

The exhaustive genome comparison effort. A quarter-century later

Steven A. Benner, Kevin M. Bradley, Stephan G. Chamberlain

A quarter century after the Benner and Gonnet groups began their collaboration in evolutionary bioinformatics, evolution-based functional genomics is a field with considerable scope. Even with the remarkable advances in computing power over this period, the explosion of data derived from genomic and protein sources have required more and more sophisticated approaches be developed and utilized. We describe here new software combined with data organization techniques and illustrate how we are harnessing these to place physiological function of protein sequence data using natural history.

1 The exhaustive genome comparison effort. A quarter-century later

2
3 Steven A. Benner,* Kevin M. Bradley, Stephen G. Chamberlain

4 Foundation for Applied Molecular Evolution

5 Gainesville Florida

6 *Correspondence author: Steven A. Benner

7 720 SW 2nd Avenue, Suite 201

8 Gainesville, FL 32601

9 sbenner@ffame.org

10 352-219-3570

11
12 **Abstract**

13 A quarter century after the Benner and Gonnet groups began their collaboration in evolutionary
14 bioinformatics, evolution-based functional genomics is a field with considerable scope. Even with the
15 remarkable advances in computing power over this period, the explosion of data derived from genomic
16 and protein sources have required more and more sophisticated approaches be developed and utilized.
17 We describe here new software combined with data organization techniques and illustrate how we are
18 harnessing these to place physiological function of protein sequence data using natural history.

19
20 **Introduction**

21 It has now exactly a quarter-century since the Benner and Gonnet groups began their
22 collaboration in evolutionary bioinformatics [Gonnet and Benner 1991], a collaboration made
23 possible when my wife (Beverly Sanders) directed me to attend a seminar that Prof. Gonnet (then
24 from Waterloo, Canada) was giving on his work with the Oxford Unabridged English
25 Dictionary. That collaboration had at first only a modest goal: to update the (by then) more than
26 20 year old amino acid substitution matrix that had been introduced in the 1960s by Margaret
27 Dayhoff and her colleagues at the National Bureau of Standards [Dayhoff et al. 1972]. However,
28 even though the age of genomic sequence had not been begun, it was clear that it would soon get
29 underway, and that it would deliver a large number of whole genome sequences that could
30 service the platform for this new field. We wanted to be prepared for this revolution in biology.

31 Fortunately, the tools that the Gonnet group developed to organize the Oxford Unabridged
32 Dictionary were applicable to manage protein sequence databases. When applied to SwissProt

PeerJ PrePrints

33 [Bairoch and Apwiler 2000] and other early databases, this yielded the first exhaustive matching
34 of a protein genome sequence database and was published in *Science* in 1992 [Gonnet et al.
35 1992]. It provided not only new Dayhoff matrices, but also a clear understanding of how patterns
36 of amino acid substitution [Gonnet et al. 1994] and gapping [Benner et al. 1993] differed as two
37 protein sequences diverged. In the three dozen papers that were to follow, interpretation of those
38 patterns formed the basis for the then-emerging field of “evolutionary-based functional
39 genomics”, including the resurrection of ancestral genes and proteins [Jermann et al. 1995], the
40 use of evolutionary analyses to predict the folded structures of proteins [Benner et al. 1997], and
41 the analysis of the natural history of two families to understand adept, drift, functional change,
42 and pathway interactions [Benner et al. 1998] [Benner 2003]. Further, the exhaustive matching
43 supported of some of the earliest efforts to infer the metabolism of very ancient organisms
44 [Benner et al. 1989], including organisms standing at the branch points of the major three
45 kingdoms, and organisms that invented protein translation [Benner et al. 1993].

46 This work was well underway, of course, before complete genome sequences were available
47 for any individual organism. As these emerged for microorganisms, the Benner group, in
48 collaboration with EraGen Biosciences, introduced a naturally organized genome database
49 [Benner et al. 2000]. Called the MasterCatalog, the database organized protein sequences by
50 evolutionary families, much as had been done by more primitive databases dating back to
51 Dayhoff herself, but also in earlier versions of computerized database such as Hovergen [Duret
52 et al. 1994]. However, the MasterCatalog also included pre-computed trees, multiple sequence
53 alignments, and probabilistic ancestral sequences at nodes of the trees. A commercial version of
54 the MasterCatalog was bundled with several dozen complete genome sequences from various
55 microorganisms that had been assembled in a commercial effort at the company Genome
56 Therapeutics. This product generated approximately \$3.4 million in sales during its lifetime.

57 In a later version, secondary structure assignments determined by protein crystallography were
58 added to these evolutionary models for individual protein families to create the Magnum
59 database [Bradley and Benner 2006]. These supported a range of tools to extract functional
60 information from evolutionary comparisons between different species.

61 The Magnum database was announced just as whole genome sequences of vertebrates are
62 becoming available. This opened an entirely new direction for the assembly of an evolution-base,
63 naturally organized database, if the families of orthologous proteins from advanced organisms

64 could be reliably inferred with few errors. Complete genomic sequences were proposed to offer,
65 as one of their outputs, the prospect of knowing what genes and proteins are *not* present in a
66 biological organism. This prospect has driven, now for 20 years, the technology to determine
67 every last nucleotide in a chromosome, close all of the chromosomes in a genome, and provide a
68 complete list of genetic components in a complex organism.

69 The advance of deep sequencing, of course, created a crisis in data management. These crises
70 were associated with a series of problems briefly outlined below.

71

72 **Data volume**

73 Even in 1998, when MasterCatalog was conceived, sequence data resources were large and
74 their size was growing almost exponentially. Early development with GenBank involved a
75 dataset with some redundancy and about 700,000 sequences, which grew to well over 2.5 million
76 sequences in the space of 3 years.

77 Today, the number of bacterial whole genomes in RefSeq (a service provided by the NCBI) is
78 in the thousands, with is combined with the dozens of vertebrate organisms that have been
79 sequenced. It would be easy to collect 20 million sequences from whole genome sources alone.
80 This is such a rapid growth in computational demands that even technological increases in
81 computer power have been unable to keep up. In particular, a naïve all-against-all comparison
82 scales with the square of the database size. While indexing and other algorithmic tools can be
83 used to cause the scaling factor to be smaller, even Moore's law would be unable to manage the
84 size of the database. Therefore, any practical computational approach to organizing this data by
85 clustering needs to balance the time and space demands.

86 Fortunately, this problem could be mitigated simply by exploiting the realities of natural
87 history. The entire protein sequence space has not been explored during that natural history, not
88 the least of which and certainly not by vertebrates. Accordingly, focusing on chordate, once
89 genomes representing each of the major branches in the chordate tree are available, it is no
90 longer necessary to do an “all-against-all” comparison of an entire database.

91 MasterCatalog was designed to manage computational challenge of the rapidly growing
92 volume of data as efficiently as possible. Minimizing the number of sequence pairs for which
93 redundancy must be computed is the first step; performed using a BLAST comparison tuned for
94 nearly identical sequences, with pairwise comparisons only for sequences of the same species.

95 One sequence (the longest one) is used as a representative for subsequence comparisons. A
96 following BLAST “all-against-all” pairwise comparison considers only non-redundant
97 sequences, storing all significant matches. If more sophisticated estimates of sequence similarity
98 are required (true for some clustering algorithms), such as optimal local or semi-global
99 alignment, these are performed last and only for those sequence pairs that are significant non-
100 redundant matches.

101 With this, individual families within complete chordate genomes that are deemed it to be
102 especially reliable (in this case, we initially used 18) can be exhaustively matched, the families
103 identified, and evolutionary models (multiple sequence alignments, trees, and inferred
104 probabilistic ancestral sequences) constructed for use family. Further, an ancestral sequence
105 standing at the top of each nuclear family can be inferred for each of these families. Then, as new
106 genome sequences become available, or even as individual sequences become available separate
107 from whole genome sequencing efforts, a new exhaustive matching is not required.

108 Rather, the family to which the new sequence(s) belong can be identified by a search against
109 the founder sequences for each of the previously identified families. Then, the new sequence(s)
110 can be “tucked into” the multiple sequence alignment for the pre-computed family, and a branch
111 within that pre-computed and rectified family can be added to indicate the point of divergence of
112 the new sequence(s). While the pre-computed sequence would presumably have a fixed (and
113 evolutionarily correct) species tree topology, the new sequence(s) might cause minor adjustments
114 of the pre-computed evolutionary model for the. Of course, as the number of members within
115 each family increases, the impact of each addition becomes smaller. Over the long-term, one can
116 expect those models to become more or less stationary, with little change occurring further as
117 additional genomic sequences are added. This strategy, therefore, brings to an end the
118 computational challenge.

119 We report here our most recent efforts constructing evolutionarily organized databases
120 following the strategy. We also report an outline of the use of the database to characterize the
121 publicly available repertoire of whole genomic sequences.

122

123 **Data quality**

124 Even with confrontational challenge in hand, further problems emerge. Of particular
125 importance is the quality of the data contained in existing publicly available genome sequence

126 databases. Data quality is a problem of broad scope caused by problems at several levels in the
127 data collection process.

