

1 **Eukaryotic antiviral immune proteins arose via**
2 **convergence, horizontal transfer, and ancient inheritance**

3 Edward M. Culbertson^a and Tera C. Levin^{*a}

4

5 ^a University of Pittsburgh, Department of Biological Sciences

6 * Address correspondence to Tera C. Levin: teralevin@pitt.edu

7 Abstract

8 Animals use a variety of cell-autonomous innate immune proteins to detect viral
9 infections and prevent replication. Recent studies have discovered that a subset of mammalian
10 antiviral proteins have homology to anti-phage defense proteins in bacteria, implying that there
11 are aspects of innate immunity that are shared across the Tree of Life. While the majority of
12 these studies have focused on characterizing the diversity and biochemical functions of the
13 bacterial proteins, the evolutionary relationships between animal and bacterial proteins are less
14 clear. This ambiguity is partly due to the long evolutionary distances separating animal and
15 bacterial proteins, which obscures their relationships. Here, we tackle this problem for three
16 innate immune families (CD-NTases [including cGAS], STINGs, and Viperins) by deeply
17 sampling protein diversity across eukaryotes. We find that Viperins and OAS family CD-NTases
18 are truly ancient immune proteins, likely inherited since the last eukaryotic common ancestor
19 and possibly longer. In contrast, we find other immune proteins that arose via at least four
20 independent events of horizontal gene transfer (HGT) from bacteria. Two of these events
21 allowed algae to acquire new bacterial viperins, while two more HGT events gave rise to distinct
22 superfamilies of eukaryotic CD-NTases: the Mab21 superfamily (containing cGAS) which has
23 diversified via a series of animal-specific duplications, and a previously undefined eSMODS
24 superfamily, which more closely resembles bacterial CD-NTases. Finally, we found that cGAS
25 and STING proteins have substantially different histories, with STINGs arising via convergent
26 domain shuffling in bacteria and eukaryotes. Overall, our findings paint a picture of eukaryotic
27 innate immunity as highly dynamic, where eukaryotes build upon their ancient antiviral
28 repertoires through the reuse of protein domains and by repeatedly sampling a rich reservoir of
29 bacterial anti-phage genes.

30 Introduction

31 As the first line of defense against pathogens, all forms of life rely on cell-autonomous
32 innate immunity to recognize threats and respond with countermeasures. Until recently, many
33 components of innate immunity were thought to be lineage-specific[1]. However, new studies
34 have revealed that an ever-growing number of proteins used in mammalian antiviral immunity
35 are homologous to bacterial immune proteins used to fight off bacteriophage infections. This list
36 includes Argonaute, CARD domains, cGAS and other CD-NTases, Death-like domains,
37 Gasdermin, NACHT domains, STING, SamHD1, TRADD-N domains, TIR domains, and Viperin,
38 among others[2–13]. This discovery has been surprising and exciting, as it implies that some
39 cellular defenses have deep commonalities spanning across the entire Tree of Life. But despite
40 significant homology, these bacterial and animal immune proteins are often distinct in their
41 molecular functions and operate within dramatically different signaling pathways (reviewed
42 here[5]). *How, then, have animals and other eukaryotes acquired these immune proteins?*

43 One common hypothesis in the field is that these immune proteins are ancient, and have
44 been inherited since the last common ancestor of bacteria and eukaryotes[5]. In other cases,
45 horizontal gene transfer (HGT) between bacteria and eukaryotes has been invoked to explain
46 the similarities[6,14]. However, because most papers in this field have focused on searching

47 genomic databases for new bacterial immune genes and biochemically characterizing them, the
48 evolution of these proteins in eukaryotes has not been as thoroughly investigated.

49 To address this knowledge gap, we turned to the EukProt database, which has been
50 specifically developed to reflect the true scope of eukaryotic diversity through the genomes and
51 transcriptomes of nearly 1,000 species, specifically selected to span the eukaryotic tree [15].
52 EukProt contains sequences from NCBI and Ensemble, plus many diverged eukaryotic species
53 not found in any other database, making it a unique resource for eukaryotic diversity[15]. While
54 it can be challenging to acquire diverse eukaryotic sequences from traditional databases due to
55 an overrepresentation of metazoan data[16], EukProt ameliorates this bias by downsampling
56 traditionally overrepresented taxa.

57 Using this database, we investigated the ancestry of three gene families that are shared
58 between animal and bacterial immunity: Stimulator of Interferon Gamma (STING), cyclic GMP-
59 AMP synthase (cGAS) and its broader family of cGAS-DncV-like nucleotidyltransferases (CD-
60 NTases), and Viperin. STING, CD-NTases, and Viperin are all interferon-stimulated genes that
61 function as antiviral immune modules, disrupting the viral life cycle by activating downstream
62 immune genes, sensing viral infection, or disrupting viral processes, respectively[17]. We found
63 eukaryotic CD-NTases arose following multiple HGT events between bacteria and eukaryotes.
64 cGAS falls within a unique, mainly metazoan clade. In contrast, OAS-like proteins were
65 independently acquired and are the predominant type of CD-NTase found across most
66 eukaryotes. Separately, we have discovered diverged eukaryotic STING proteins that bridge the
67 evolutionary gap between metazoan and bacterial STINGs, as well as two separate instances
68 where bacteria and eukaryotes have acquired similar proteins via convergent domain shuffling.
69 Finally, we find that Viperin is likely to be truly ancient, with both broad representation across
70 the eukaryotic tree of life and evidence of two additional HGT events where eukaryotes recently
71 acquired new bacterial viperins. Overall, our results demonstrate that immune proteins shared
72 between bacteria and eukaryotes are evolutionarily dynamic, with eukaryotes taking multiple
73 routes to acquire and deploy these ancient immune modules.

74 Results

75 **Discovering immune homologs across the eukaryotic tree of life**

76 The first step to understanding the evolution of CD-NTases, STINGs, and viperins was
77 to acquire sequences for these proteins from across the eukaryotic tree. To search for diverse
78 immune homologs, we employed a hidden Markov model (HMM) strategy, which has high
79 sensitivity, a low number of false positives, and the ability to separately analyze multiple
80 (potentially independently evolving) domains in the same protein[18–20]. To broaden our
81 searches from initial animal homologs to eukaryotic sequences more generally, we used
82 iterative HMM searches of the EukProt database, incorporating the hits from each search into
83 the subsequent HMM. After using this approach to create pan-eukaryotic HMMs for each protein
84 family, we then added in bacterial homologs to generate universal HMMs (Fig. 1A and Supp.
85 Fig. 1), continuing our iterative searches until we either failed to find any new protein sequences
86 or began finding proteins outside of the family of interest (Supp. Fig. 1).

87 Our searches for CD-NTases, STINGs, and viperins recovered hundreds of eukaryotic
88 proteins from each family, including a particularly large number of metazoan sequences (red

89 bars, Fig. 1B). It is not surprising that we found so many metazoan homologs, as each of these
90 proteins was discovered and characterized in metazoans and these animal genomes tend to be
91 of higher quality than other taxa (Supp. Fig 2). We also recovered homologs from other species
92 spread across the Eukaryotic tree, demonstrating that our approach could successfully identify
93 deeply diverged homologs (Fig. 1B). However, outside of Metazoa, these homologs were
94 sparsely distributed, such that for most species in our dataset (711/993), we did not recover
95 proteins from the three immune families examined (white space, lack of colored bars, Fig. 1B).
96 We believe this pattern reflects a pattern of ongoing, repeated gene losses across eukaryotes,
97 as has been found for other innate immune proteins[21–23] and other types of gene families
98 surveyed across eukaryotes[22,24–26]. We found that BUSCO completeness scores and data
99 type (genomes vs. transcriptomes) were insufficient to explain the pattern of gene loss (Supp.
100 Fig. 2). Thus, although it is always possible that our approach has missed some homologs, we
101 believe the resulting data represents a fair assessment of the true diversity across eukaryotes.
102 To analyze these genes, we aligned the homologs with MAFFT and MUSCLE and then
103 generated phylogenetic trees with IQtree, FastTree, and RaxML-ng (see Materials and
104 Methods). We considered our results to be robust if they were concordant across the majority of
105 six trees generated per gene.

107 **Eukaryotes acquired CD-NTases from bacteria through at least three independent HGT** 108 **events**

109 We next studied the evolution of the innate immune proteins, beginning with cGAS and
110 its broader family of CD-NTase enzymes, which generate diverse oligonucleotides. In addition
111 to the well-studied cGAS, a number of other eukaryotic CD-NTases have been previously
112 described: 2'-5'-Oligoadenylate Synthetase 1/2/3 (OAS1/2/3), Male abnormal 21-Like 1/2/3/4
113 (MAB21L1/2/3/4), Mab-21 domain containing protein 2 (MB21D2), Mitochondrial dynamics
114 protein 49/51 (MID49/51), and Inositol 1,4,5 triphosphate receptor-interacting protein 1/2
115 (ITPRILP/1/2)[27]. Of these, cGAS and OAS1 are the best characterized and both play roles in
116 immune signaling. cGAS, Mab21L1, and MB21D2 are all cGAS-like receptors (cGLRs), and
117 recent work has shown that cGLRs are present in nearly all metazoan taxa and generate
118 diverse cyclic dinucleotide signals[28]. However, the immune functions of Mab21L1 and
119 MB21D2 remain unclear, although they have been shown to be important for development[29–
120 31].

121 Following infections or cellular damage, cGAS binds cytosolic DNA and generates cyclic
122 GMP-AMP (cGAMP)[32–35], which then activates downstream immune responses via STING
123 [34,36–38]. OAS1 synthesizes 2',5'-oligoadenylates which bind and activate Ribonuclease L
124 (RNase L)[39]. Activated RNase L is a potent endoribonuclease that degrades both host and
125 viral RNA species, reducing viral replication (reviewed here[40,41]). Some bacterial CD-NTases
126 such as *DncV* behave similar to animal cGAS; they are activated by phage infection and
127 produce cGAMP[8,42,43]. Other bacterial CD-NTases generate a wide variety of dinucleotides,
128 cyclic trinucleotides, and cyclic oligonucleotides[11]. These CD-NTases are commonly found
129 within cyclic oligonucleotide-based anti-phage signaling systems (CBASS) across many
130 bacterial phyla and even archaea[8,27,43].

131 To understand the evolutionary history of CD-NTases we used the Pfam domain
132 PF03281 as a eukaryotic starting point. As representative bacterial CD-NTases, we used 6,132

133 bacterial sequences, representing a wide swath of CD-NTase diversity[43]. Following our
134 iterative HMM searches, we recovered 313 sequences from 109 eukaryotes, of which 34 were
135 metazoans (Supplemental Data and Fig. 1B). Most eukaryotic sequences clustered into one of
136 two distinct superfamilies, which we name here for their highest scoring PFAM domain: Mab21
137 (Mab21: PF03281) or OAS (OAS1-C: PF10421) (Fig. 2A). Bacterial CD-NTases typically had
138 sequences matching the HMM for the Second Messenger Oligonucleotide or Dinucleotide
139 Synthetase domain (SMODS: PF18144).

140 The Mab21 superfamily is composed almost entirely of metazoan sequences, with only a
141 few homologs from Amoebozoa, choanoflagellates, and other eukaryotes (Fig. 2A). Indeed, the
142 majority of animal CD-NTases (cGAS, Mid51, Mab21, Mab21L1/2/3/4, Mb21d2, ITPRI) are
143 paralogs within the Mab21 superfamily, which arose from repeated animal-specific
144 duplications[44] (Supp. Fig 6). In contrast, unlike the animal-dominated Mab21 superfamily, the
145 OAS superfamily spans a broad group of eukaryotic taxa, with OAS-like homologs present in
146 8/12 eukaryotic supergroups. This distribution makes OAS proteins the most common CD-
147 NTases found across eukaryotes and implies that they arose very early in eukaryotic history,
148 possibly within the last eukaryotic common ancestor (LECA).

149 Given the connections between cGAS and STING in both animals and some
150 bacteria[3,43,45], we asked whether species that encode STING also have Mab21 and/or OAS
151 proteins. Because the Mab21 superfamily is largely animal-specific, we performed this analysis
152 separately in either Metazoa or with all non-metazoan eukaryotes (Fig. 2B). In animal species
153 where we found a STING homolog, we also typically found Mab21 superfamily sequence
154 (32/34), and a cGAS homolog in (26/34) species (Fig. 2B), consistent with the consensus that
155 these proteins are functionally linked. We also observed 19 metazoan species that had a
156 Mab21-like sequence with no detectable STING homolog. Almost half of these species (10/19)
157 were arthropods, agreeing with prior findings of STING sparseness among arthropods[45].
158 Outside of animals, we found that species with a STING homolog typically did not have a
159 detectable CD-NTase protein from either superfamily (22/34). While it remains possible that
160 these STING proteins function together with a to-be-discovered CD-NTase that was absent from
161 our dataset, we therefore hypothesize that many eukaryotes outside of metazoans and their
162 close relatives[46] use STING and CD-NTase homologs independently of each other.

