution of the Nematode

452     Dauer Signalling Pathways: Refinement through Gene Gain and Pathway Co-Option." *BMC*

453     *Genomics* 17 (1): 476. doi:10.1186/s12864-016-2770-7.

454     Han, Mira V., Gregg W C Thomas, Jose Lugo-Martinez, and Matthew W. Hahn. 2013.

455     "Estimating Gene Gain and Loss Rates in the Presence of Error in Genome Assembly and

456     Annotation Using CAFE 3." *Molecular Biology and Evolution* 30 (8): 1987–97.

457     doi:10.1093/molbev/mst100.

458     Hill, A V, and J S Wainscoat. 1986. "The Evolution of the Alpha- and Beta-Globin Gene

459     Clusters in Human Populations." *Human Genetics* 74 (1). Germany: 16–23.

460     Howe, Kevin L., Bruce J. Bolt, Scott Cain, Juancarlos Chan, Wen J. Chen, Paul Davis, James

461       Done, et al. 2016. "WormBase 2016: Expanding to Enable Helminth Genomic Research."

462       *Nucleic Acids Research* 44 (D1): D774–80. doi:10.1093/nar/gkv1217.

463    Huang, Ren-E, Xiaoliang Ren, Yifei Qiu, and Zhongying Zhao. 2014. "Description of

464       Caenorhabditis Sinica Sp. N. (Nematoda: Rhabditidae), a Nematode Species Used in

465       Comparative Biology for C. Elegans." *PloS One* 9 (11): e110957.

466       doi:10.1371/journal.pone.0110957.

467    Hurley, I., M. E. Hale, and V. E. Prince. 2005. "Duplication Events and the Evolution of

468       Segmental Identity." In *Evolution and Development*, 7:556–67. doi:10.1111/j.1525-

469       142X.2005.05059.x.

470    Kapitulnik, J, and F J Gonzalez. 1993. "Marked Endogenous Activation of the CYP1A1 and

471       CYP1A2 Genes in the Congenitally Jaundiced Gunn Rat." *Mol.Pharmacol.* 43 (5): 722–25.

472       http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation

473       &list_uids=8502229%5Cnpapers2://publication/uuid/296206A2-54D7-44C5-AEB9-

474       8FC88A45018D.

475    Kaplan, Ray M., and Anand N. Vidyashankar. 2012. "An Inconvenient Truth: Global Worming

476       and Anthelmintic Resistance." *Veterinary Parasitology* 186 (1–2): 70–78.

477       doi:10.1016/j.vetpar.2011.11.048.

478    Librado, P., F. G. Vieira, and J. Rozas. 2012. "BadiRate: Estimating Family Turnover Rates by

479       Likelihood-Based Methods." *Bioinformatics* 28 (2): 279–81.

480       doi:10.1093/bioinformatics/btr623.

481    Lynch, M, and J S Conery. 2000. "The Evolutionary Fate and Consequences of Duplicate

482       Genes." *Science (New York, N.Y.)* 290 (5494): 1151–55.

483    doi:10.1126/science.290.5494.1151.

484    Ma, Jian, Aakrosh Ratan, Brian J Raney, Bernard B Suh, Louxin Zhang, Webb Miller, and David

485        Haussler. 2008. "DUPCAR: Reconstructing Contiguous Ancestral Regions with

486        Duplications." *Journal of Computational Biology : A Journal of Computational Molecular*

487        *Cell Biology* 15 (8): 1007–27. doi:10.1089/cmb.2008.0069.

488    Mirarab, Siavash, Md Shamsuzzoha Bayzid, and Tandy Warnow. 2016. "Evaluating Summary

489        Methods for Multilocus Species Tree Estimation in the Presence of Incomplete Lineage

490        Sorting." *Systematic Biology* 65 (3): 366–80. doi:10.1093/sysbio/syu063.

491    Nei, M, X Gu, and T Sitnikova. 1997. "Evolution by the Birth-and-Death Process in Multigene

492        Families of the Vertebrate Immune System." *Proceedings of the National Academy of*

493        *Sciences of the United States of America* 94 (15): 7799–7806. doi:10.1073/pnas.94.15.7799.

494    Niimura, Y, and M Nei. 2005. "Evolutionary Changes of the Number of Olfactory Receptor

495        Genes in the Human and Mouse Lineages." *Gene* 346: 23–28. doi:S0378-1119(04)00591-8

496        [pii]\r10.1016/j.gene.2004.09.027.