128 First, the sequence data delivered to genomic database can be unrepresentative of the actual
129 sequence of the providing organism. Sometimes, this is the result of low coverage shotgun
130 sequencing. Parts of the assembly may have coverage only from a single clone. Even worse,
131 circumstances exist when the shotgun sequences have been assembled against a template from
132 another source – often a different organism. Where short gaps exist arising from no coverage,
133 the template sequence is used instead of “N”, yielding a complete false impression of the quality
134 of the data. While such problems can be easily identified when examining the primary assembly,
135 it is impractical to do this on a genomic scale, and software working with genome-wide
136 comparisons accepts the bulk annotation without any effort to assess its validity.

137 A second problem with data quality is simply that the gene calling is inaccurate. Occasionally
138 this caused by gaps in the genome assembly, where whole exons (and introns) are missing from
139 the primary sequence data. More commonly, the gene calling shows variations in the start or end
140 point of genes that appear reasonably justified when the genome is examined in isolation, but are
141 obviously wrong when orthologous genes are examined together. This is usually limited to errors
142 in the end of one exon and/or the start of another, but can occasionally be much more
143 pronounced, with genes that would be expected to be orthologs (as judged by an alignment made
144 by MUMmer) being reported with wildly different transcripts that share few common exons.

145 As a third fact that creates problems when high quality alignments are required, alternative
146 transcripts are typically not consistently reported from one whole genome project to another.
147 Sometimes this may be because of differences in reporting criteria by the authors, but the quality
148 of EST data upon which gene predictions are based must also be a factor. Although we must
149 consider that EST data to be more reliable than raw DNA sequence data (it is, after all,
150 experimental evidence for the expression of particular transcripts), it is rare to find transcripts in
151 one organism that are supported by EST data to be disallowed due to mutation of alternative
152 splice sites in closely related orthologs (say amongst mammals); one infers that absence of an
153 alternative splice variant prediction is very weak evidence that the transcript is *not* present in the
154 second organism.

155 A less important problem is that functional annotations (the linguistic statement reporting to
156 describe the contribution of the protein to the fitness of the host organism) of genes in current

157 genomic databases can be quite wrong. While this makes functional annotation provided as
158 linguistic statements of limited use, it also clearly illustrates how much the process of annotation
159 (gene and functional) draws information from an arbitrary, but previously annotated homolog to
160 make predictions. It is sometimes possible to see which species has been used as a template for
161 another in the annotation process, simply by looking at the functional annotations and the gene
162 calls.

163

164 **Sequence clustering starting with 18 selected chordate families**

165 Clustering can be an expensive process, not merely because of the number of sequences that
166 are typically being organized, but because some clustering methods seek to cluster sequence
167 regions instead of whole sequences. The primary challenge in building clusters is how to use
168 (and weight) the similarities between sequences to generate desirable clusters. Many schemes
169 have been explored over the last 20 years. Two methods have been developed internally during
170 the development of MasterCatalog, while the GUI framework is capable of allowing the
171 examination of arbitrary clustering (called Catalogs).

172 We illustrate here one early run that examined the Ensembl65 sequences for 18 chordate
173 species (**Table 1**). The resulting catalog contained 8,199 individual families (with only clusters
174 with four or more sequences being considered a family); on average, each family had 44
175 members.

176

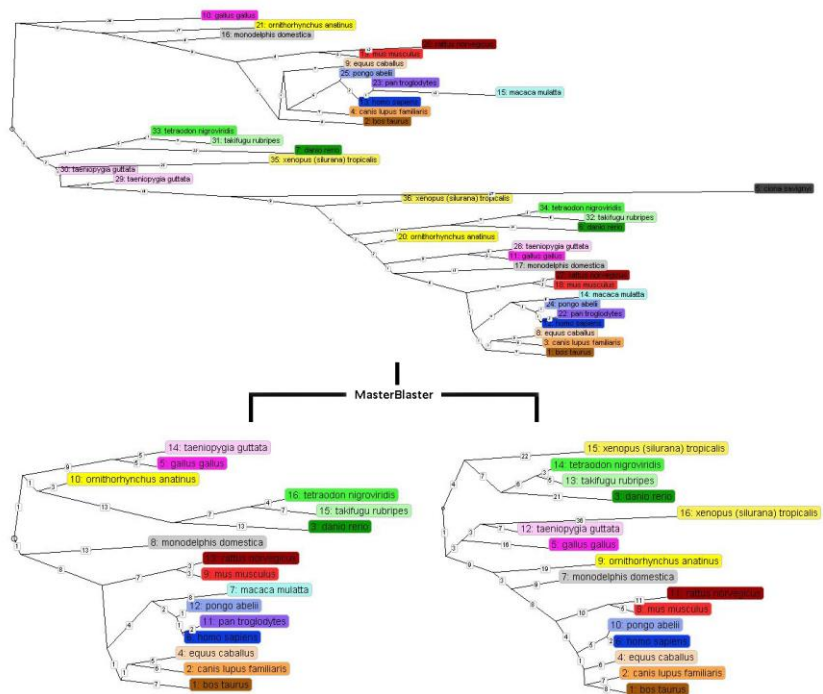
177 **Table 1.** The 18 “Whole” Chordate Genomes Examined Here

178 Species	Common Name
179 <i>Bos taurus</i>	Cow
180 <i>Canis lupus familiaris</i>	Dog
181 <i>Ciona savignyi</i>	Ciona
182 <i>Danio rerio</i>	Zebrafish
183 <i>Equus caballus</i>	Horse
184 <i>Gallus gallus</i>	Chicken
185 <i>Homo sapiens</i>	Human
186 <i>Macaca mulatta</i>	Rhesus monkey
187 <i>Monodelphis domestica</i>	Opossum
188 <i>Mus musculus</i>	Mouse
189 <i>Ornithorhynchus anatinus</i>	Platypus
190 <i>Pan troglodytes</i>	Chimpanzee
191 <i>Pongo abelii</i>	Orangutan
192 <i>Rattus norvegicus</i>	Rat
193 <i>Taeniopygia guttata</i>	Finch

194 *Takifugu rubripes* Fugu
 195 *Tetraodon nigroviridis* Pufferfish
 196 *Xenopus (Silurana) tropicalis* Frog

197
 198 Upon use of this catalog for individual family analyses, two problems became quickly
 199 apparent. First, protein sequences that had sufficient similarity to cluster in one family often had
 200 very different lengths. This caused problems in the creation of the MSA, which needed multiple
 201 gaps, often substantial in length, to accommodate sequences of very different lengths. This, in
 202 turn, corrupted distance metrics, which never score gaps correctly.

203 Second, when the clustering threshold was too low, very large families were created. In this
 204 case, 183 families contained more than 200 sequences. To address this problem, we created a
 205 tool to recluster families, internally referred to as MasterBlaster. This filter works by first
 206 examining the average sequence length with a family and removing any sequences that are of
 207 significantly different length from the average. Then, families were reclustered using BlastClust
 208 with more stringent criteria to allow larger families to be broken into smaller, more “natural”
 209 units, as well as further removing sequences that are distant enough to cause issues with the
 210 MSA. Application of MasterBlaster to the 18-genome database resulted in the creation of 14,058
 211 individual families with an average of 20 sequences per family and only 6 families containing
 212 more than 200 sequences.



213

214 **Figure 1.** Application of MasterBlaster to Family 1923 (HGNC symbol *LETMI*) derived from
215 the 18-genome comparison. Note that the root is placed internal to the tree. The MasterBlaster
216 splits a cumbersome family into two separate trees, as well as removing troublesome sequences
217 such as *Ciona savignyi* (colored above as black). While the *Ciona* sequence is conveniently used
218 to root the tree as the most primitive chordate outgroup, the time since its divergence from
219 vertebrates (~1 billion years in both branches) causes sufficient sequence divergence and
220 gapping to render distance metrics imprecise.

221

222 **Missing Data Analysis**

223 For an overwhelming majority of the families, one or more of the 18 "complete" genomes
224 examined did not have a representative within the tree. In fact, 13,146 of 14,058 families were
225 missing a sequence for at least one of the 18 species (93.5%). For example, **Figure 2** shows
226 Family 5041 (cysteine dioxygenase type1, *CDOI*). This family had an apparent ortholog in 17 of
227 the 18 species; however, no ortholog was found in the zebrafish. This could, of course, mean that
228 the gene was lost during the episode of natural history following the divergence of zebrafish
229 from other fish species. Alternatively, it could mean that (i) the whole genome sequence was less
230 "whole" than desired, with the DNA segment encoding that gene missed in the sequencing effort,
231 or (ii) the DNA segment encoding the *CDOI* gene was actually sequenced and present in the
232 database, but the bioinformatic gene finding tool failed to find it.