163 What was the evolutionary origin of eukaryotic CD-NTases? Interestingly, the Mab21
164 and OAS superfamilies are only distantly related to one another. Each lies nested within a
165 different, previously defined, bacterial CD-NTase clade (Fig. 2 C and D). The OAS superfamily
166 falls within bacterial Clade C (with the closest related bacterial CD-NTases being those of
167 subclade C02-C03, Fig. 2C), while the metazoan Mab21 superfamily lies within bacterial Clade
168 D (subclade D12) (Fig. 2D). We note that in this tree (Fig. 2D), Clade D does not form a single
169 coherent clade, as was also true in the phylogeny that originally defined the bacterial CD-NTase
170 clades [11].

171 We also observed a number of eukaryotic sequences scattered across different bacterial
172 CD-NTase clades (Fig. 2A, colored branches within gray clades). While some of these may
173 reflect additional HGT events, others may come from technical artifacts such as bacterial
174 contamination of eukaryotic sequences. To minimize such false positive HGT calls, we took a
175 conservative approach in our analyses, considering potential bacteria-eukaryote HGT events to
176 be trustworthy only if: 1) eukaryotic and bacterial sequences branched near one another with

177 strong support (bootstrap values >70); 2) the eukaryotic sequences formed a distinct subclade,
178 represented by at least 2 species from the same eukaryotic supergroup; 3) the eukaryotic
179 sequences were produced by at least 2 different studies; and 4) the position of the horizontally
180 transferred sequences was robust across all alignment and phylogenetic reconstruction
181 methods used (Supp. Fig. 3). While these restrictions limit our attention to relatively old HGT
182 events, they also give us confidence these events are likely to be real.

183 The Mab21 superfamilies passed all four of these HGT thresholds, as did another
184 eukaryotic clade of CD-NTases that were all previously undescribed. We name this clade the
185 eukaryotic SMODS (eSMODS) superfamily, because the top scoring domain from hmmscan for
186 each sequence in this superfamily was the SMODS domain (PF18144), which is typically found
187 only in bacterial CD-NTases (Supplementary Data). This sequence similarity suggests that
188 eSMODS arose following a recent HGT from bacteria and/or that these CD-NTases have
189 diverged from their bacterial predecessors less than the eukaryotic OAS and Mab21 families
190 have. Additionally, all of these sequences were predicted to have a Nucleotidyltransferase
191 domain (PF01909), and (8/12) had a Polymerase Beta domain (PF18765), which are features
192 shared with many bacterial CD-NTases in Clades D,E, and F (Supplementary Data). The
193 eSMODS superfamily is made up of sequences from Amoebozoa, choanoflagellates,
194 Ancryomonadida, and one animal (the sponge *Oscarella pearse*), which clustered robustly and
195 with high support within bacterial Clade D (e.g. subclade D04, CD-NTase 22 from *Myxococcus*
196 *xanthus*) (Supp. Fig 4). The eSMODS placement on the tree, which was robust to all alignment
197 and phylogenetic algorithms used (Supp. Fig. 3), suggesting that eSMODS represent an
198 additional, independent acquisition of CD-NTases from bacteria.

199 CD-NTases from bacterial Clade C and Clade D are the only CD-NTases to produce
200 cyclic trinucleotides, producing cyclic tri-Adenylate and cAAG, respectively[11,47–49].
201 Interestingly, OAS produces linear adenylates, which is one step away from the cAAA product
202 made by Class C CD-NTases, and similarly cGAMP (made by cGAS) is one adenylate away
203 from the class D product cAAG. As of this writing, the Clade D CD-NTases closest to the
204 eSMODS and Mab21 superfamilies (D04 and D12, respectively), have not been well
205 characterized. Therefore we argue that these CD-NTases should be a focus of future studies,
206 as they may hint at the evolutionary stepping stones that allow eukaryotes to acquire bacterial
207 immune proteins.

208

209 **Diverged eukaryotic STINGs bridge the gap between bacteria and animals**

210 We next turned to analyze Stimulator of Interferon Gamma (STING) proteins. In animals,
211 STING is a critical cyclic dinucleotide sensor, important during viral, bacterial, and parasitic
212 infections (reviewed here[50]). Structurally, most metazoan STINGs consist of an N-terminal
213 transmembrane domain (TM), made of 4 alpha helices fused to a C-terminal STING domain[51].
214 Canonical animal STINGs show distant homology with STING effectors from the bacterial cyclic
215 oligonucleotide-based antiphage signaling system (CBASS), with major differences in protein
216 structure and pathway function between these animal and bacterial defenses. For example, in
217 bacteria, the majority of STING proteins are fusions of a STING domain to a TIR
218 (Toll/interleukin-1 receptor) domain (Fig. 3A). Bacterial STING proteins recognize cyclic di-GMP
219 and oligomerize upon activation, which promotes TIR enzymatic activity[3,52,53]. Some
220 bacteria, such as *Flavobacteriaceae*, encode proteins that fuse a STING domain to a

221 transmembrane domain, although it is unclear how these bacterial TM-STINGs function[3].
222 Other bacteria have STING domain fusions with deoxyribohydrolase, α/β - hydrolase, or trypsin
223 peptidase domains[14]. In addition to eukaryotic TM-STINGs, a few eukaryotes such as the
224 oyster *Crassostrea gigas* have TIR-STING fusion proteins, although the exact role of their TIR
225 domain remains unclear[3,54,55].

226 Given these major differences in domain architectures, ligands, and downstream
227 immune responses, how have animals and bacteria evolved their STING-based defenses, and
228 what are the relationships between them? Prior to this work, the phylogenetic relationship
229 between animal and bacterial STINGs has been difficult to characterize with high support[14].
230 Indeed, when we made a tree of previously known animal and bacterial STING domains, we
231 found that the metazoan sequences were separated from the bacterial sequences by one very
232 long branch, along which many changes had occurred (Fig. 3B).

233 To improve the phylogeny through the inclusion of a greater diversity of eukaryotic
234 STING sequences, we began by carefully identifying the region of STING that was homologous
235 between bacterial and animal STINGs, as we expected this region to be best conserved across
236 diverse eukaryotes. Although Pfam domain PF15009 (TMEM173) is commonly used to define
237 animal STING domains, this HMM includes a portion of STING's transmembrane domain which
238 is not shared by bacterial STINGs. Therefore, we compared the crystal structures of HsSTING
239 (6NT5), *Flavobacteriaceae* sp. STING (6WT4) and *Crassostrea gigas* STING (6WT7) to define
240 a core "STING" domain. We used the region corresponding to residues 145-353 of 6NT5 as an
241 initial HMM seed alignment of 15 STING sequences from PF1500915 ("Reviewed" sequences
242 on InterPro). Our searches yielded 146 eukaryotic sequences from 64 species, which included
243 STING homologs from 34 metazoans (Supplemental Data and Fig. 1). Using maximum
244 likelihood phylogenetic reconstruction, we identified STING-like sequences from 26 diverse
245 microeukaryotes that clustered in between bacterial and metazoan sequences, breaking up the
246 long branch. We name these sequences the bacteria-like STINGs (bSTINGs) because they
247 were the only eukaryotic group of STINGs with a bacteria-like Prok_STING domain (PF20300)
248 and due to the short branch length (0.86 vs. 1.8) separating them from bacterial STINGs on the
249 tree (Fig. 3C). While a previous study reported STING domains in two eukaryotic species (one
250 in Stramenopiles and one in Haptista) [14], we were able to expand this set to additional species
251 and also recover bSTINGs from Amoebozoa, Rhizaria and choanoflagellates. This diversity
252 allowed us to place the sequences on the tree with high confidence, recovering a substantially
253 different tree than previous work[14]. As for CD-NTases, the tree topology we recovered was
254 robust across multiple different alignment and phylogenetic tree construction algorithms (Supp.
255 Fig. 3).

256 Given the similarities between the STING domains of the bSTINGs and bacterial
257 STINGs, we next asked whether the domain architectures of these proteins were similar using
258 Hmmscan and AlphaFold. The majority of the new eukaryotic bSTINGs were predicted to have
259 four N-terminal alpha helices (Fig. 3A, and Supplementary Data), similar to human STING.
260 While bacterial TM-STINGs were superficially similar with N-terminal transmembrane domains,
261 these proteins were predicted to have only two alpha helices and in 5/6 phylogenetic trees
262 bacterial TM-STINGs were more similar to other bacterial STINGs than to eukaryotic homologs
263 (Supp. Fig. 3). These results suggest that eukaryotes and bacteria independently converged on
264 a common TM-STING domain architecture through domain shuffling.

265 Interestingly, a similar pattern of convergent domain shuffling appears to have occurred
266 a second time with the TIR-STING proteins. It was previously known that some eukaryotes such
267 as the oyster *C. gigas*, have a TIR-STING fusion protein[3,54,55]. The STING domain of these
268 TIR-STINGs clustered closely to other metazoan STINGs, suggesting an animal origin (Fig. 3B).
269 We also investigated the possibility that *C. gigas* acquired the TIR-domain of its TIR-STING
270 protein via HGT from bacteria, however this analysis also suggested an animal origin for the TIR
271 domain (Supp. Fig. 7). Eukaryotic TIR-STINGs are rare, further supporting the hypothesis that
272 this protein resulted from recent convergence, where animals independently fused STING and
273 TIR domains to make a protein resembling bacterial TIR-STINGs, consistent with previous
274 reports[14]. Overall, we find that the TM-STING and TIR-STING proteins represent at least two
275 independent examples of convergent evolution, where bacteria and eukaryotes have created
276 similar proteins through the reuse of ancient protein domains. Our work also identified a number
277 of non-metazoan STINGs (the bISTINGs) that have a domain architecture similar to animal
278 STINGs but a STING domain more similar to bacterial STINGs.

279

280 **Viperin is an ancient and widespread immune family**

281 Viperins are innate immune proteins that restrict the replication of a diverse array of
282 viruses by conversion of nucleotides into 3'-deoxy-3',4'didehydro- (ddh) nucleotides[4,56–58].
283 Incorporation of these ddh nucleotides into a nascent RNA molecule leads to chain termination,
284 blocking RNA synthesis and inhibiting viral replication[56,59]. While metazoan viperin
285 specifically catalyzes CTP to ddhCTP[56], homologs from archaea and bacteria can generate
286 ddhCTP, ddhGTP, and ddhUTP[4,60]. Previous structural and phylogenetic analysis showed
287 that eukaryotic viperins are highly conserved at both the sequence and structural level and that,
288 phylogenetically, animal and fungal viperins form a distinct monophyletic clade compared to
289 bacterial viperins[4,57,60].

290 As viperin proteins consist of a single Radical SAM protein domain, we iteratively
291 searched EukProt beginning with domain PF04055 (Radical_SAM). The 194 viperin-like
292 proteins we recovered came from 158 species spanning the full range of eukaryotic diversity,
293 including organisms from all of the major eukaryotic supergroups, as well as some orphan taxa
294 whose taxonomy remains open to debate (Fig. 1). When we constructed phylogenetic trees
295 from these sequences, we found that the large majority of the eukaryotic viperins cluster
296 together in a single, monophyletic clade, separate from bacterial or archaeal viperins (Fig. 4).
297 Within the eukaryotic viperin clade, sequences from more closely related eukaryotes often
298 clustered together (Fig. 4, colored blocks), as would be expected if viperins were present and
299 vertically inherited within eukaryotes for an extended period of time. The vast species diversity
300 and tree topology both strongly support the inference that viperins are a truly ancient immune
301 module and have been present within the eukaryotic lineage likely dating back to the last
302 eukaryotic common ancestor (LECA).

303 In addition to this deep eukaryotic ancestry, we also uncovered two examples of
304 bacteria-eukaryote HGT that have occurred much more recently, both in Chloroplastida, a group
305 within Archaeplastida. The first of these consists of a small clade of Archaeplastida (Clade A)
306 consisting of marine algae such as *Chlorocladus australicus* and *Nemeris dumetosa*. These
307 algal viperins cluster closely with the marine cyanobacteria *Anabaena cylindrica* and
308 *Plankthriodies* (Fig. 4 and Supp. Fig. 6). The second clade (Clade B) includes four other

309 Archaeplastida green algal species, mostly *Chlamydomonas spp.* In some of our trees the
310 Clade B viperins branched near to eukaryotic sequences from other eukaryotic supergroups,
311 however the placement of the neighboring eukaryotic sequences varied depending on the
312 algorithms we used; only the Archaeplastida placement was consistent. (Fig. 4 and Supp. Fig. 3
313 & 6). Taken together, we conclude that viperins represent a class of ancient immune proteins
314 that have likely been present in eukaryotes since the LECA. Yet, we also find ongoing
315 evolutionary innovation in viperins via HGT, both among eukaryotes and between eukaryotes
316 and bacteria.