497    Opazo, Juan C, Federico G Hoffmann, and Jay F Storz. 2008. "Genomic Evidence for

498        Independent Origins of Beta-like Globin Genes in Monotremes  and Therian Mammals."

499        *Proceedings of the National Academy of Sciences of the United States of America* 105 (5).

500        United States: 1590–95. doi:10.1073/pnas.0710531105.

501    Page, Antony P, and Iain L Johnstone. 2007. "The Cuticle." In *WormBook*, 1–15.

502        doi:10.1895/wormbook.1.138.1.

503    Rasmussen, Matthew D., and Manolis Kellis. 2007. "Accurate Gene-Tree Reconstruction by

504        Learning Gene- and Species-Specific Substitution Rates across Multiple Complete

505        Genomes." *Genome Research* 17 (12): 1932–42. doi:10.1101/gr.7105007.

506        ———. 2011. "A Bayesian Approach for Fast and Accurate Gene Tree Reconstruction."

507        *Molecular Biology and Evolution* 28 (1): 273–90. doi:10.1093/molbev/msq189.

508    Rousseeuw, Peter J. 1987. "Silhouettes: A Graphical Aid to the Interpretation and Validation of

509        Cluster Analysis." *Journal of Computational and Applied Mathematics* 20 (C): 53–65.

510        doi:10.1016/0377-0427(87)90125-7.

511    Saunders, David, and Chanthap Lon. 2016. "Combination Therapies for Malaria Are Failing-

512        What Next?" *The Lancet Infectious Diseases*. doi:10.1016/S1473-3099(15)00525-3.

513    Scally, Aylwyn, Julien Y Dutheil, LaDeana W Hillier, Gregory E Jordan, Ian Goodhead, Javier

514        Herrero, Asger Hobolth, et al. 2012. "Insights into Hominid Evolution from the Gorilla

515        Genome Sequence." *Nature* 483 (7388): 169–75. doi:10.1038/nature10842.

516    Sievers, Fabian, Andreas Wilm, David Dineen, Toby J Gibson, Kevin Karplus, Weizhong Li,

517        Rodrigo Lopez, et al. 2011. "Fast, Scalable Generation of High-Quality Protein Multiple

518        Sequence Alignments Using Clustal Omega." *Molecular Systems Biology* 7 (1): 539.

519        doi:10.1038/msb.2011.75.

520    Stamatakis, Alexandros. 2014. "RAxML Version 8: A Tool for Phylogenetic Analysis and Post-

521        Analysis of Large Phylogenies." *Bioinformatics* 30 (9): 1312–13.

522        doi:10.1093/bioinformatics/btu033.

523    Stein, Lincoln D., Zhirong Bao, Darin Blasiar, Thomas Blumenthal, Michael R. Brent, Nansheng

524        Chen, Asif Chinwalla, et al. 2003. "The Genome Sequence of Caenorhabditis Briggsae: A

525     Platform for Comparative Genomics." *PLoS Biology* 1 (2).

526     doi:10.1371/journal.pbio.0000045.

527   Stolzer, Maureen, Han Lai, Minli Xu, Deepa Sathaye, Benjamin Vernot, and Dannie Durand.

528     2012. "Inferring Duplications, Losses, Transfers and Incomplete Lineage Sorting with

529     Nonbinary Species Trees." *Bioinformatics* 28 (18): 409–15.

530     doi:10.1093/bioinformatics/bts386.

531   Su, C, I Jakobsen, X Gu, and M Nei. 1999. "Diversity and Evolution of T-Cell Receptor Variable

532     Region Genes in Mammals and Birds." *Immunogenetics* 50 (5–6): 301–8.

533     doi:10.1007/s002510050606.

534   The C. elegans Sequencing Consortium. 1998. "Genome Sequence of the Nematode *C. Elegans*:

535     A Platform for Investigating Biology." *Science (New York, N.Y.)* 282 (5396): 2012–18.

536     doi:10.1126/science.282.5396.2012.

537   Thomas, James H. 2007. "Rapid Birth-Death Evolution Specific to Xenobiotic Cytochrome P450

538     Genes in Vertebrates." *PLoS Genetics* 3 (5): 720–28. doi:10.1371/journal.pgen.0030067.

539   Thomas, James H, Joanna L Kelley, Hugh M Robertson, Kim Ly, and Willie J Swanson. 2005.

540     "Adaptive Evolution in the SRZ Chemoreceptor Families of Caenorhabditis Elegans and

541     Caenorhabditis Briggsae." *Proceedings of the National Academy of Sciences of the United*

542     *States of America* 102 (12): 4476–81. doi:10.1073/pnas.0406469102.