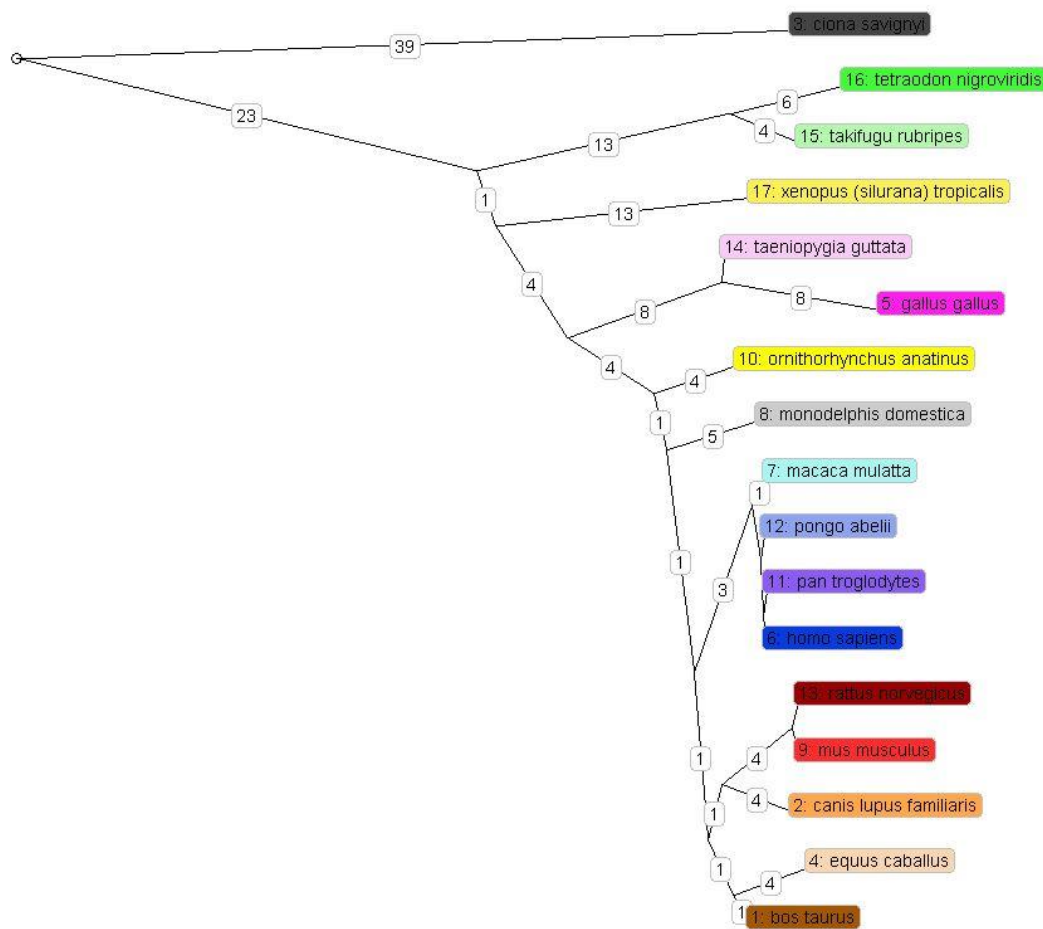
233 To assess the relative likelihood of these alternatives, we constructed a tool that would search
234 various sequence databases, including genomic, mRNA, and protein databases, in the event that
235 an ortholog was missing or truncated. Both mRNA and protein sequences from species present
236 in the tree were queried in hopes of finding missing sequences.

237 When this search tool was applied, we discovered a few alternative transcripts/translations that
238 were missed. For example, for Family 5041, we found that Ensembl protein
239 ENSDARP00000124052 from zebrafish would have clustered within this family and provided
240 the missing ortholog. However, an alternate sequence from zebrafish, ENSDARP00000085212,
241 was marked as the canonical protein sequence by the MasterCatalog; its sequence was too
242 divergent to have made it into the family.

243 Only a minority of the families could be completed using this strategy. Specifically, we
244 searched for 5292 sequences that were missing or truncated (defined as being < 50% of the

245 average length of sequences within the family) within 2843 families that had up to 3 of the 18
 246 species without representative sequences. In all, we found alternate transcripts accounted for
 247 17.1% of the missing sequences. We could further find 9.6% of the missing sequences in the
 248 genomic databases, and 7.7% within Ensembl's protein sequence database. However, even when
 249 missing sequences were found in the genomic or protein databases, these most often contained
 250 enough truncation or internal deletions to explain why they were not annotated by automated
 251 processes. For the remaining 65.6% of the sequences, our tool failed to find any reasonable trace
 252 that could complete the families. Thus an alternate approach was used to complete the families.

253



254

255 **Figure 2.** Family 5041, representing the *CDOI* gene, contains sequences for 17 of the 18
 256 species within MasterCatalog. Zebrafish (*danio rerio*) was not present in the initial clustering of
 257 this family.

258

259 **Inclusion of Ensembl Compara Data**

260 In 2008, Ensembl began to adopt several of the MasterCatalog innovations within a public
261 database, the Ensembl Compara database [Vilella et al. 2009]. With their extensive computing
262 resources, the Ensembl75 release delivered clusters and phylogenetic trees for proteins from 66
263 “whole” genome species. Accordingly, we exploit this data to bypass the protein clustering, by
264 far the most compute-expensive tasks, to allow this data to be analyzed within MasterCatalog.
265 Our use of the Ensembl75 clustering data helped reduce issues associated with gene finding,
266 including problems arising from alternate transcript/translation, and the missing sequence
267 problem, both by having the capability to do more exhaustive analysis of sequences and by
268 providing enough redundancy of species within tree branches to make an occasional missing
269 sequence acceptable.

270 The core databases, containing DNA and protein sequences, for each of these 66 species were
271 downloaded from Ensembl and loaded into the MasterCatalog database using built-in import
272 functionality. The Compara database was then downloaded and mirrored as a local database.

273 Using custom scripts, two catalogs were created using the Ensembl data. The first used the
274 protein clustering provided by Ensembl, but allowed the MasterCatalog to produce the MSAs
275 and trees. The second used the clustering, MSAs, and reconciled trees produced by Ensembl.
276 These two catalogs allow families to be viewed both with and without trees reconciled with the
277 expected species tree, allowing a better overall view of the data and the error contained within.

278 For both catalogs, MasterCatalog is able to calculate K_a/K_s values for each ancestral node
279 within a tree (**Figure 3**), a MasterCatalog innovation that has not yet been copied by the Ensembl
280 Compara database. It was, of course, implemented in the TAED database [Liberles et al. 2001],
281 which has been maintained and updated by the Liberles group. So far, the Ensembl Compara
282 database simply has these values indicated by leaf-leaf comparisons between contemporary
283 sequence pairs.

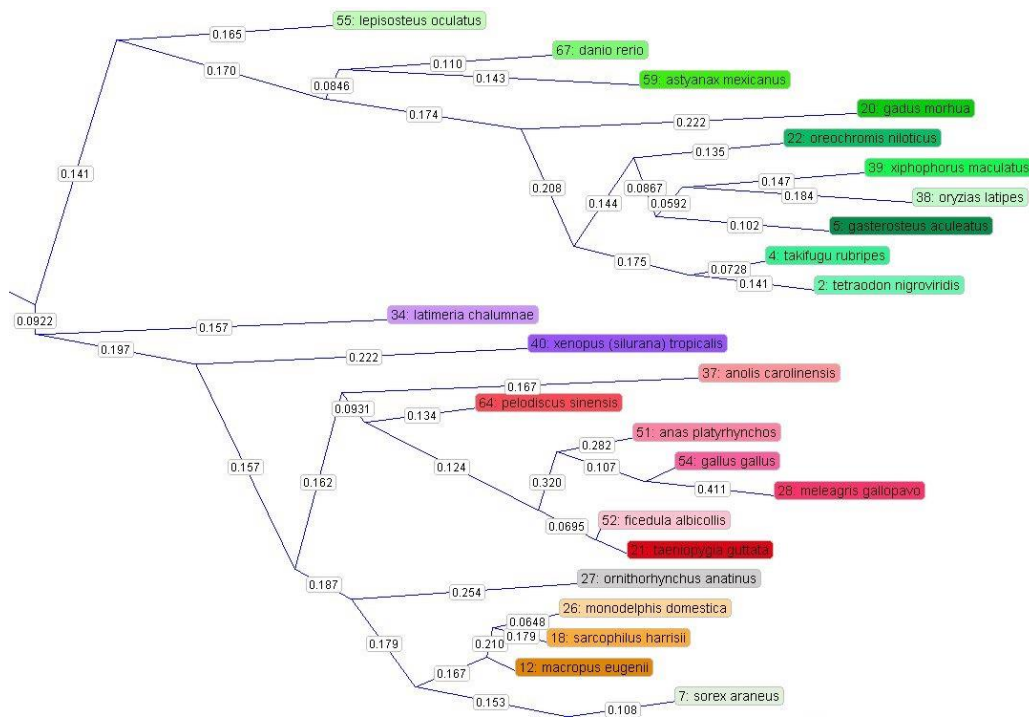


Figure 3: Computed K_a/K_s values are shown on ancestral branches for a portion of Family 4000 (HGNC symbol *FAM206A*)