317 Discussion

318 The recent discoveries that bacteria and mammals share mechanisms of innate
319 immunity have been surprising, because they imply that there are similarities in immunity that
320 span the Tree of Life. But how did these similarities come to exist? Here we uncover several
321 evolutionary trajectories that have led animals and bacteria to share homologous immune
322 proteins (summarized in Fig. 5). We found that Viperin is truly ancestral, dating back to at least
323 the Last Eukaryotic Common Ancestor (LECA), and likely further. We also uncovered examples
324 of convergence, as in STING, where the shuffling of ancient domains has led animals and
325 bacteria to independently arrive at similar protein architectures. Finally, we found evidence of
326 multiple examples of bacteria-eukaryote HGTs that have given rise to immune protein families.
327 An essential part of our ability to make these discoveries was the analysis of data from nearly
328 1000 diverse eukaryotic taxa. These organisms allowed us to distinguish between proteins
329 found across eukaryotes vs. animal-specific innovations, to document both recent and ancient
330 HGT events from bacteria that gave rise to eukaryotic immune protein families (Fig. 2 & 4), and
331 to identify STING proteins with eukaryotic domain architectures but more bacteria-like domains
332 (bSTINGs, Fig. 3). Because these diverged eukaryotic STINGs were found in organisms where
333 we typically did not find any CD-NTase proteins, we hypothesize that bSTINGs may detect and
334 respond to exogenous cyclic nucleotides, such as those generated by pathogens. In contrast to
335 the STINGs, the eukaryotic CD-NTases had substantially different evolutionary histories, with
336 multiple major CD-NTase superfamilies each emerging from within larger bacterial clades. While
337 these analyses cannot definitively determine the directionality of the transfer, we favor the most
338 parsimonious explanation that these components came into the eukaryotic lineage from
339 bacterial origins.

340 While not as prevalent as in bacteria, HGT in eukaryotes represents a significant force in
341 eukaryotic evolution, especially for unicellular eukaryotes[61–64]. In this study, our criteria for
342 ‘calling’ HGT events was relatively strict, meaning that our estimate of HGT events is almost
343 certainly an underestimate. Importantly, this pattern suggests that the bacterial pan-genome has
344 been a rich reservoir that eukaryotes have repeatedly sampled to acquire novel innate immune
345 components. Some of these HGT events have given rise to new eukaryotic superfamilies (e.g.
346 eSMODS) that have never been characterized and could represent novel types of eukaryotic
347 immune proteins. We speculate that the eSMODS superfamily CD-NTases and the bSTINGs
348 may function more similarly to their bacterial homologs, potentially producing and responding to
349 a variety of cyclic di- or tri-nucleotides[11] Similarly, bacterial viperins have been shown to
350 generate ddhCTP, ddhGTP, and ddhUTP, whereas animal viperins only make

351 ddhCTP[4,56,60]. Thus, the two algal viperin clades arising from HGT may have expanded
352 functional capabilities as well. A caveat of this work is that such strictly bioinformatic
353 investigations are insufficient to reveal protein biochemical functions, nor can they determine
354 whether diverse homologs have been co-opted for non-immune functions. We therefore urge
355 future, functional studies to focus on these proteins to resolve the questions of 1) whether/how
356 bLSTINGs operate in the absence of CD-NTases, 2) whether/how the functions of algal viperins
357 and eSMODS changed following their acquisition from bacteria, and 3) whether the homologs
358 truly function in immune defense.

359 In addition to these instances of gene gain, eukaryotic gene repertoires have been
360 dramatically shaped by losses. Even for viperins, which likely date back to the eukaryotic last
361 common ancestor, these proteins were sparsely distributed across eukaryotes and were absent
362 from the majority of species we surveyed. While some of this finding may be due to technical
363 limitations, such as dataset incompleteness or inability of the HMMs to recover distant
364 homologs, we believe this explanation is insufficient to fully explain the sparseness, as many
365 plant, fungal, and amoebozoan species are represented by well-assembled genomes where
366 these proteins are certifiably absent (Supp. Fig. 2). Instead, we propose that the sparse
367 distribution likely arises from ongoing and repeated gene loss, as has been previously
368 documented for other gene families across the Tree of Life[22,24–26].

369 Overall, our results yield a highly dynamic picture of immune protein evolution across
370 eukaryotes, wherein multiple mechanisms of gene gain are offset by ongoing losses.
371 Interestingly, this pattern mirrors the sparse distributions of many of these immune homologs
372 across bacteria[65–67], as anti-phage proteins tend to be rapidly gained and lost from genomic
373 defense islands[68,69]. It will be interesting to see if some eukaryotes evolve their immune
374 genes in similarly dynamic islands, particularly in unicellular eukaryotes that undergo more
375 frequent HGT[70].

376 We expect that our examination of STING, CD-NTases, and Viperin represents just the
377 tip of the iceberg when it comes to the evolution of eukaryotic innate immunity. New links
378 between bacterial and animal immunity continue to be discovered and other immune families
379 and domains such as Argonaute, Gasdermins, NACHT domains, CARD domains, TIR domains,
380 and SamHD1 have been shown to have bacterial roots[2,6,7,9,10]. To date, the majority of
381 studies have focused on proteins specifically shared between metazoans and bacteria. We
382 speculate that there are probably many other immune components shared between bacteria
383 and eukaryotes outside of animals. Further studies of immune defenses in microeukaryotes are
384 likely to uncover new mechanisms of cellular defense and to better illustrate the origins and
385 evolution of eukaryotic innate immunity.
386

387 Methods

388 Iterative HMM Search

389 The goal of this work was to search the breadth of EukProt v3 for immune proteins from
390 the CD-NTase, STING, and viperin families that span the gap between metazoan and bacterial
391 immunity. Our overall strategy was to first search with eukaryotes alone (starting from mainly
392 Metazoa). Then we added in bacterial sequences and searched with a mixed bacterial-

393 eukaryotic HMM search until we either found no new hits, or until we began getting hits from an
394 outgroup gene family. In parallel, we also performed bacteria-only and eukaryote-only searches,
395 to ensure that we found as many homologs as possible (schematized in Fig 1A, and further in
396 Supp. Fig. 1A).

397 Phase 1: Eukaryotic searches To begin, hidden Markov model (HMM) profiles from Pfam
398 (for CD-NTases and Viperin) or an HMM profile generated from a multiple sequence alignment
399 (for STING) were used to search EukProt V3[15], for diverse eukaryotic sequences. For CD-
400 NTases and Viperin, HMM profiles of Pfams PF03281 and PF04055 were used respectively.

401 For STING, where the Pfam profile includes regions of the protein outside of the STING
402 domain, we generated a new HMM for the initial search. First, we aligned crystal structures of
403 HsSTING (6NT5), *Flavobacteriaceae* sp. STING (6WT4) and *Crassostrea gigas* STING (6WT7)
404 to define a core “STING” domain. Then we aligned 15 eukaryotic sequences from PF15009
405 (“Reviewed” sequences on InterPro) with MAFFT(v7.4.71)[71] and manually trimmed the
406 sequences down to the boundaries defined by our crystal alignment (residues 145-353 of
407 6NT5). We then trimmed the alignment with TrimAl (v1.2)[72] with options -gt 0.2. The trimmed
408 MSA was then used to generate an HMM profile with hmmbuild from the hmmer (v3.2.1)
409 package (hmmer.org).

410 HMM profiles were used to search EukProt via hmmsearch (also from hmmer v3.2.1)
411 with a statistical cutoff value of 1e-3 and -hit parameter set to 10 (i.e. the contribution of a single
412 species to the output list is capped at 10 sequences). The resulting hits from this search were
413 then aligned with hmalign (included within hmmer) and used to generate a new HMM profile
414 with hmmbuild. This profile was used to search EukProt v3 again and the process was repeated
415 for a total of 3-4 eukaryotic searches.

416 Phase 2: combining eukaryotic and bacterial sequences into an HMM After the
417 eukaryotic searches reached saturation (i.e. no additional eukaryotic sequences were recovered
418 after additional searches), bacterial sequences were acquired from previous literature (Viperins
419 from[4], CD-NTases from[11], and STINGs from[3,8,43]). To ensure the combined HMM did not
420 have an overrepresentation of either bacterial or eukaryotic sequences, we downsampled the
421 bacterial sequences and eukaryotic sequences to obtain 50 phylogenetically diverse sequences
422 of each, and then combined the two downsampled lists. To do this, eukaryotic and bacterial
423 sequences were each separately aligned with MAFFT, phylogenetic trees were built with
424 FastTree (v2.1.10)[73], and the Phylogenetic Diversity Analyzer (pda/1.0.3)[74] software was
425 used to downsample the sequences while maximizing remaining sequence diversity.

426 The combined bacterial-eukaryotic sequence list was then aligned with hmalign and
427 used to construct a new HMM profile with hmmbuild. This HMM profile was used to search
428 EukProt v3. The eukaryotic hits from this search were then aligned with MAFFT, and a tree was
429 constructed with FastTree. From this tree the sequences were then downsampled with PDA and
430 once again combined with the bacterial list, aligned, used to generate a new HMM, and a new
431 search. This process was iterated 3-5 times until saturation or until the resulting sequence hits
432 included other gene families that branched outside of the sequence diversity defined by the
433 metazoan and bacterial homologs.

434 Phase 3: Searching with a bacteria-only or existing eukaryote-only HMM profiles

435 To search EukProt v3 with a bacteria-only HMM for each protein family, we aligned the
436 full set of published bacterial sequences with MAFFT, trimmed with TrimAl, and hmmbuild was

437 used to generate an HMM profile which was used to search EukProt v3. As a point of
438 comparison, we also searched the database with only the starting, animal-dominated Pfam
439 (PF15009) for STING.

440 Phase 4: Combining all hits into a single list and scanning for domains

441 Sequences from all iterative searches were combined to generate a total hits FASTA file
442 for STING, CD-NTase, and Viperin. First, duplicate sequences were removed, then the fasta
443 files were scanned using hmmscan (also from hmmer v3.2.1) against the Pfam database (Pfam-
444 A.hmm) and all predicted domains with an E-value <1e-3 were considered. Next, we generated
445 phylogenetic trees (first by aligning with MAFFT and then building a tree with FastTree), and
446 used these trees along with the hmmscan domains to determine in-group and out-group
447 sequences. Out-group sequences were removed from the fasta file. We determined outgroup
448 sequences by these criteria: 1) if the sequence clustered outside of known outgroup sequences
449 (e.g. Poly(A) RNA polymerase (PAP) sequences for the CD-NTases, and molybdenum cofactor
450 biosynthetic enzyme (MoaA) for Viperin), or 2) if sequence did not have at least one of the
451 relevant domains (Mab21/OAS1-C/SMODS for CD-NTases, TMEM173/Prok_STING for STING,
452 and Radical_SAM for Viperin). These three FASTA files were used for the final alignments and
453 phylogenetic trees. To identify protein domains in each sequence, the FASTA files were
454 scanned using hmmscan (also from hmmer v3.2.1) against the Pfam database (Pfam-A.hmm)
455 and all predicted domains with an E-value <1e-3 were considered. See Supplemental Data for
456 the hmmscan results of all included homologs.

457

458 **Final Alignment and Tree Building**

459 To generate final phylogenetic trees, all eukaryotic search hits and bacterial sequences
460 were aligned using MAFFT. We downsampled the CD-NTase bacterial sequences from ~6000
461 down to 500 as described above, to facilitate more manageable computation times on
462 alignments and tree construction. For the STING and Viperin trees we included all bacterial
463 sequences. These initial alignments were first trimmed manually in Geneious (v2023.1.2) to
464 remove unaligned N- and C-terminal regions, and then realigned with MAFFT or
465 MUSCLE(v5.1)[75] and trimmed with TrimAl (v1.2)[72]. MUSCLE was used with the “-super5”
466 option, and otherwise default parameters. These alignments were used to generate
467 phylogenetic trees using three tree inference softwares: FastTree (v2.1.10)[73], IQtree
468 (2.2.2.7)[76] and RaxML-ng (v0.9.0)[77]. FastTree was utilized with default settings. IQtree was
469 used to determine the appropriate evolutionary model, and was run with 1000 ultrafast
470 bootstraps (IQtree settings: -s, -bb 1000, -m TEST, -nt AUTO). RaxML-ng trees were produced
471 with 100 bootstraps using the molecular model specified from the IQtree analysis (Raxml-ng
472 settings: --all, --model [specified by IQtree], --tree pars{10} --bs-trees 100). Phylogenetic trees
473 were visualized with iTOL[78].

474

475 **TIR Domain Alignment and Tree**

476 We used hmmscan to identify the coordinates of TIR domains in a list of 203 TIR domain
477 containing-sequences from InterPro (Family: IPR015032) and 104 bacterial TIR-STING proteins
478 (the same TIR-STING proteins used in Fig. 3)[3]. Next, we trimmed the sequences down to the
479 TIR coordinates and aligned the TIR domains with MUSCLE. We trimmed the alignments with
480 TrimAL and built a phylogenetic tree with IQtree.