543   Torgerson, Warren S. 1952. "Multidimensional Scaling: I. Theory and Method." *Psychometrika*

544     17 (4): 401–19. doi:10.1007/BF02288916.

545   Vernot, Benjamin, Maureen Stolzer, Aiton Goldman, and Dannie Durand. 2008. "Reconciliation

546    with Non-Binary Species Trees." *Journal of Computational Biology : A Journal of*

547    *Computational Molecular Cell Biology* 15 (8): 981–1006. doi:10.1089/cmb.2008.0092.

548    Wasmuth, James D., Viviana Pszenny, Simon Haile, Emily M. Jansen, Alexandra T. Gast, Alan

549    Sher, Jon P. Boyle, Martin J. Boulanger, John Parkinson, and Michael E. Grigg. 2012.

550    "Integrated Bioinformatic and Targeted Deletion Analyses of the SRS Gene Superfamily

551    Identify SRS29C as a Negative Regulator of Toxoplasma Virulence." *mBio* 3 (6): e00321-

552    12. doi:10.1128/mBio.00321-12.

553    Zhang, Jianzhi. 2003. "Evolution by Gene Duplication: An Update." *Trends in Ecology and*

554    *Evolution*. doi:10.1016/S0169-5347(03)00033-8.

555    Zhou, S F, L P Yang, Z W Zhou, Y H Liu, and E Chan. 2009. "Insights into the Substrate

556    Specificity, Inhibitors, Regulation, and Polymorphisms and the Clinical Impact of Human

557    Cytochrome P450 1A2." *Aaps J* 11 (3): 481–94. doi:10.1208/s12248-009-9127-y [doi].

558

# Figure 1

MIPhy results interface.

MIPhy display for the 628 vertebrate Cyps from (Thomas, 2007). The MIGs are listed in the table on the left as well as indicated by the light orange shapes on the interior of the tree. The instability of each cluster is visualized by the bar charts around the outside of the tree. The colours of the band just inside of the circle match the colours of the tree nodes, and represent the originating species of each sequence.
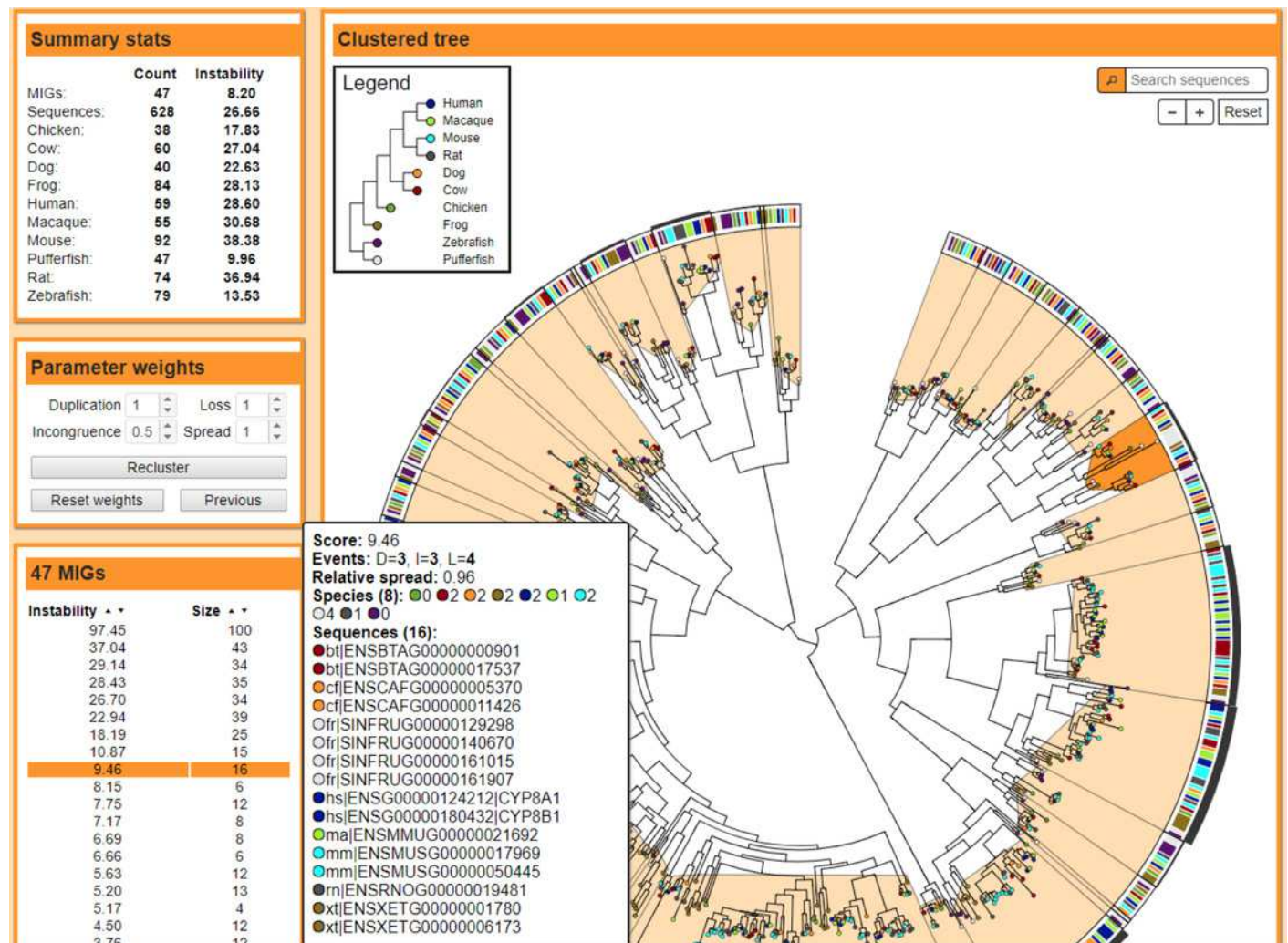
# Figure 2