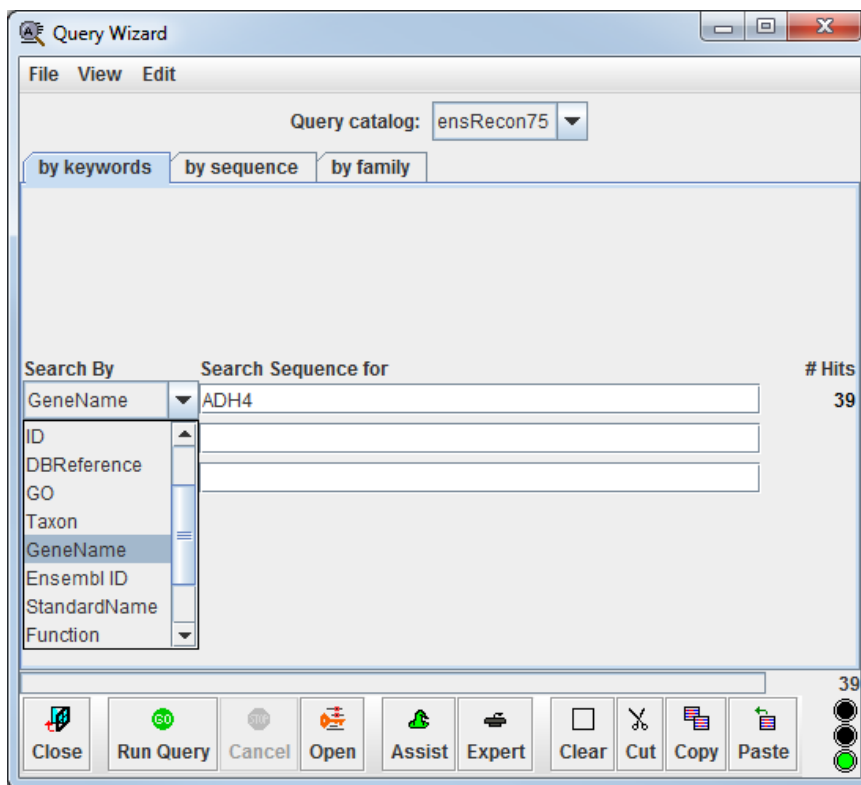
Using K_a/K_s values computed for ancestral nodes [Messier and Stewart 1997], we can apply many evolutionary functional genomic tools, including identifying groups of proteins that appear to have undergone active change at given times within evolutionary history. This is a powerful tool for finding groups of proteins that emerged or significantly changed function at the time that significant evolutionary changes occurred, such as the development of the breast and prostate (see below).

The Graphical User Interface

Evolutionary functional genomics is greatly facilitated by graphical user interfaces that allow scientists to "surf the genome", examining the pre-computed evolutionary models for individual protein families quickly. Accordingly, a major advantage of MasterCatalog is its high-level graphical user interface. This interface has been used to generate all the figures presented so far in this paper.

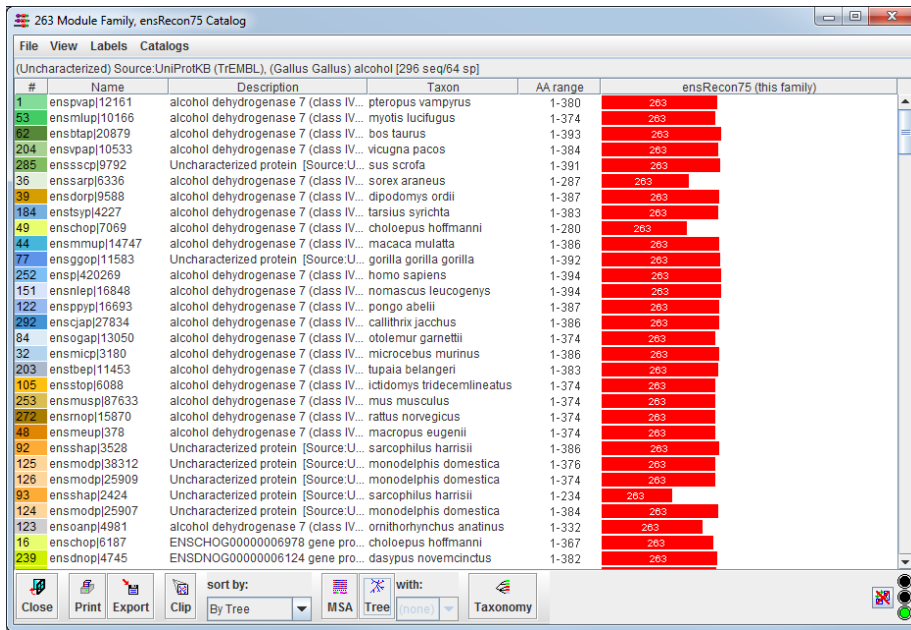
The interface is written in Java, which allows it to be compatible with all major operating systems. Although a full description of this interface would be outside the scope of this writing, we briefly describe some of the features here.

303 First, to view the underlying data, the scientist must be able to find the genes and protein
304 products of interest. To this end, MasterCatalog contains the ability to search by many fields,
305 including gene name, product name, gene description, and various external IDs, including
306 Ensembl's, OMIM, and GO (**Figure 4**). Families can also be found via sequences comparison,
307 with either protein or DNA/RNA sequence for the query. Once a family is found, it is displayed
308 to the user in tabular format with a graphical representation of the protein sequence (**Figure 5**).
309 From here, the user may choose to view these sequences as a multiple sequence alignment
310 (**Figure 6**) or the phylogenetic tree can be explored. The tree viewer has many features to help
311 the user explore the data. The entire tree can be fit to the window to gain a sense of the structure
312 of large families (**Figure 7**), while the scale can be quickly customized in both the X and Y
313 direction to focus on a specific region of the tree (**Figure 8**).
314

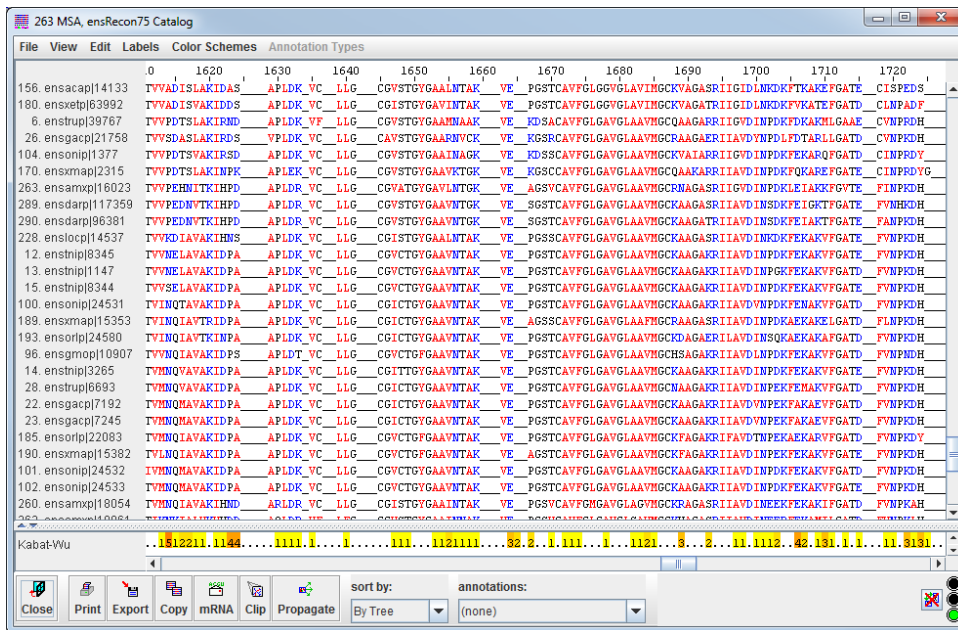


315
316 **Figure 4.** The query window for MasterCatalog allows the user to search by multiple fields to
317 assist in finding the correct sequence family.

318



319
320 **Figure 5.** A tabular view of the sequences in a MasterCatalog family with graphical
321 representation of sequences.