481

482 **Venn Diagrams**

483 Venn diagrams were generated via DeepVenn[79] using presence/absence information
484 for Mab21, OAS, and STING from each eukaryotic species that encoded at least one of these
485 proteins.

486

487 **Protein Structure Modeling**

488 To model 3-D protein structures for STING homologs without a published crystal
489 structure, we ran AlphaFold (v2.1.1)[80,81]. We generated 5 ranked models for STINGs from
490 *Flavobacteriaceae* (IMG ID: 2624319773), *Nitzschia sp.* (EukProt ID: P007051), and
491 *Caveostelium apophysatum* (EukProt ID: P019191). Fig. 2C shows highest ranked models only.

492

493 **Acknowledgements**

494 We thank Daniel Richter for his feedback, encouragement, and scientific guidance. Maureen
495 Stolzer, Kevin Forsberg, Patrick Mitchell, and members of the Levin lab also provided helpful
496 input on the project and manuscript. Thanks to Andrew VanDemark for helpful discussions
497 about 3-D modeling and to Jacob Durrant for help running AlphaFold. This research was
498 supported in part by the University of Pittsburgh Center for Research Computing,
499 RRID:SCR_022735, through the resources provided. Specifically, this work used the HTC
500 cluster, which is supported by NIH award number S10OD028483. Ed Culbertson was supported
501 by NSF Postdoctoral fellowship 2208971 and Tera Levin was supported by NIH R00AI139344
502 and R35GM150681.

503 **References:**

- 504 1. Litman GW, Cannon JP, Dishaw LJ. Reconstructing immune phylogeny: new
505 perspectives. *Nat Rev Immunol.* 2005;5: 866–879.
- 506 2. Johnson AG, Wein T, Mayer ML, Duncan-Lowey B, Yirmiya E,
507 Oppenheimer-Shaanan Y, et al. Bacterial gasdermins reveal an ancient
508 mechanism of cell death. *bioRxiv.* 2021. p. 2021.06.07.447441.
509 doi:10.1101/2021.06.07.447441
- 510 3. Morehouse BR, Govande AA, Millman A, Keszei AFA, Lowey B, Ofir G, et al.
511 STING cyclic dinucleotide sensing originated in bacteria. *Nature.* 2020;586:
512 429–433.
- 513 4. Bernheim A, Millman A, Ofir G, Meitav G, Avraham C, Shomar H, et al.
514 Prokaryotic viperins produce diverse antiviral molecules. *Nature.* 2021;589:
515 120–124.

- 516 5. Wein T, Sorek R. Bacterial origins of human cell-autonomous innate immune
517 mechanisms. *Nat Rev Immunol*. 2022. doi:10.1038/s41577-022-00705-4
- 518 6. Kibby EM, Conte AN, Burroughs AM, Nagy TA, Vargas JA, Whalen LA, et al.
519 Bacterial NLR-related proteins protect against phage. *Cell*. 2023.
520 doi:10.1016/j.cell.2023.04.015
- 521 7. Swarts DC, Makarova K, Wang Y, Nakanishi K, Ketting RF, Koonin EV, et al.
522 The evolutionary journey of Argonaute proteins. *Nat Struct Mol Biol*.
523 2014;21: 743–753.
- 524 8. Cohen D, Melamed S, Millman A, Shulman G, Oppenheimer-Shaanan Y,
525 Kacén A, et al. Cyclic GMP–AMP signalling protects bacteria against viral
526 infection. *Nature*. 2019;574: 691–695.
- 527 9. Ofir G, Herbst E, Baroz M, Cohen D, Millman A, Doron S, et al. Antiviral
528 activity of bacterial TIR domains via immune signalling molecules. *Nature*.
529 2021;600: 116–120.
- 530 10. Tal N, Millman A, Stokar-Avihail A, Fedorenko T, Leavitt A, Melamed S, et
531 al. Bacteria deplete deoxynucleotides to defend against bacteriophage
532 infection. *Nat Microbiol*. 2022;7: 1200–1209.
- 533 11. Whiteley AT, Eaglesham JB, de Oliveira Mann CC, Morehouse BR, Lowey
534 B, Nieminen EA, et al. Bacterial cGAS-like enzymes synthesize diverse
535 nucleotide signals. *Nature*. 2019;567: 194–199.
- 536 12. Kaur G, Burroughs AM, Iyer LM, Aravind L. Highly regulated, diversifying
537 NTP-dependent biological conflict systems with implications for the
538 emergence of multicellularity. *Elife*. 2020;9. doi:10.7554/eLife.52696
- 539 13. Wein T, Johnson AG, Millman A, Lange K, Yirmiya E, Hadary R, et al.
540 CARD-like domains mediate anti-phage defense in bacterial gasdermin
541 systems. *bioRxiv*. 2023. p. 2023.05.28.542683.
542 doi:10.1101/2023.05.28.542683
- 543 14. Burroughs AM, Aravind L. Identification of Uncharacterized Components of
544 Prokaryotic Immune Systems and Their Diverse Eukaryotic Reformulations.
545 *J Bacteriol*. 2020;202. doi:10.1128/JB.00365-20
- 546 15. Richter DJ, Berney C, Strassert JFH, Poh Y-P, Herman EK, Muñoz-Gómez
547 SA, et al. EukProt: a database of genome-scale predicted proteins across
548 the diversity of eukaryotes. *bioRxiv*. bioRxiv; 2020.
549 doi:10.1101/2020.06.30.180687

- 550 16. del Campo J, Sieracki ME, Molestina R, Keeling P, Massana R, Ruiz-Trillo I.
551 The others: our biased perspective of eukaryotic genomes. *Trends Ecol*
552 *Evol.* 2014;29: 252–259.
- 553 17. Shaw AE, Hughes J, Gu Q, Behdenna A, Singer JB, Dennis T, et al.
554 Fundamental properties of the mammalian innate immune system revealed
555 by multispecies comparison of type I interferon responses. *PLoS Biol.*
556 2017;15: e2004086.
- 557 18. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, et al.
558 Pfam: the protein families database. *Nucleic Acids Res.* 2014;42: D222–30.
- 559 19. Eddy SR. Accelerated Profile HMM Searches. *PLoS Comput Biol.* 2011;7:
560 e1002195.
- 561 20. Eddy SR. Profile hidden Markov models. *Bioinformatics.* 1998;14: 755–763.
- 562 21. Demuth JP, De Bie T, Stajich JE, Cristianini N, Hahn MW. The evolution of
563 mammalian gene families. *PLoS One.* 2006;1: e85.
- 564 22. Albalat R, Cañestro C. Evolution by gene loss. *Nat Rev Genet.* 2016;17:
565 379–391.
- 566 23. Wang X, Grus WE, Zhang J. Gene losses during human origins. *PLoS Biol.*
567 2006;4: e52.
- 568 24. Koumandou VL, Wickstead B, Ginger ML, van der Giezen M, Dacks JB,
569 Field MC. Molecular paleontology and complexity in the last eukaryotic
570 common ancestor. *Crit Rev Biochem Mol Biol.* 2013;48: 373–396.
- 571 25. Prokopchuk G, Butenko A, Dacks JB, Speijer D, Field MC, Lukeš J. Lessons
572 from the deep: mechanisms behind diversification of eukaryotic protein
573 complexes. *Biol Rev Camb Philos Soc.* 2023. doi:10.1111/brv.12988
- 574 26. Farkas Z, Kovács K, Sarkadi Z, Kalapis D, Fekete G, Birtyik F, et al. Gene
575 loss and compensatory evolution promotes the emergence of morphological
576 novelties in budding yeast. *Nat Ecol Evol.* 2022;6: 763–773.
- 577 27. Kranzusch PJ. cGAS and CD-NTase enzymes: structure, mechanism, and
578 evolution. *Curr Opin Struct Biol.* 2019;59: 178–187.
- 579 28. Li Y, Slavik KM, Toyoda HC, Morehouse BR, de Oliveira Mann CC, Elek A,
580 et al. cGLRs are a diverse family of pattern recognition receptors in innate
581 immunity. *Cell.* 2023. doi:10.1016/j.cell.2023.05.038

- 582 29. de Oliveira Mann CC, Kiefersauer R, Witte G, Hopfner K-P. Structural and
583 biochemical characterization of the cell fate determining
584 nucleotidyltransferase fold protein MAB21L1. *Sci Rep.* 2016;6: 27498.
- 585 30. Chow KL, Hall DH, Emmons SW. The mab-21 gene of *Caenorhabditis*
586 *elegans* encodes a novel protein required for choice of alternate cell fates.
587 *Development.* 1995;121: 3615–3626.
- 588 31. Yamada R, Mizutani-Koseki Y, Hasegawa T, Osumi N, Koseki H, Takahashi
589 N. Cell-autonomous involvement of Mab21l1 is essential for lens placode
590 development. *Development.* 2003;130: 1759–1770.
- 591 32. Keating SE, Baran M, Bowie AG. Cytosolic DNA sensors regulating type I
592 interferon induction. *Trends Immunol.* 2011;32: 574–581.
- 593 33. Hornung V, Latz E. Intracellular DNA recognition. *Nat Rev Immunol.*
594 2010;10: 123–130.
- 595 34. Ablasser A, Goldeck M, Cavlar T, Deimling T, Witte G, Röhl I, et al. cGAS
596 produces a 2'-5'-linked cyclic dinucleotide second messenger that activates
597 STING. *Nature.* 2013;498: 380–384.
- 598 35. Gao P, Ascano M, Wu Y, Barchet W, Gaffney BL, Zillinger T, et al. Cyclic
599 [G(2',5')pA(3',5')p] Is the Metazoan Second Messenger Produced by DNA-
600 Activated Cyclic GMP-AMP Synthase. *Cell.* 2013. pp. 1094–1107.
601 doi:10.1016/j.cell.2013.04.046
- 602 36. Sun L, Wu J, Du F, Chen X, Chen ZJ. Cyclic GMP-AMP synthase is a
603 cytosolic DNA sensor that activates the type I interferon pathway. *Science.*
604 2013;339: 786–791.
- 605 37. Burdette DL, Monroe KM, Sotelo-Troha K, Iwig JS, Eckert B, Hyodo M, et al.
606 STING is a direct innate immune sensor of cyclic di-GMP. *Nature.* 2011;478:
607 515–518.
- 608 38. Gui X, Yang H, Li T, Tan X, Shi P, Li M, et al. Autophagy induction via
609 STING trafficking is a primordial function of the cGAS pathway. *Nature.*
610 2019;567: 262–266.
- 611 39. Dong B, Silverman RH. A bipartite model of 2-5A-dependent RNase L. *J Biol*
612 *Chem.* 1997;272: 22236–22242.
- 613 40. Silverman RH. Viral encounters with 2',5'-oligoadenylate synthetase and
614 RNase L during the interferon antiviral response. *J Virol.* 2007;81: 12720–