The phylogenetic instability of the 59 human Cyp proteins.

The vertical dashed line separates the stable from the unstable sequences. 'Substrate' indicates those proteins with primarily endogenous roles (filled squares), primarily xenobiotic roles (empty circles), both xenobiotic and endogenous roles (empty squares), and pseudogenes (P). 'Selection' indicates which of the 18 sequences tested showed evidence of positive selection (+), or no positive selection (-). In the 'Clusters' row, the solid lines indicate those genes that are located in tandem arrays in the human genome or are syntenic with a tandem array in the mouse genome (S). All substrate, selection, and clustering data were taken from (Thomas, 2007).
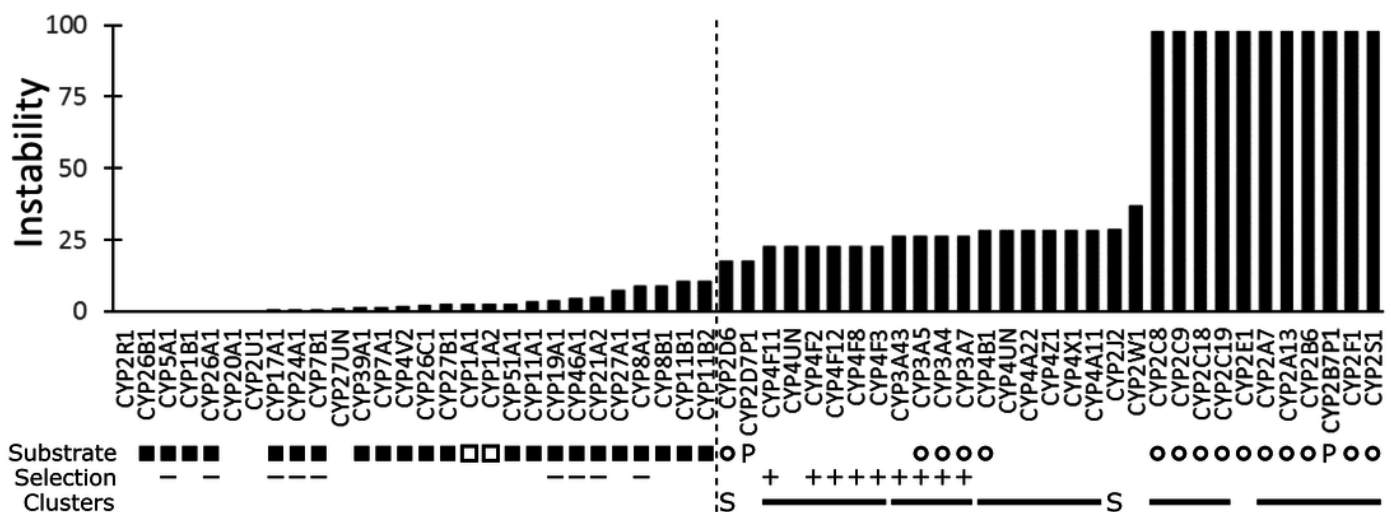
# Figure 3

The phylogenetic instability of the 151 *C. elegans* collagen MIGs.

The MIGs containing genes with an observable knock-down phenotype are indicated by triangles.