322
323 **Figure 6.** A view of MasterCatalog's MSA display window.
324

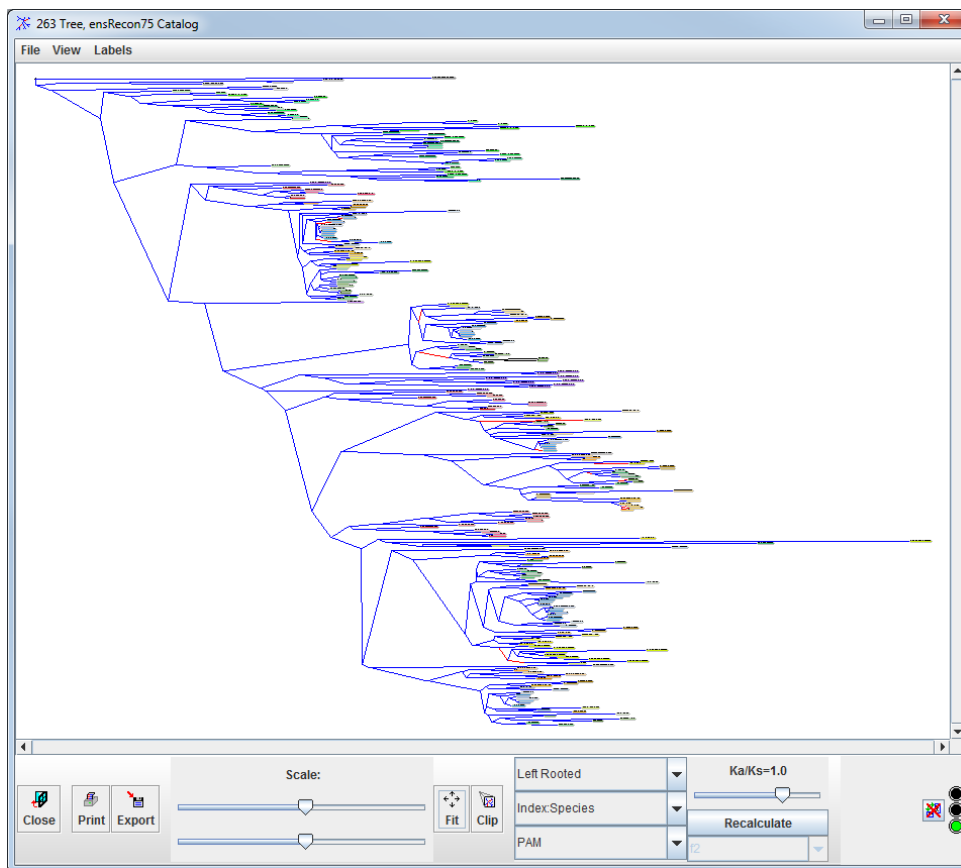
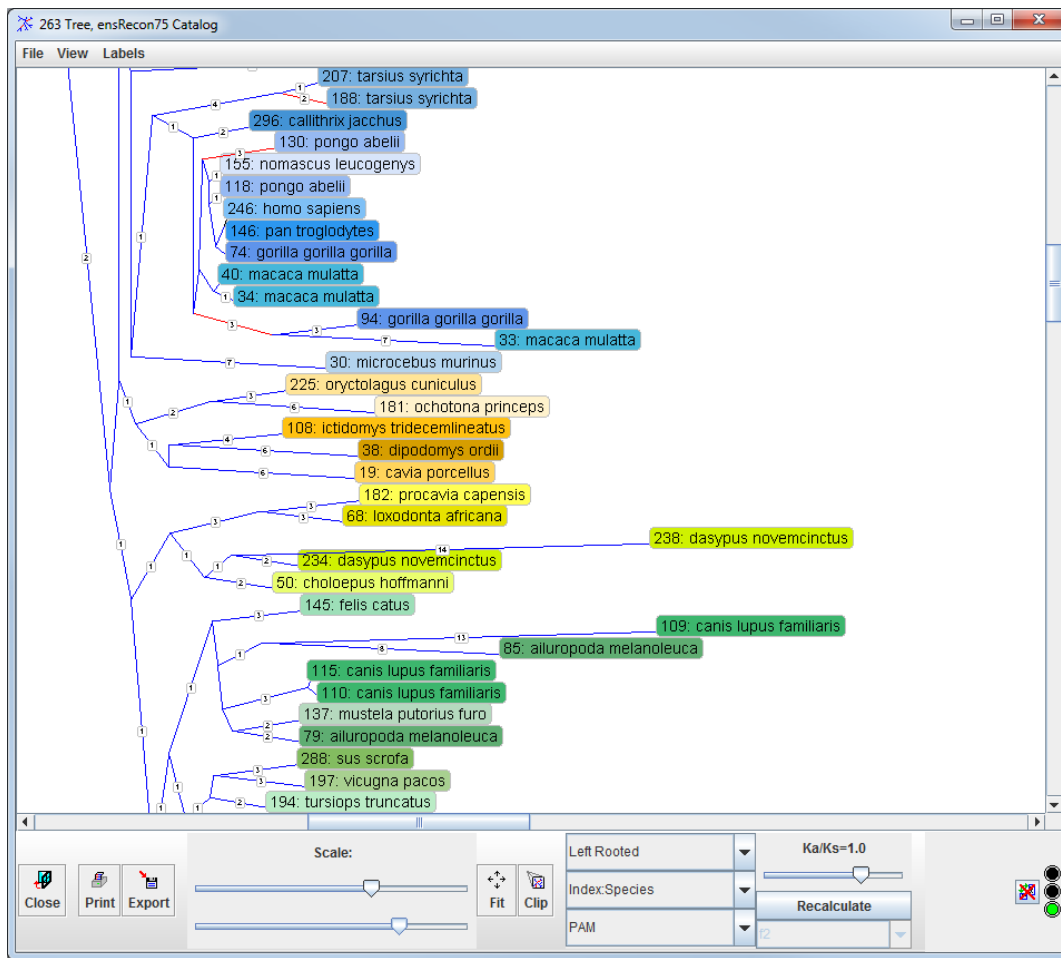
325
326
327

Figure 7. A top level view of a large phylogenetic tree in MasterCatalog (here displaying the sequence for alcohol dehydrogenase).



328
329 **Figure 8.** A zoomed in view of alcohol dehydrogenase sequences.

330
331 The user may also change tree structures (left rooted, top rooted, and unrooted), change the
332 description attached to leaves (species name, sequence name, sequence description, etc.), and
333 select the data displayed on the branches (PAM distance or K_a/K_s values). One advanced feature
334 within MasterCatalog is the ability to “clip” data, which allows for subsets of sequences within a
335 family to be displayed in either MSA or tree form. This clipped data can also have K_a/K_s values
336 recalculated, which can be helpful when removing error-filled sequences that are obviously
337 interfering with such calculations. Another advanced feature is the ability to view sequences in
338 multiple catalogs at the same time. This allows comparison between different clustering criteria.
339 Such comparisons can include sequence clustering vs. structural clustering, internal clustering vs.
340 Ensembl clustering, and comparisons of trees reconciled against the species trees with non-
341 reconciled trees. Such comparisons are unique to the MasterCatalog and allows the user a much
342 greater insight into the data. With this comes the ability to select groups of sequences in one

343 window and have those sequences select across all open datasets, helping the user quickly
344 navigate through a wealth of information.

345

346 **Examples of use**

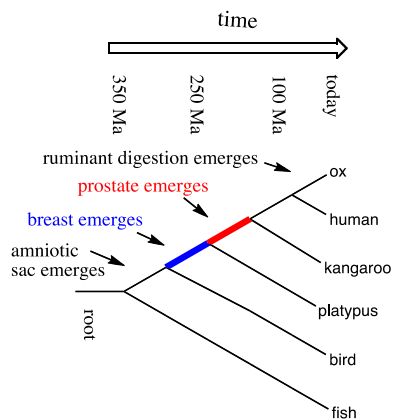
347 The database allows us to explore “high level” questions in biology using genome-supported
348 evolutionary analyses. For example, breast cancers display an intriguing mixture of
349 characteristics, each having an associated diagnostic/prognostic/therapeutic problem. For
350 example, improvements in screening (including ductal lavage) have allowed ductal carcinoma *in*
351 *situ* to be detected quite often. However, most ductal carcinoma do not become invasive, and it is
352 not understood why. Understanding of the “why” could minimize unnecessary interventions,
353 providing immediate improvements in the management of breast cancer.

354 Each of these characteristics presumably has one or more genetic/epigenetic correlates,
355 suggesting that if associated genes could be found, they might be sequenced in individual
356 patients to identify markers that would help diagnose the primary cancer, prognosis its course,
357 and choose a preferred therapy. Already this is done for estrogen-sensitive and insensitive
358 cancers. However, to date, only about half of breast cancers are explained by the broadest set of
359 risk factors; the most commonly used risk marker (BRCA1) covers perhaps only a fifth of total
360 cancer incidence. Likewise, early exposure to radiation changes the spectrum of cancer risk,
361 presumably by mutating genes, some perhaps not yet identified. Last, even if a breast cancer
362 patient has survived for five years, a good chance remains of recurrence, again with uncertain
363 etiology.

364 Thus, the ultimate overarching challenge (the elimination of mortality associated with breast
365 cancer) is associated with a challenge: Can we identify a spectrum of genes that, if sequence-
366 analyzed, can guide the genetic counselor, diagnostician, and physician in understanding the risk
367 of breast cancer in individual patients, identify consequential cancers at an early stage,
368 distinguish between early-stage aggressive from indolent cancers, and choose therapies that are
369 neither too much nor too little? Analyses of these genes in combination will, we hope, cover all
370 breast cancers, much as genetic analyses of BRCA or estrogen receptors offer similar guidance
371 for a fraction of those cancers.