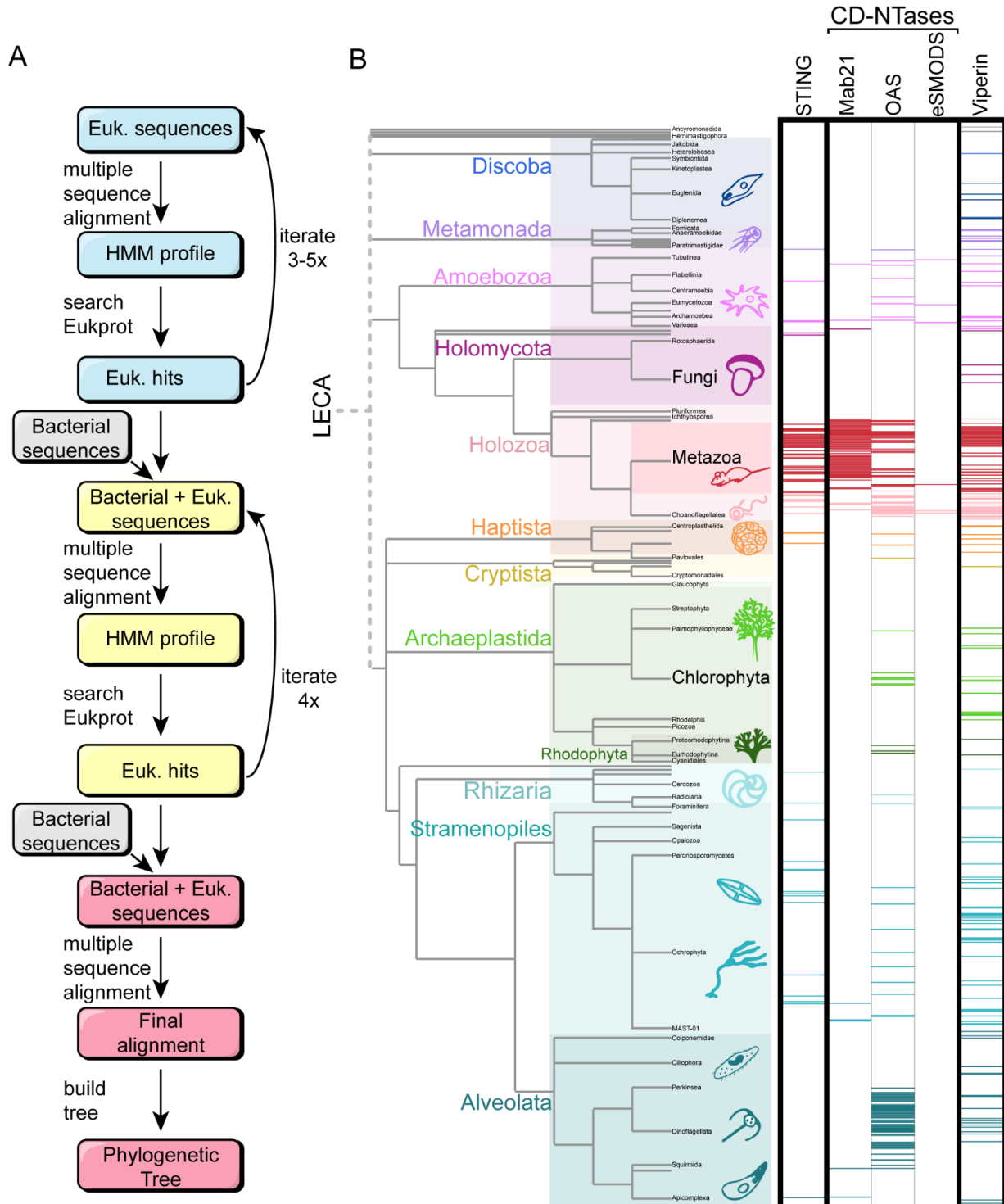
- 615 12729.
- 616 41. Kristiansen H, Gad HH, Eskildsen-Larsen S, Despres P, Hartmann R. The
617 oligoadenylate synthetase family: an ancient protein family with multiple
618 antiviral activities. *J Interferon Cytokine Res.* 2011;31: 41–47.
- 619 42. Severin GB, Ramliden MS, Hawver LA, Wang K, Pell ME, Kieninger A-K, et
620 al. Direct activation of a phospholipase by cyclic GMP-AMP in El Tor *Vibrio*
621 *cholerae*. *Proc Natl Acad Sci U S A.* 2018;115: E6048–E6055.
- 622 43. Millman A, Melamed S, Amitai G, Sorek R. Diversity and classification of
623 cyclic-oligonucleotide-based anti-phage signalling systems. *Nat Microbiol.*
624 2020;5: 1608–1615.
- 625 44. Tsang WH, Shek KF, Lee TY, Chow KL. An evolutionarily conserved nested
626 gene pair - Mab21 and Lrba/Nbea in metazoan. *Genomics.* 2009;94: 177–
627 187.
- 628 45. Wu X, Wu F-H, Wang X, Wang L, Siedow JN, Zhang W, et al. Molecular
629 evolutionary and structural analysis of the cytosolic DNA sensor cGAS and
630 STING. *Nucleic Acids Res.* 2014;42: 8243–8257.
- 631 46. Woznica A, Kumar A, Sturge CR, Xing C, King N, Pfeiffer JK. STING
632 mediates immune responses in the closest living relatives of animals. *Elife.*
633 2021;10. doi:10.7554/eLife.70436
- 634 47. Ye Q, Lau RK, Mathews IT, Birkholz EA, Watrous JD, Azimi CS, et al.
635 HORMA Domain Proteins and a Trip13-like ATPase Regulate Bacterial
636 cGAS-like Enzymes to Mediate Bacteriophage Immunity. *Molecular Cell.*
637 2020. pp. 709–722.e7. doi:10.1016/j.molcel.2019.12.009
- 638 48. Govande AA, Duncan-Lowey B, Eaglesham JB, Whiteley AT, Kranzusch PJ.
639 Molecular basis of CD-NTase nucleotide selection in CBASS anti-phage
640 defense. *Cell Rep.* 2021;35: 109206.
- 641 49. Hogrel G, Guild A, Graham S, Rickman H, Grüşchow S, Bertrand Q, et al.
642 Author Correction: Cyclic nucleotide-induced helical structure activates a TIR
643 immune effector. *Nature.* 2023;614: E15.
- 644 50. Du Y, Hu Z, Luo Y, Wang HY, Yu X, Wang R-F. Function and regulation of
645 cGAS-STING signaling in infectious diseases. *Front Immunol.* 2023;14:
646 1130423.
- 647 51. Ishikawa H, Barber GN. STING is an endoplasmic reticulum adaptor that

- 648 facilitates innate immune signalling. *Nature*. 2008;455: 674–678.
- 649 52. Nishimura MT. TIR Domains of Plant Immune Receptors are NAD plus
650 Consuming Enzymes that Promote Cell Death. *MOLECULAR PLANT-*
651 *MICROBE INTERACTIONS*. AMER PHYTOPATHOLOGICAL SOC 3340
652 PILOT KNOB ROAD, ST PAUL, MN 55121 USA; 2019. pp. 217–218.
- 653 53. Horsefield S, Burdett H, Zhang X, Manik MK, Shi Y, Chen J, et al. NAD+
654 cleavage activity by animal and plant TIR domains in cell death pathways.
655 *Science*. 2019;365: 793–799.
- 656 54. Margolis SR, Wilson SC, Vance RE. Evolutionary Origins of cGAS-STING
657 Signaling. *Trends Immunol*. 2017;38: 733–743.
- 658 55. Zhang G, Fang X, Guo X, Li L, Luo R, Xu F, et al. The oyster genome
659 reveals stress adaptation and complexity of shell formation. *Nature*.
660 2012;490: 49–54.
- 661 56. Gizzi AS, Grove TL, Arnold JJ, Jose J, Jangra RK, Garforth SJ, et al. A
662 naturally occurring antiviral ribonucleotide encoded by the human genome.
663 *Nature*. 2018;558: 610–614.
- 664 57. Fenwick MK, Li Y, Cresswell P, Modis Y, Ealick SE. Structural studies of
665 viperin, an antiviral radical SAM enzyme. *Proc Natl Acad Sci U S A*.
666 2017;114: 6806–6811.
- 667 58. Rivera-Serrano EE, Gizzi AS, Arnold JJ, Grove TL, Almo SC, Cameron CE.
668 Viperin Reveals Its True Function. *Annu Rev Virol*. 2020;7: 421–446.
- 669 59. Seifert M, Bera SC, van Nies P, Kirchdoerfer RN, Shannon A, Le T-T-N, et
670 al. Inhibition of SARS-CoV-2 polymerase by nucleotide analogs from a
671 single-molecule perspective. *Elife*. 2021;10. doi:10.7554/eLife.70968
- 672 60. Lachowicz JC, Gizzi AS, Almo SC, Grove TL. Structural Insight into the
673 Substrate Scope of Viperin and Viperin-like Enzymes from Three Domains of
674 Life. *Biochemistry*. 2021;60: 2116–2129.
- 675 61. Eme L, Gentekaki E, Curtis B, Archibald JM, Roger AJ. Lateral Gene
676 Transfer in the Adaptation of the Anaerobic Parasite *Blastocystis* to the Gut.
677 *Curr Biol*. 2017;27: 807–820.
- 678 62. Richards TA, Soanes DM, Jones MDM, Vasieva O, Leonard G, Paszkiewicz
679 K, et al. Horizontal gene transfer facilitated the evolution of plant parasitic
680 mechanisms in the oomycetes. *Proc Natl Acad Sci U S A*. 2011;108: 15258–

- 681 15263.
- 682 63. Gabaldón T. Patterns and impacts of nonvertical evolution in eukaryotes: a
683 paradigm shift. *Ann N Y Acad Sci.* 2020;1476: 78–92.
- 684 64. Leger MM, Eme L, Stairs CW, Roger AJ. Demystifying Eukaryote Lateral
685 Gene Transfer (Response to Martin 2017 DOI: 10.1002/bies.201700115).
686 *Bioessays.* 2018;40: e1700242.
- 687 65. Bernheim A, Sorek R. The pan-immune system of bacteria: antiviral defence
688 as a community resource. *Nat Rev Microbiol.* 2020;18: 113–119.
- 689 66. Koonin EV, Makarova KS, Wolf YI. Evolutionary Genomics of Defense
690 Systems in Archaea and Bacteria. *Annu Rev Microbiol.* 2017;71: 233–261.
- 691 67. van Houte S, Buckling A, Westra ER. Evolutionary Ecology of Prokaryotic
692 Immune Mechanisms. *Microbiol Mol Biol Rev.* 2016;80: 745–763.
- 693 68. Hochhauser D, Millman A, Sorek R. The defense island repertoire of the
694 *Escherichia coli* pan-genome. *PLoS Genet.* 2023;19: e1010694.
- 695 69. LeGault KN, Hays SG, Angermeyer A, McKitterick AC, Johura F-T, Sultana
696 M, et al. Temporal shifts in antibiotic resistance elements govern phage-
697 pathogen conflicts. *Science.* 2021;373. doi:10.1126/science.abg2166
- 698 70. Keeling PJ, Palmer JD. Horizontal gene transfer in eukaryotic evolution. *Nat*
699 *Rev Genet.* 2008;9: 605–618.
- 700 71. Katoh K, Standley DM. MAFFT multiple sequence alignment software
701 version 7: improvements in performance and usability. *Mol Biol Evol.*
702 2013;30: 772–780.
- 703 72. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for
704 automated alignment trimming in large-scale phylogenetic analyses.
705 *Bioinformatics.* 2009;25: 1972–1973.
- 706 73. Price MN, Dehal PS, Arkin AP. FastTree 2--approximately maximum-
707 likelihood trees for large alignments. *PLoS One.* 2010;5: e9490.
- 708 74. Chernomor O, Minh BQ, Forest F, Klaere S, Ingram T, Henzinger M, et al.
709 Split diversity in constrained conservation prioritization using integer linear
710 programming. *Methods in Ecology and Evolution.* 2015. pp. 83–91.
711 doi:10.1111/2041-210x.12299
- 712 75. Edgar RC. MUSCLE v5 enables improved estimates of phylogenetic tree

- 713 confidence by ensemble bootstrapping. bioRxiv. 2021. p.
714 2021.06.20.449169. doi:10.1101/2021.06.20.449169
- 715 76. Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von
716 Haeseler A, et al. IQ-TREE 2: New Models and Efficient Methods for
717 Phylogenetic Inference in the Genomic Era. *Mol Biol Evol.* 2020;37: 1530–
718 1534.
- 719 77. Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. RAxML-NG: a fast,
720 scalable and user-friendly tool for maximum likelihood phylogenetic
721 inference. *Bioinformatics.* 2019;35: 4453–4455.
- 722 78. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v5: an online tool for
723 phylogenetic tree display and annotation. *Nucleic Acids Res.* 2021;49:
724 W293–W296.
- 725 79. Hulsen T. DeepVenn -- a web application for the creation of area-
726 proportional Venn diagrams using the deep learning framework
727 Tensorflow.js. arXiv [cs.HC]. 2022. Available:
728 <http://arxiv.org/abs/2210.04597>
- 729 80. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al.
730 Highly accurate protein structure prediction with AlphaFold. *Nature.*
731 2021;596: 583–589.
- 732 81. Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, Yordanova G, et
733 al. AlphaFold Protein Structure Database: massively expanding the
734 structural coverage of protein-sequence space with high-accuracy models.
735 *Nucleic Acids Res.* 2022;50: D439–D444.
- 736 82. Shang G, Zhang C, Chen ZJ, Bai X-C, Zhang X. Cryo-EM structures of
737 STING reveal its mechanism of activation by cyclic GMP–AMP. *Nature.*
738 2019;567: 389–393.