372 Following an evolution-based functional genomics strategy, we begin by recognizing that the
373 breast (as a tissue) emerged only recently in the history of Earth, approximately 300 million

374 years ago (**Figure 9**). This episode is indicated in the tree in **Figure 9** by a blue line, the episode
 375 when suckling vertebrates emerged via divergence from other amniotes, most notably non-
 376 suckling birds and reptiles, which diverged still earlier from amphibians and, even earlier, from
 377 fish. The assignment of this time in natural history as the time when the breast emerged is, of
 378 course, identical to the statement from biological systematics that mammals form a true
 379 vertebrate class. The episode is recent enough in history to avoid much of the ambiguity that
 380 arises when bioinformatics tools model more ancient events.



381
 382 **Figure 9.** A schematic outlining the evolution of vertebrate tissues. Time is in million years.
 383 The red and blue lines indicate the episodes for the emergence of the prostate and breast
 384 (respectively).

385
 386 A second tissue new in mammals emerged in the episode immediately following (the red line):
 387 the prostate. This is indicated by the lack of a prostate in the platypus, but its presence in
 388 marsupials. The prostate is also a tissue that appears to generate cancer without obvious “insult”
 389 (although environmental factors can increase the incidence of prostate cancer, as they can breast
 390 cancer). Again, the red episode is sufficiently recent to avoid many ambiguities that make
 391 difficult bioinformatics analysis of more ancient events. Further, it is convenient to have two
 392 tissues from opposite genders equally susceptible to cancer and equally accessible to
 393 evolutionary analysis, as they can serve as controls (of a sort) for each other.

394 It is axiomatic in evolutionary developmental biology that the emergences of the breast and
 395 prostate in the blue and red episodes were associated with genetic changes. Further, Bayesian
 396 analyses are well known to be able to infer genetic events in the historical past through the
 397 analysis of modern gene sequences [Yang 1997]. Thus, we (and others) have long inferred the

398 sequences of ancestral genes and proteins from ancient genomes by analyzing the sequences of
399 their descendants. In a field invented in the Benner laboratory [Benner 2007], paleogenetics can
400 go still further, resurrecting inferred ancestral sequences from extinct animals by recombinant
401 DNA technology, making ancient proteins available for study in the laboratory. Since maximum
402 likelihood DNA and protein sequences at nodes in an evolutionary tree can be inferred using
403 Bayesian analysis, probabilistic changes can be assigned to individual branches in a tree like that
404 shown in **Figure 9**. Therefore, when applied to entire genomes, protein family by protein family,
405 we can say what amino acids were replaced, inserted, or deleted during the episode when the
406 breast emerged, or when the prostate emerged. While we agree that non-coding regions are also
407 important to an "evo-devo" analysis, these are not addressed here because of the greater
408 difficulty in inferring their histories. Further, numerous examples suggest that when a gene is
409 recruited to perform a new role, a signature of recruitment and its associated adaptive evolution
410 can be inferred by examining what amino acids are replaced, inserted, or deleted.

411 More sophisticated analyses can be applied across whole families and whole pathways. Thus,
412 we can suggest a central hypothesis: *To identify genes and proteins involved in the emergence of*
413 *the breast and/or prostate, we might go stepwise, family by family, through the genomic history*
414 *of vertebrates to find those that carry signatures of functional adaptation at the time when the*
415 *breast emerged and/or when the prostate emerged*. This work will deliver this family-by-family
416 analysis.

417 Various hypotheses give such analyses medical relevance. First, we hypothesize that genes
418 involved in the emergence of the breast (and, as a control, prostate) are likely candidates for
419 genes involved in regulating the growth, development, and functioning of these tissues in
420 modern mammals. Further, we hypothesize that mutations in these genes create susceptibility to
421 these cancers, determine the types of cancers that result, control the likelihood that those cancers
422 will progress and metastasize, and govern susceptibility of the resulting cancers to different
423 therapies. In this view, our evolutionary analysis will complement "classical" approaches to the
424 same goal, such as "deep sequencing" of multiple specimens of breast cancer tissue in search of
425 mutations with etiological significance, in established cancer-linked proteins, or the use of large-
426 scale typing of genetic markers in a case/control study format.

427 A naturally organized database can help, especially if provided with semi-automated tools that
428 will address error in genome annotation, heuristic development, and expert analysis [Benner et

449 prostate. In addition the work will integrate and exploit the wealth of available new data from
450 “modern biology” (e.g. functional genomics)

451
452 These are dated using the TREx clock [Li et al. 2006] to have occurred 31 ± 5 million years
453 ago (**Figure 10b**). One of the duplications is associated with a relatively high (0.93) K_a/K_s ratio
454 (**Figure 10a**). While this ratio is not greater than unity, it is large compared to ratios in other
455 branches of the tree (which are typically ~ 0.2) (**Figure 10a**). Therefore, this high *relative* ratio
456 suggests that this family is undergoing functional change. The very survival of paralogs, of
457 course, also suggests adaptation.

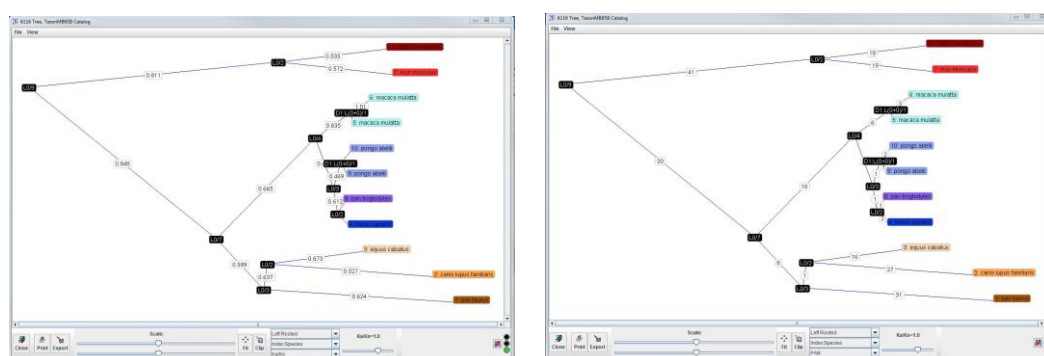
458 Had we stopped here, this case might have been just one many disputed examples in a
459 literature containing many of these. However, the next step in the multi-metric approach noted
460 that the amino acids replaced during the episode with a high K_a/K_s ratio were not randomly
461 distributed across the structure of the protein (**Figure 10d**). Rather, they were clustered near the
462 substrate binding site and the co-reduction binding site (**Figure 10d**). This suggests that during
463 this episode, details of the structure of the substrate had changed, as did the co-reductant. This
464 led to experimental work that showed that the different paralogs had different substrate
465 specificities (**Figure 10e**) [Corbin et al. 1999].

466 But what does it mean functionally? A cladogram based on fossil records (**Figure 10c**)
467 suggested that this episode of adaptive change occurs near the time in which pig litters went from
468 one piglet (with occasional twinning) to five or more piglets. This generated the hypothesis that
469 this gene triplication emerged to manage a new reproductive physiology in pigs (large litter size).
470 This drove an analysis of the molecular physiology, which confirmed this inference Corbin et al.
471 2004] [Kao et al. 2000] [Conley et al. 2001].

472 To date, multi-metric analyses have been generally developed case-by-case. From the Benner
473 group, these include analyses of dehydrogenases [Benner 1989], ribonucleases [Sassi and Benner
474 2007], leptins [Gaucher et al. 2003], sulfotransferases [Bradley and Benner 2005], inflammatory
475 proteins [Benner 2002], hypertension [Johnson et al. 2008], SARS [Benner et al. 2003], cystic
476 fibrosis [Gaucher et al. 2006], uterin serpins [Peltier et al. 2000], ribonucleotide reductase [Tauer
477 and Benner 1997], and elongation factors [Gaucher et al. 2001], among others [Benner et al.
478 2002]. In some cases, we have been interested in developing statistical heuristics that assess the
479 number of free variables that should be used to model adaptive divergence [Sassi et al. 2007]. In

480 other cases, we have explored the use of heterotachy (see below) to identify episodes of
 481 functional adaptation [Gaucher et al. 2002]. In other cases, we have asked how codon models
 482 [Benner 2012], scoring tools [Gonnet et al. 2000], and gapping models [Benner et al. 1993]
 483 [Chang and Benner 2004] improve multiple sequence alignments, or the impact of homoplasy in
 484 corrupting gene trees with short branch lengths [Carrigan et al. 2012]. In other cases, we have
 485 used the approach to improve the models upon which Bayesian inference relies [Gonnet et al.
 486 1994]. In others, the analysis has been the start of paleogenetics experiments, where ancestral
 487 proteins from now-extinct organisms are resurrected in the laboratory for study [Stackhouse et al.
 488 1990] [Jermann et al. 1995] [Ciglic et al. 1998] [Opitz et al. 1998] [Gaucher et al. 2003]
 489 [Thomson et al. 2005].

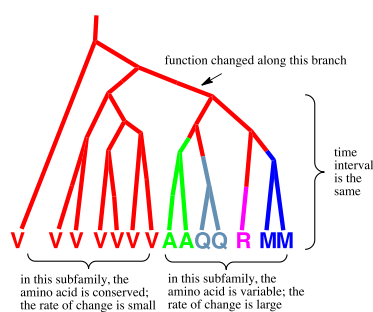
490 Again a principal challenge doing such work by automated methods arises from annotation and
 491 gene finding errors in whole genome sequence databases. These are illustrated in **Figure 11** for
 492 the primate BRCA1 gene family. Purely automated assembly of the family encounters situations
 493 where (in this case) it appears as if gene duplications created two paralogs in *Macaca* and *Pongo*
 494 (the rhesus monkey and orangutan, respectively). Of course, it is conceivable that this gene
 495 actually did suffer duplication independently in the two lineages leading to *Pongo* and *Macaca*.
 496 If so, this duplication would (like in the aromatase case) indicate functional adaptation,
 497 especially as the half-life for survival of nonfunctional duplicates that have not acquired a new
 498 function is about 11 million years [Trabesinger-Ruer et al. 1996]. In fact, this is not the case; the
 499 apparent paralogs in *Pongo* and *Macaca* are evidently the result of mistaken gene finding.



501
 502 **Figure 11.** The 18 vertebrate whole genomes within the MasterCatalog are illustrated here for
 503 the BRCA1 nuclear family for primates. This example illustrates “false paralogization”, where
 504 the *Pongo* and *Macaca* genomes have extra sequences that suggest duplicates where none exist.

505 These families will be rectified to remove such errors, when they corrupt the interpretation about
506 historical events at the time where the breast and prostate emerged.

507
508 We can further apply advanced metrics to detect functional change that complement the K_a/K_s
509 ratio and duplications [Benner 2004]. For example, heterotachy recognizes that two different
510 branches of a tree whose respective members have different functions also have different site-by-
511 site constraints on functional divergence (**Figure 12**). In lay language, that means that in proteins
512 having different functions, different sites evolve more rapidly while other sites evolve more
513 slowly. Likewise, homoplasy can indicate specific sites having specifically changing functional
514 roles. Other tools include an analysis of compensatory changes [Fukami-Kobayashi et al. 2002]
515 and crystallographic clustering, both of which bring crystallographic data to bear on an analysis
516 of functional divergence [Benner et al. 1997]. For example, amino acids being replaced during
517 an episode of relatively high non-synonymous substitution are often not distributed randomly
518 across the three-dimensional structure, but rather are clustered, perhaps near a substrate binding
519 or regulatory site. This crystallographic clustering is strong evidence for adaptive change and, as
520 in the aromatase example, can guide specific experiments to confirm/deny a hypothesis of
521 changing function.



522
523 **Figure 12.** “Heterotachy” is a change in the *rate* of amino acid substitution at a site that indicates
524 a change in function. It requires whole family analysis to detect. Shown the amino acids reside at
525 a site in a hypothetical protein. Purifying selection retained a valine at this position in the left
526 branch, but did not retain any amino acid in the right branch. The change in functional
527 constraints at this site indicates that the function in the protein changed in the episode indicated
528 by the arrow.

529
530 **Summary**

531 A quarter century has passed since Gaston Gonnet began to help us use evolutionary analysis
532 to extract function in ever-growing sequence databases. The results of this collaboration are now
533 having impact throughout biomedical research, much no longer acknowledged by (or even
534 known to) today's beneficiaries of a research program that began so long ago. However, the hour
535 spent 25 years ago in a seminar that Gaston gave on the Oxford Unabridged English Dictionary
536 was more than well spend. Thanks again to Beverly for making me aware of, and encouraging
537 me to attend, it.

538

539

- 540 Gonnet, G. H., Benner, S. A. (1991) Computational Biochemistry Research at ETH. *Technical*
541 *Report 154, Departement Informatik, March*
- 542 Dayhoff, M. O., Eck, R. V., Park, C. M. (1972) *A Atlas of Protein Sequence and Structure*. Vol.
543 5 (ed. M. O. Dayhoff)
- 544 Bairoch, A., Apweiler, R. (2000) The SWISS-PROT protein sequence database and its
545 supplement TrEMBL in 2000. *Nucl. Acids Res.* **28**, 45-48
- 546 Gonnet, G. H., Cohen, M. A., Benner, S. A. (1992) Exhaustive matching of the entire protein
547 sequence database. *Science* **256**, 1443-1445
- 548 Gonnet, G. H., Cohen, M. A., Benner, S. A. (1994) Analysis of amino acid substitution during
549 divergent evolution. The 400 by 400 dipeptide mutation matrix. *Biochem. Biophys. Res.*
550 *Comm.* **199**, 489-496
- 551 Benner, S. A., Cohen, M. A., Gonnet, G. H. (1993) Empirical and structural models for
552 insertions and deletions in the divergent evolution of proteins. *J. Mol. Biol.* **229**, 1065-
553 1082
- 554 Jermann, T. M., Opitz, J. G., Stackhouse, J., Benner, S. A. (1995) Reconstructing the
555 evolutionary history of the artiodactyl ribonuclease superfamily. *Nature* **374**, 57-59
- 556 Benner, S. A., Cannarozzi, G., Turcotte, M., Gerloff, D., Chelvanayagam, G., (1997) *Bona fide*
557 predictions of protein secondary structure using transparent analyses of multiple sequence
558 alignments. *Chem. Rev.* **97**, 2725-2843
- 559 Benner, S. A., Trabesinger-Ruef, N., Schreiber, D. R. (1998) Post-genomic science. Converting
560 primary structure into physiological function. *Adv. Enzyme Regul.* **38**, 155-180
- 561 Benner, S. A. (2003) Interpretive proteomics. Finding biological meaning in genome and
562 proteome databases. *Adv. Enzyme Regul.* **43** 271-359
- 563 Benner, S. A., Ellington, A. D., Tauer, A. (1989) Modern metabolism as a palimpsest of the
564 RNA world. *Proc. Nat. Acad. Sci.* **86**, 7054-7058
- 565 Benner, S. A., Cohen, M. A., Gonnet, G. H., Berkowitz, D. B., Johnsson, K. (1993) Reading the
566 palimpsest. Contemporary biochemical data and the RNA world. in *The RNA World*.
567 Gesteland. R. F., Atkins, J. F. editors, Cold Spring Harbor Laboratory Press, 27-70
- 568 Benner, S. A., Chamberlin, S. G., Liberles, D. A., Govindarajan, S., Knecht, L. (2000)
569 Functional inferences from reconstructed evolutionary biology involving rectified

570 databases. An evolutionarily-grounded approach to functional genomics. *Research*
571 *Microbiol.* **151**, 97-106.

572 Duret, L., Mouchiroud, D., Gouy, M. (1994) Hovergen, a database of homologous vertebrate
573 genes. *Nucleic Acids Res.* **22**, 2360-2365

574 Bradley, M. E., Benner, S. A. (2006) Integrating protein structures and precomputed genealogies
575 in the Magnum database. Examples with cellular retinoid binding protein. *BMC Evol.*
576 *Biol.* **7**, Art. No. 89

577 Vilella, A. J., Severink J., Ureta-Vidal, A., Heng, L., Durbin, R., Birney, E. (2009)
578 EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in
579 vertebrates. *Genome Res* **19**, 327-335

580 Liberles, D. A., Schreiber, D. R., Govindarajan, S., Chamberlin, S. G., Benner, S. A. (2001) The
581 adaptive evolution database (TAED) *Genome Biol.* **2**, 0028.1-0028.6

582 Messier, W., Stewart, C. B. (1997) Episodic adaptive evolution of primate lysozymes. *Nature*
583 **385**, 151-154

584 Yang, Z. (1997a) PAML: A program package for phylogenetic analysis by maximum likelihood.
585 *Comput. Appl. Biosci.* **13**, 555-556

586 Benner, S. A. (2007) The early days of paleogenetics. Connecting molecules to the planet.
587 Experimental Paleogenetics, D. A. Liberles, ed. Academic Press, pp. 3-19

588 Benner, S. A., Jermann, T. M., Opitz, J. G., Stackhouse, J., Knecht, L. J., Gonnet, G. H. (1995)
589 Uncertainty in ancient phylogenies. *Nature* **377**, 109-110

590 Liberles, D. A., Schreiber, D. R., Govindarajan, S., Chamberlin, S. G., Benner, S. A. (2001) The
591 adaptive evolution database (TAED). *Genome Biol.* **2**, 0003.1-0003.18

592 Gaucher, E. A., Graddy, L. G., Simmen, R. C. M., Simmen, F. A., Kowalski, A. A., Schreiber,
593 D. R., Liberles, D. A., Janis, C. M., Chamberlin, S. G., Benner, S. A. (2004) The
594 planetary biology of cytochrome P450 aromatases. *BMC Biology.* **2**, Art. No. 19

595 Li, T. *et al.* (2006) Analysis of transitions at two-fold redundant sites in mammalian genomes.
596 Transition redundant approach-to-equilibrium (TReX) distance metrics. *BMC*
597 *Evolutionary Biology*, **6**, 25.241

598 Corbin, C. J., Trant, J. M., Walters, K. W., Conley, A. J. (1999) Changes in testosterone
599 metabolism associated with the evolution of placental and gonadal isozymes of porcine
600 aromatase cytochrome P450. *Endocrinology* **140**, 5202-5210.