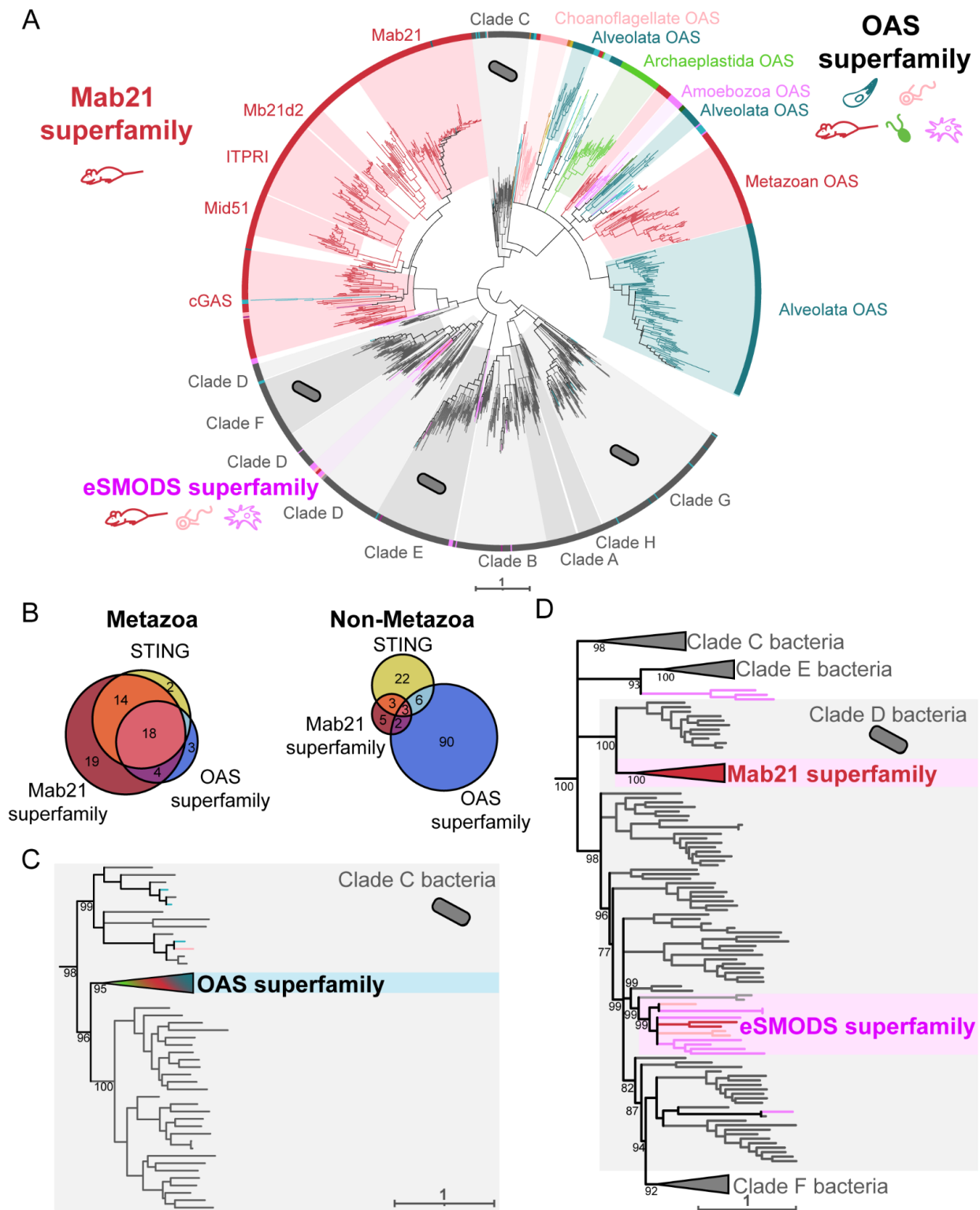
739 **Figures**



740

741 **Figure 1: HMM searches to find homologs across the eukaryotic Tree of Life**

742 (A) A schematic of the HMM search process. Starting from initial, animal-dominated HMM
743 profiles for each protein family, we used iterative HMM searches of the EukProt database to
744 generate pan-eukaryotic HMMs. These were combined with bacterial sequences to enable
745 discovery of bacteria-like homologs in eukaryotes. Each set of searches was repeated 3-5 times
746 until few or no additional eukaryotic sequences were recovered. (B) Phylogenetic tree of
747 eukaryotes, with major supergroups color-coded. The height of the colored rectangles for each
748 group is proportional to its species representation in EukProt. Horizontal, colored bars mark
749 each eukaryotic species in which we found homologs of STINGs, CD-NTases, or Viperins.
750 White space indicates species where we did not recover any homologs. The CD-NTase hits are
751 divided into the three eukaryotic superfamilies defined in Fig. 2.



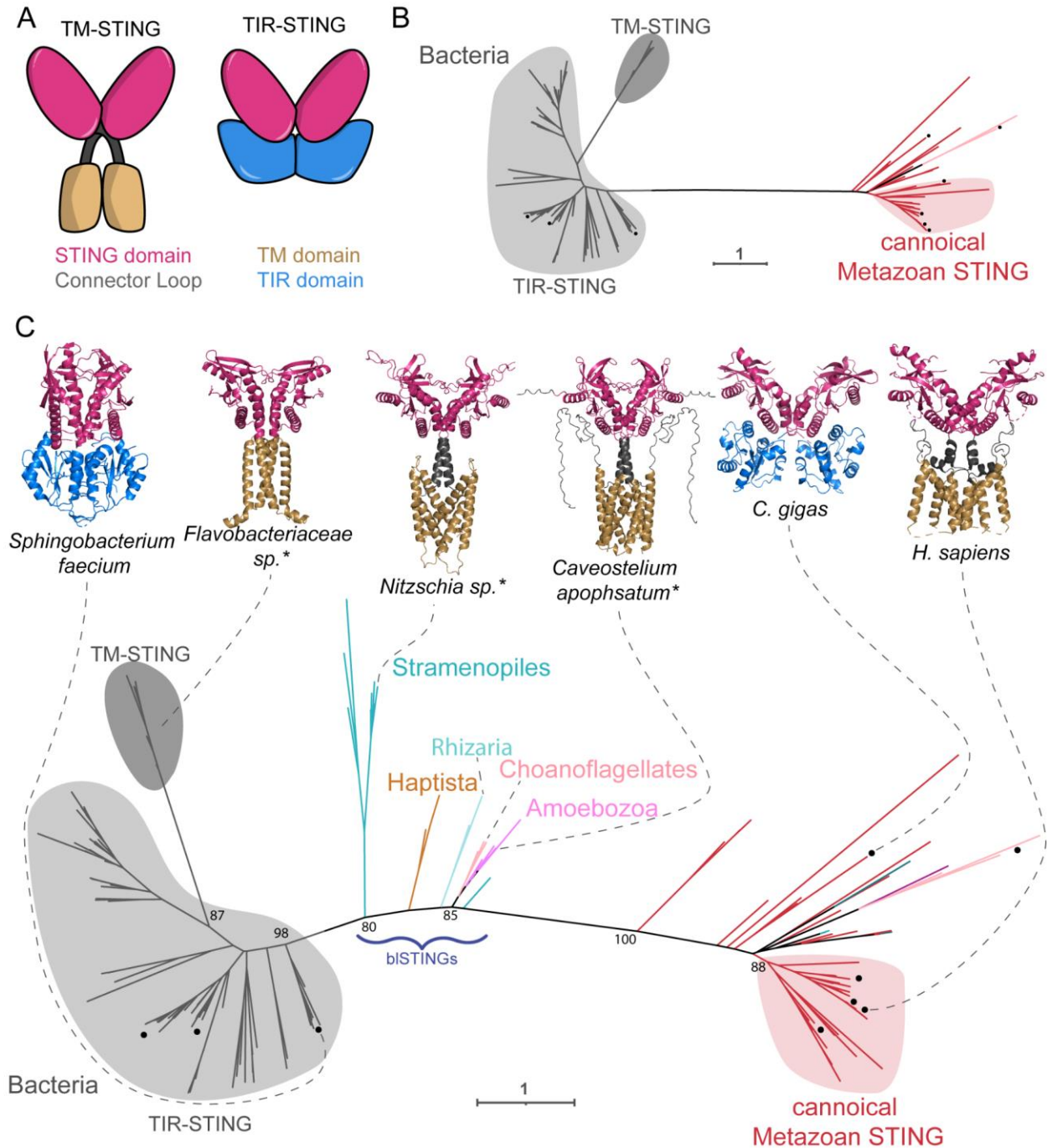
752

753 **Figure 2: Independent HGT events gave rise to multiple CD-NTase superfamilies**

754 (A) Maximum likelihood phylogenetic tree generated by IQtree of CD-NTases spanning

755 eukaryotic and bacterial diversity. The Mab21 superfamily (red, top left) is largely an animal-

756 specific innovation, with many paralogs including cGAS. In contrast, most other eukaryotic
757 lineages encode CD-NTases from the OAS superfamily (multicolor, top right). The relatively
758 small eSMODS superfamily (pink, bottom left) is a recent HGT between clade D bacteria and
759 eukaryotes. Bacterial CD-NTase sequences shown in gray. Eukaryotic sequences are colored
760 according to eukaryotic group as in Fig. 1. The tree is arbitrarily rooted on a branch separating
761 clades A, B, G, and H, which did not typically have associated eukaryotic sequences, from the
762 rest of the bacterial CD-NTases. (B) Venn diagrams showing the number of species where we
763 detected at least one homolog of STING, Mab21 superfamily CD-NTases, and/or OAS
764 superfamily CD-NTases in Metazoa (left) or non-metazoan eukaryotes (right). (C) Within the
765 CD-NTase phylogenetic tree in A, the OAS superfamily branches within clade C bacterial CD-
766 NTases (gray branches). (D) Clade D CD-NTases (gray branches) have been horizontally
767 transferred into eukaryotes multiple times, giving rise to both the Mab21 superfamily and the
768 eSMODS superfamily. Ultrafast bootstraps determined by IQtree shown at key nodes. See
769 Supplementary Figure 4 for full CD-NTase phylogenetic tree.

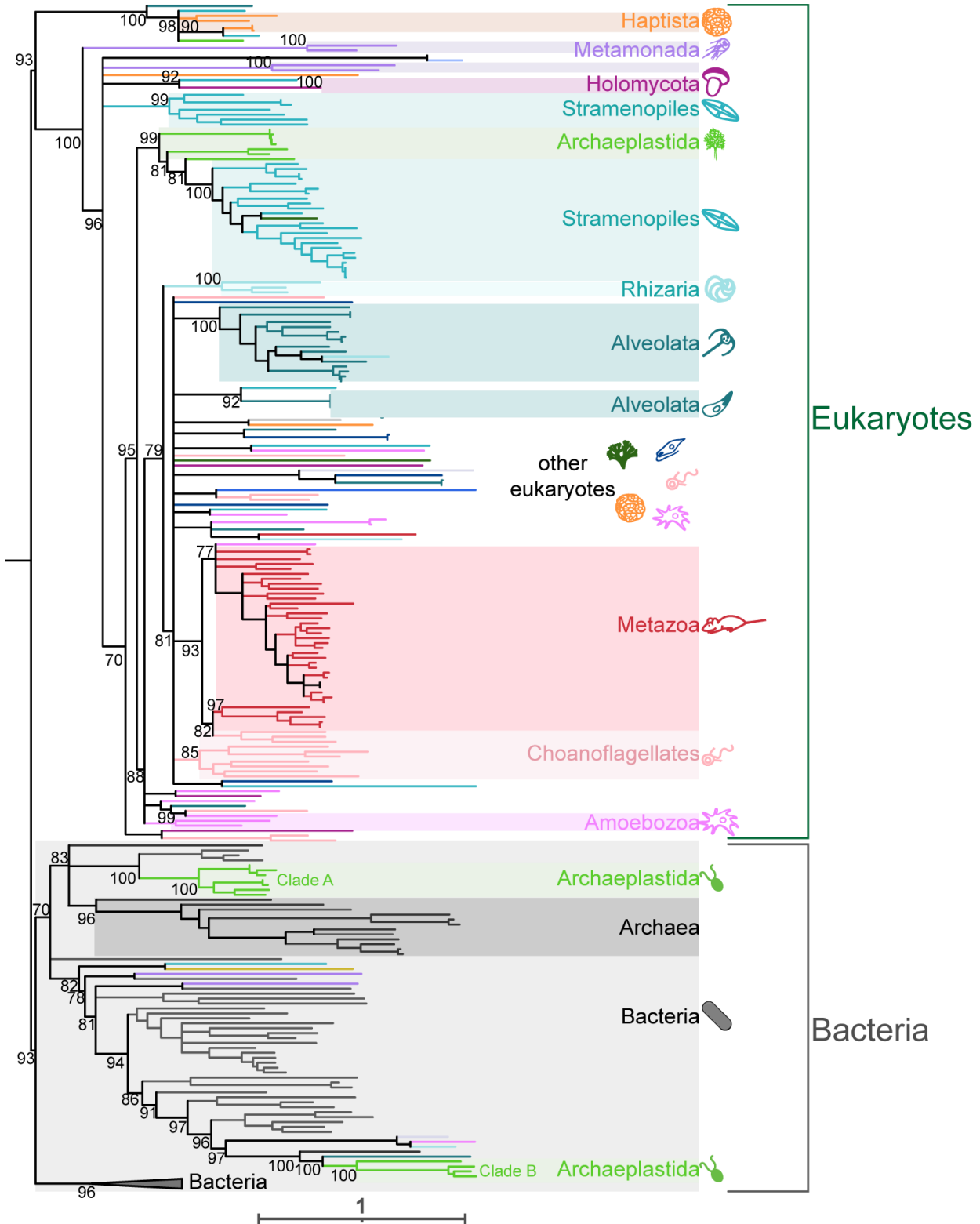


770
771
772
773
774
775
776
777
778
779

Figure 3: Diverse eukaryotic STING proteins bridge the gap between metazoans and bacteria

(A) Graphical depiction of common domain architectures of STING proteins. (B) Maximum likelihood unrooted phylogenetic tree of STING domains from Metazoa and bacteria, which are separated by one long branch. Black dot (•) indicates proteins that have been previously experimentally characterized. Bacterial sequences are in gray and animal sequences are in red. (C) Maximum likelihood unrooted phylogenetic tree of hits from iterative HMM searches for diverse eukaryotic STING domains. The STING domains from bacteria-like STINGs (bISTINGs) from diverse eukaryotes break up the long branch between bacterial and animal STINGs.