- 601 Corbin, C. J., Mapes, S. M., Marcos, J., Shackleton, C. H., Morrow, D., Safe, S., Wise, T., Ford,
602 J. J., Conley, A. J. (2004) Paralogues of porcine aromatase cytochrome p450: A novel
603 hydroxylase activity is associated with the survival of a duplicated gene. *Endocrinology*
604 **145**, 2157-2164.
- 605 Kao, Y. C., Higashiyama, T., Sun, X., Okubo, T., Yarborough, C., Choi, I., Osawa, Y., Simmen,
606 F. A., Chen, S. (2000) Catalytic differences between porcine blastocyst and placental
607 aromatase isozymes. *Eur, J, Biochem.* **267**, 6134-6139.
- 608 Conley, A., Mapes, S., Corbin, C. J., Greger, D., Walters, K., Trant, J., Graham, S. (2001) A
609 comparative approach to structure-function studies of mammalian aromatases. *J Steroid*
610 *Biochem.* **79**, 289-297.
- 611 Benner, S. A. (1989) Patterns of divergence in homologous proteins as indicators of tertiary and
612 quaternary structure. *Adv. Enz. Regul.* **28**, 219-236
- 613 Sassi, S. O., Benner, S. A. (2007) The resurrection of ribonucleases from mammals. From
614 ecology to medicine. *Experimental Paleogenetics*, D. A. Liberles, ed., NY, Academic
615 Press, pp 208-224
- 616 Gaucher, E. A., Miyamoto, M. M., Benner, S. A. (2003) Evolutionary, structural and
617 biochemical evidence for a new interaction site of the leptin obesity protein *Genetics* **163**,
618 1549-1553
- 619 Bradley, M. E., Benner, S. A. (2005) Phylogenomic approaches to common problems
620 encountered in the analysis of low copy repeats: The sulfotransferase 1A gene family
621 example. *BMC Evolutionary Biology* **5**, Art. No. 22
- 622 Benner, S. A. (2002) The past as the key to the present. Resurrection of ancient proteins from
623 eosinophils. *Proc. Natl. Acad. Sci. USA* **99**, 4760-4761
- 624 Johnson, R.J., Gaucher, E. A., Sautin, Y. Y., Henderson, G. N., Angerhofer, A. J., Benner, S. A.
625 (2008) The planetary biology of ascorbate and uric acid and their relationship with the
626 epidemic of obesity and cardiovascular disease. *Medical Hypotheses* **71**, 22-31
- 627 Benner, S. A., Gaucher, E. A., Li, T. (2003) Post-genomic evolutionary analyses of the Severe
628 Acute Respiratory Syndrome (SARS) virus genome using the MasterCatalog interpretive
629 proteomics platform. *Pharmagenomics –Application Notebook* **2003**, 23

- 630 Gaucher, E. A. De Kee, D. W., Benner, S. A. (2006) Application of DETECTER, an
631 evolutionary genomic tool to analyze genetic variation, to the cystic fibrosis gene family.
632 *BMC Genomics* **7**, No. 44
- 633 Peltier, M.R., Raley, L.C., Liberles, D. A., Benner, S.A., Hansen, P.J. (2000) Evolutionary
634 history of the uterine serpins. *J. Exp. Zool. (Mol. Devel. Evol.)* **288**, 165-174
- 635 Tauer, A., Benner, S. A. (1997) The B12-dependent ribonucleotide reductase from the
636 archaeobacterium *Thermoplasma acidophila*. An evolutionary conundrum. *Proc. Natl.*
637 *Acad. Sci.* **94**, 53-58
- 638 Gaucher, E. A., Miyamoto, M. M., Benner, S. A. (2001) Function-structure analysis of proteins
639 using covarion-based evolutionary approaches. Elongation factors. *Proc. Natl. Acad. Sci.*
640 *USA* **98**, 548-552
- 641 Benner, S. A., Caraco, M. D., Thomson, J. M., Gaucher, E. A. (2002) Planetary biology.
642 Paleontological, geological, and molecular histories of life. *Science* **293**, 864-868
- 643 Sassi, S. O., Braun, E. L., Benner, S. A. (2007) The evolution of seminal ribonuclease.
644 Pseudogene reactivation or multiple gene inactivation events? *Mol. Biol. Evol.* **24**, 1012-
645 1024
- 646 Gaucher, E. A., Gu, X., Miyamoto, M. M., Benner, S. A. (2002) Predicting functional divergence
647 in protein evolution by site-specific rate shifts. *Trends Biochem. Sci.* **27**, 315-321
- 648 Benner, S. A. (2012) Use of codon models in molecular dating and functional analysis. In *Codon*
649 *Evolution. Mechanisms and Models*. Oxford University Press, 133-144
- 650 Gonnet, G., Korostensky, C., Benner, S. (2000) Evaluation measures of multiple sequence
651 alignments. *J. Comput. Biol.* **7**, 261-276
- 652 Benner, S. A., Cohen, M. A., Gonnet, G. H. (1993) Empirical and structural models for
653 insertions and deletions in the divergent evolution of proteins. *J. Mol. Biol.* **229**, 1065-
654 1082
- 655 Chang, M., Benner, S. A. (2004) Empirical analysis of insertions and deletions in protein
656 sequence evolution. *J. Mol. Biol.* **341**, 617-631
- 657 Carrigan, M. A., Uryasev, O., Davis, R. P., Zhai, L-M., Hurley, T. D., Benner, S. A. (2012) The
658 natural history of Class I primate alcohol dehydrogenases includes gene duplication, gene
659 loss, and gene conversion. *PLOS 1* **7** (7) e41175

660 Gonnet, G. H., Cohen, M. A., Benner, S. A (1994) Analysis of amino acid substitution during
661 divergent evolution. The 400 by 400 dipeptide mutation matrix. *Biochem. Biophys. Res.*
662 *Comm.* **199**, 489-496

663 Stackhouse, J., Presnell, S. R., McGeehan, G. M., Nambiar, K. P., Benner, S. A. (1990) The
664 ribonuclease from an extinct bovid. *FEBS Lett.* **262**, 104-106

665 Jermann, T. M., Opitz, J. G., Stackhouse, J., Benner, S. A. (1995) Reconstructing the
666 evolutionary history of the artiodactyl ribonuclease superfamily. *Nature* **374**, 57-59

667 Ciglic, M. I., Jackson, P. J., Raillard-Yoon, S.-A., Haugg, M., Jermann, T. M., Opitz, J. G.,
668 Benner, S. A. (1998) Origin of dimeric structure in the ribonuclease superfamily.
669 *Biochemistry* **37**, 4008-4022

670 Opitz, J. G., Ciglic, M. I., Haugg, M., Trautwein-Fritz, K., Raillard-Yoon, S.-A., Jermann, T. M.,
671 Moore, R., Benner, S. A. (1998) Origin of the catalytic activity of bovine seminal
672 ribonuclease against double-stranded RNA. *Biochemistry* **37**, 4023-4033

673 Gaucher, E. A., Thomson, J. M., Burgan, M. F., Benner, S. A. (2003) Inferring the
674 paleoenvironment of ancient bacteria on the basis of resurrected proteins. *Nature* **425**,
675 285-288

676 Thomson, J. M., Gaucher, E. A., Burgan, M. F., De Kee, D.W., Li, T., Aris, J. P., Benner, S. A.
677 (2005) Resurrecting ancestral alcohol dehydrogenases from yeast. *Nature Genetics* **37**,
678 630-635

679 Trabesinger-Ruef, N., Jermann, T. M., Zankel, T. R., Durrant, B., Frank, G., Benner, S. A.
680 (1996) Pseudogenes in ribonuclease evolution. A source of new biomacromolecular
681 function? *FEBS Lett.* **382**, 319-322

682 Benner, S. A. (2004) Evolution-based functional proteomics. US Patent Application.
683 US20040204861 A1

684 Fukami-Kobayashi, K., Schreiber, D. R., Benner, S.A. (2002) Detecting compensatory
685 covariation signals in protein evolution using reconstructed ancestral sequences. *J. Mol.*
686 *Biol.* **319**, 729-743

687 Benner, S. A., Cannarozzi, G., Gerloff, D. L., Turcotte, M., Chelvanayagam, G. (1997) *Bona*
688 *fide* predictions of protein secondary structure using transparent analyses of multiple
689 sequence alignments. *Chem. Rev.* **97**, 2725-2843

690