780 Structures of the indicated STING proteins are shown above, with those predicted by AlphaFold
781 indicated by an asterisk. Homologs with X-ray crystal structures are from[3,82]. Colored regions
782 show two domain architectures in bacteria and eukaryotes (STING linked to a TIR domain and
783 STING linked to a transmembrane domain), each of which have evolved convergently in
784 bacteria and eukaryotes. Ultrafast bootstraps determined by IQtree shown at key nodes. See
785 Supplementary Figure 5 for full STING phylogenetic tree.



786

787

Figure 4: Viperin is a deeply conserved innate immune module

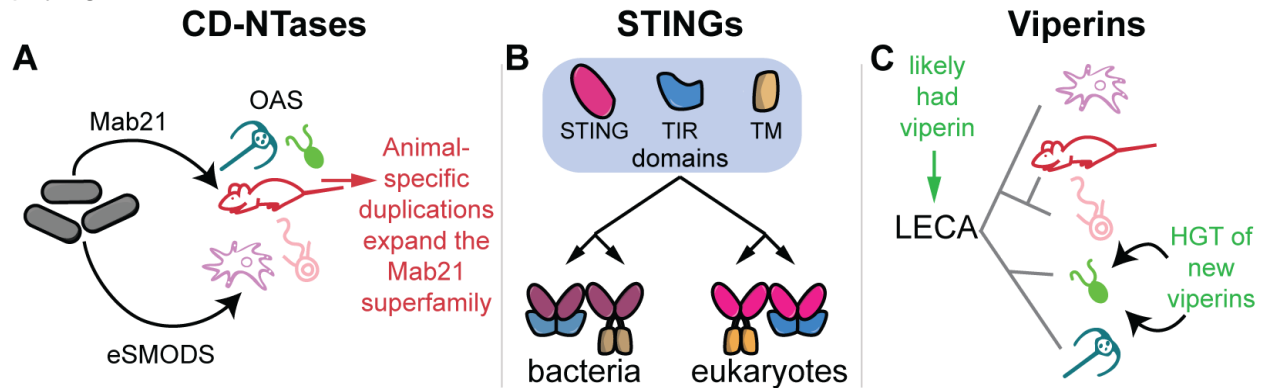
788

(A) Maximum likelihood phylogenetic tree generated by IQtree of Viperins from eukaryotes,

789

bacteria, and archaea. All major eukaryotic supergroups have at least two species that encode

790 a Viperin homolog (colored supergroups). Bacterial Viperin sequences shown in gray and
791 archaeal sequences in dark gray. There are two clades of Chloroplastida (a group within
792 Archaeplastida) sequences that branch robustly within the bacteria clade. Ultrafast bootstraps
793 determined by IQtree shown at key nodes. Tree is arbitrarily rooted between the major
794 eukaryotic and bacterial clades. See Supplementary Figure 6 for fully annotated Viperin
795 phylogenetic tree.

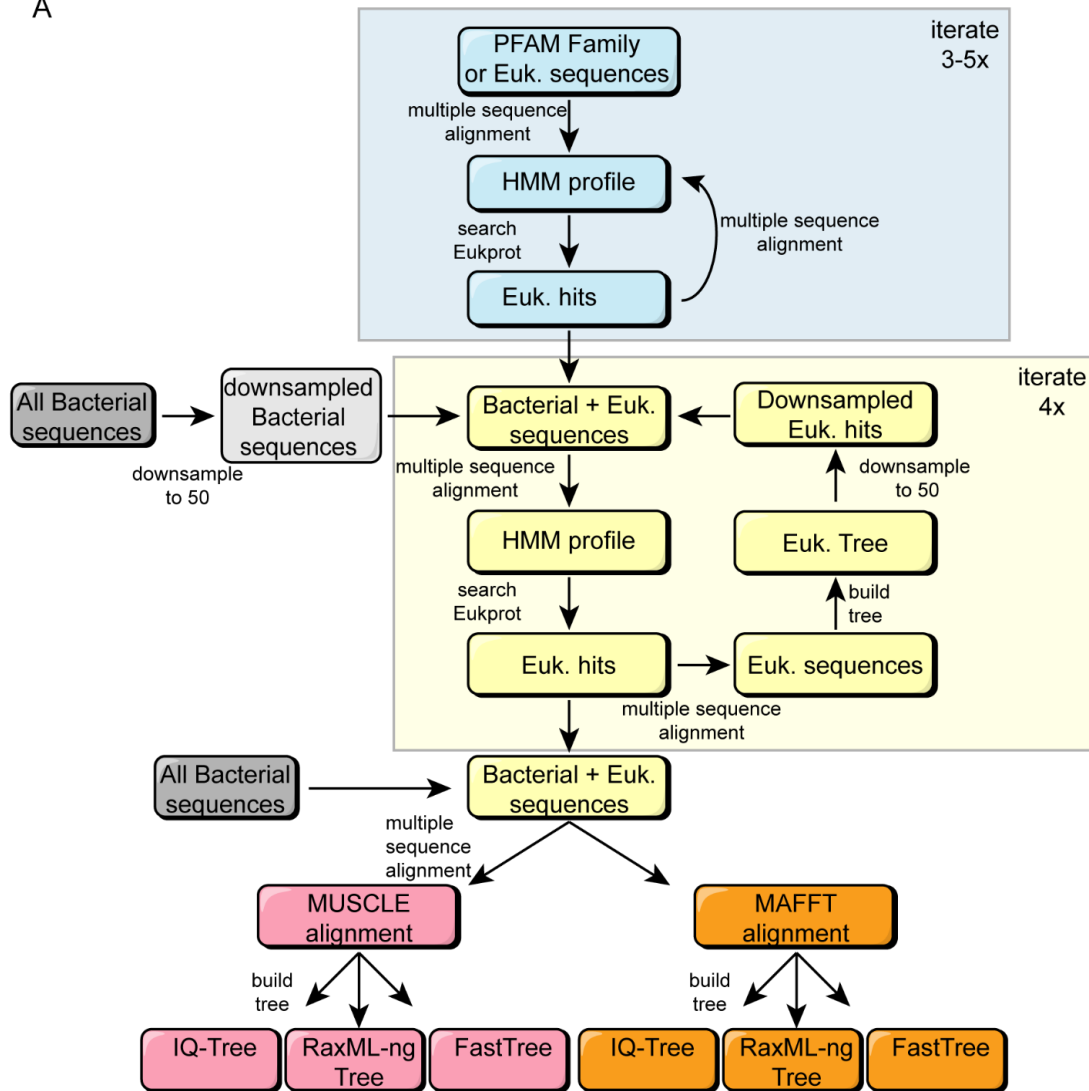


796
797

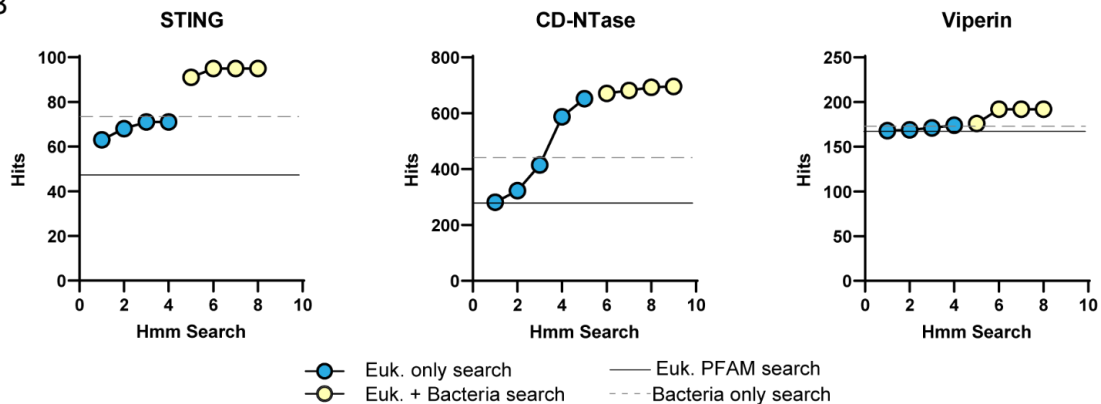
798 **Figure 5: Proposed model of evolutionary history of CD-NTases, STING, and Viperin**

799 Proposed summary models of the evolutionary history of innate immune components. (A) We
800 define two distinct superfamilies of CD-NTases that likely arose from bacteria-eukaryote HGT:
801 eSMODS and Mab21. Within the Mab21 superfamily (which contains cGAS), a number of
802 animal-specific duplications gave rise to numerous paralogs. The OAS superfamily of CD-
803 NTases are abundant across diverse eukaryotic taxa and were likely present in the LECA. (B)
804 Drawing on a shared ancient repertoire of protein domains that includes STING, TIR, and
805 transmembrane (TM) domains, bacteria and eukaryotes have convergently evolved similar
806 STING proteins through domain shuffling. (C) Viperins are widespread across the eukaryotic
807 tree and likely were present in the LECA. In addition, two sets of recent HGT events from
808 bacteria have equipped algal species with new viperins.

A



B



809

810

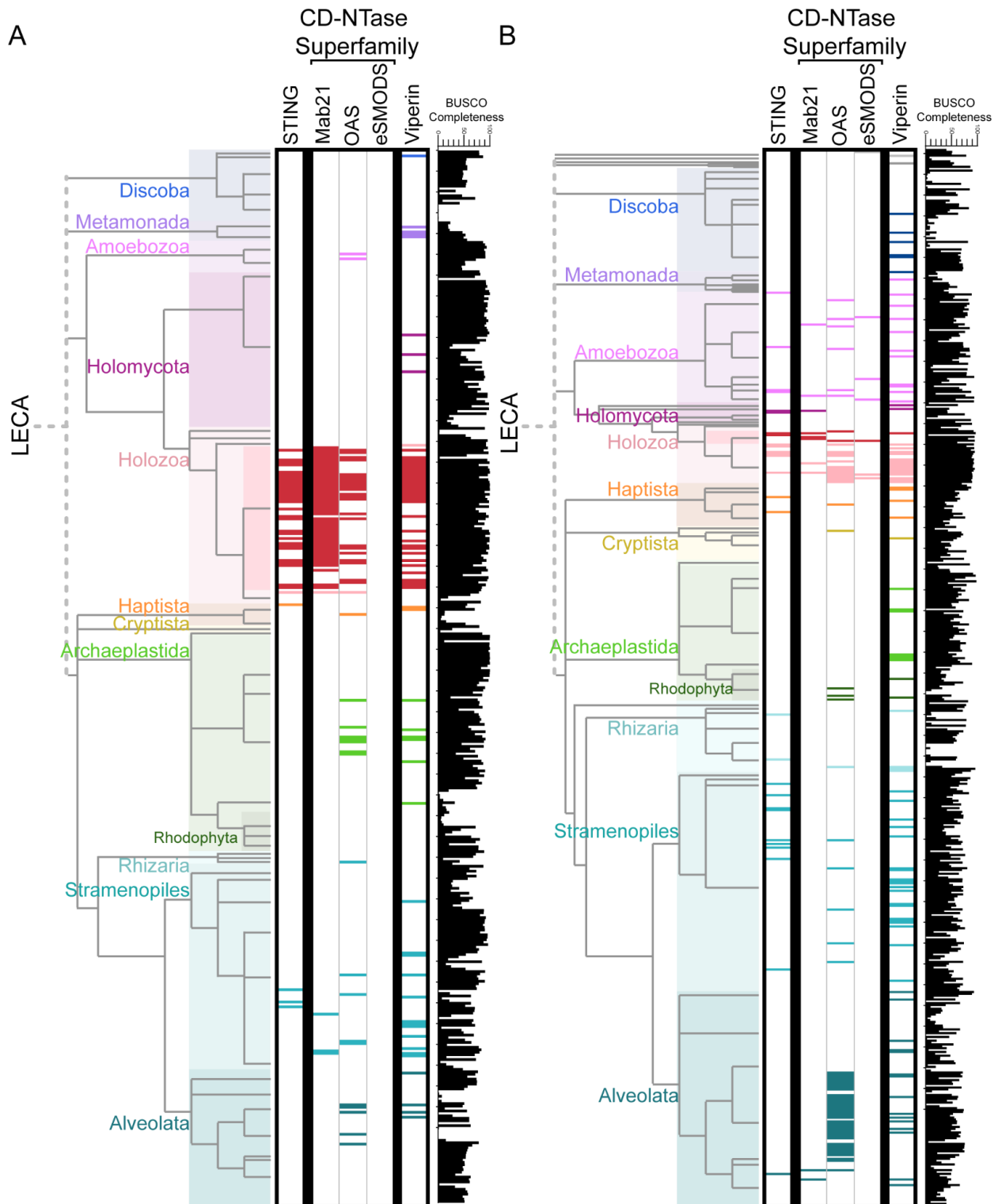
811

812

Supplementary Figure 1: Collectors curves and full search strategy

(A) Detailed schematic outlining the iterative HMM search strategy. Blue boxes and blue shaded region show eukaryotic-only searches to create pan-eukaryotic HMMs and yellow indicates

813 eukaryotic-bacterial searches to create universal HMMs. For the combined bacterial/eukaryotic
814 searches (yellow box), bacterial and eukaryotic sequences were each downsampled to 50
815 sequences (phylogenetic tree downsampled via PDA) to maintain equal contributions from
816 bacteria and eukaryotic sequences. Separately, bacterial sequences were aligned and used to
817 make an HMM which was used to search EukProt as a 'bacteria only search' and for STING we
818 searched with PF15009 for a comparable Eukaryotic PFAM search (not shown in flowchart). We
819 did this extra search for STING as PF15009 contains part of the eukaryotic STING
820 transmembrane domain and so our first search with STING was with a STING-domain-only
821 HMM (See Materials and Methods). Pink (MUSCLE) and orange (MAFFT) boxes show the final
822 alignments and phylogenetic trees that were constructed. (B) STING, CD-NTase, and Viperin
823 collector's curves showing the number of cumulative protein sequences that were found after
824 each iterative search. Results from eukaryotic searches are shown in blue and the combined
825 searches in yellow. Solid black line indicates the number of hits from the starting Pfam HMM
826 alone and the dotted gray line shows the number of hits from a bacteria-only HMM. Note that
827 some searches yielded hits that were members of more distant protein families, which were
828 later removed from the analysis and are not counted here.



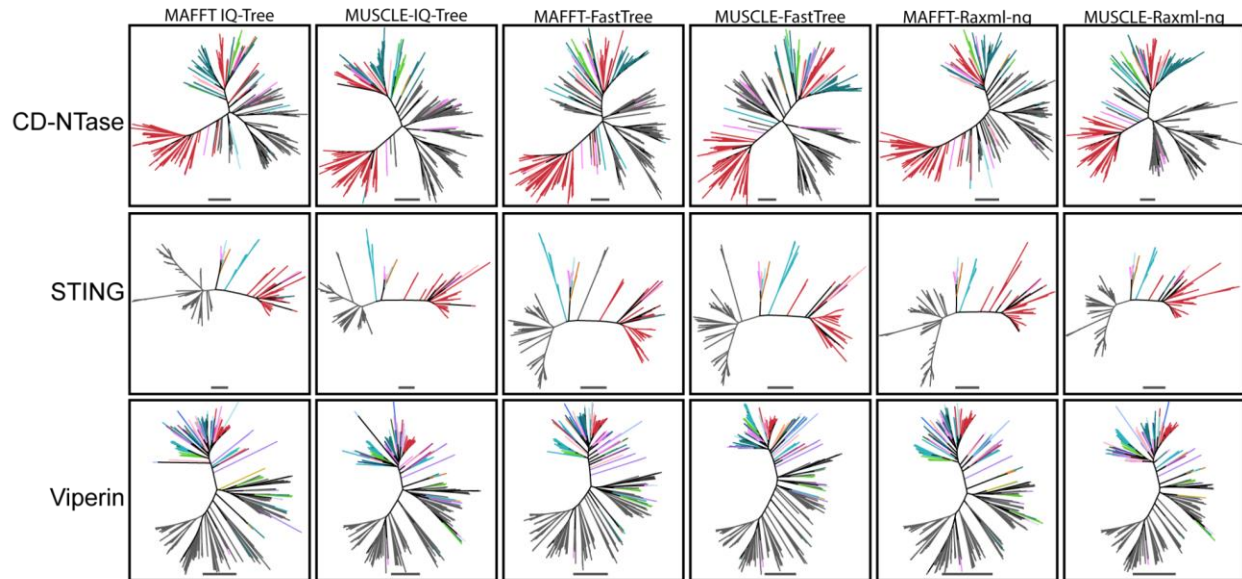
829

830 **Supplementary Figure 2: Phylogenetic trees of EukProt species by data type**

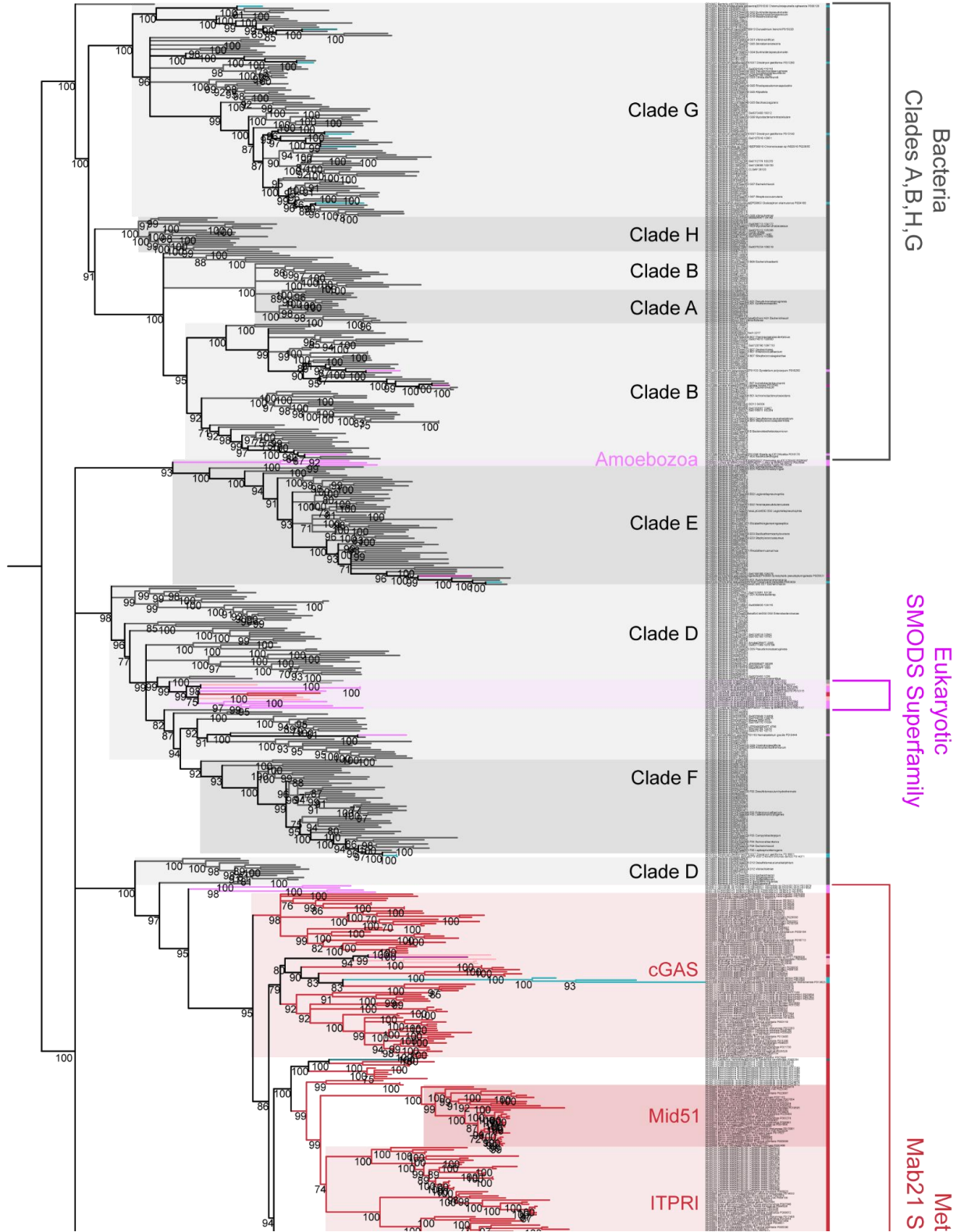
831 Phylogenetic trees derived from Fig. 1 separating species represented in EukProt v3 by

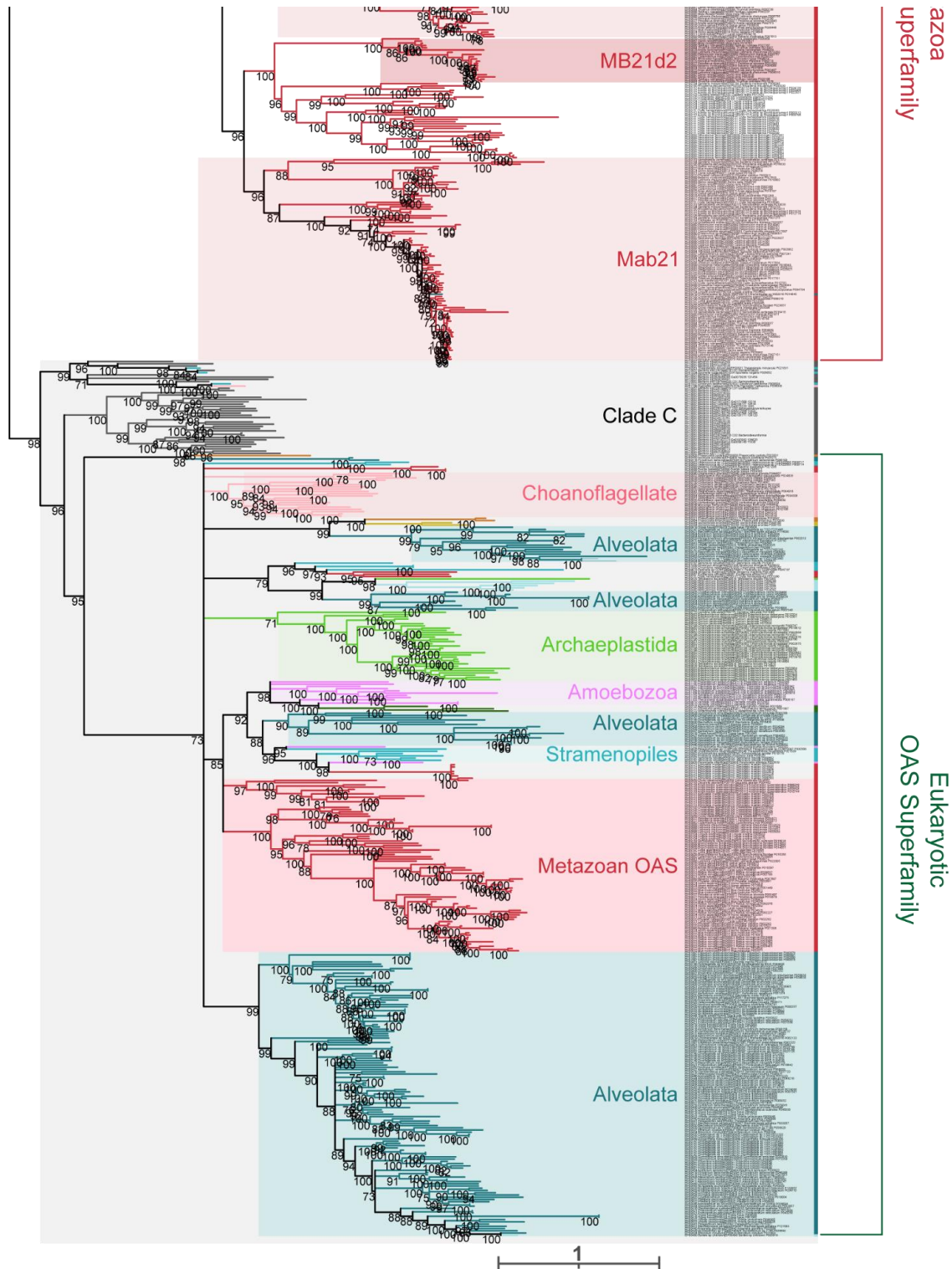
832 genomes (A) or transcriptomes (B). Supergroups are color-coded as in Figure 1. Colored bars

833 mark each eukaryotic species in which the HMM search found a homolog sequence of STING,
834 CD-NTase, or Viperin. Black bar chart shows BUSCO completeness score for each
835 genome/transcriptome.
836
837
838



839 **Supplementary Figure 3: Phylogenetic trees from different alignments and tree building**
840 **methods show robust topologies**
841 Unrooted maximum likelihood phylogenetic trees generated from two separate alignments
842 (MUSCLE and MAFFT) and with three different tree inference programs (FastTree, IQtree, and
843 RaxML-ng). Scale bar of 1 shown beneath each tree represents the number of amino acid
844 substitutions per position in the underlying alignment. Colored branches show eukaryotic
845 sequences with the same color scheme as Fig. 1, while gray lines are bacterial sequences. For
846 the majority of relationships discussed here, we recovered the same tree topology at key nodes
847 regardless of alignment or tree reconstruction algorithm used. The only exception was in the
848 STING FastTree phylogenies, wherein the TM-STING clade moved to multiple positions in the
849 phylogeny, depending on alignment algorithm used.
850



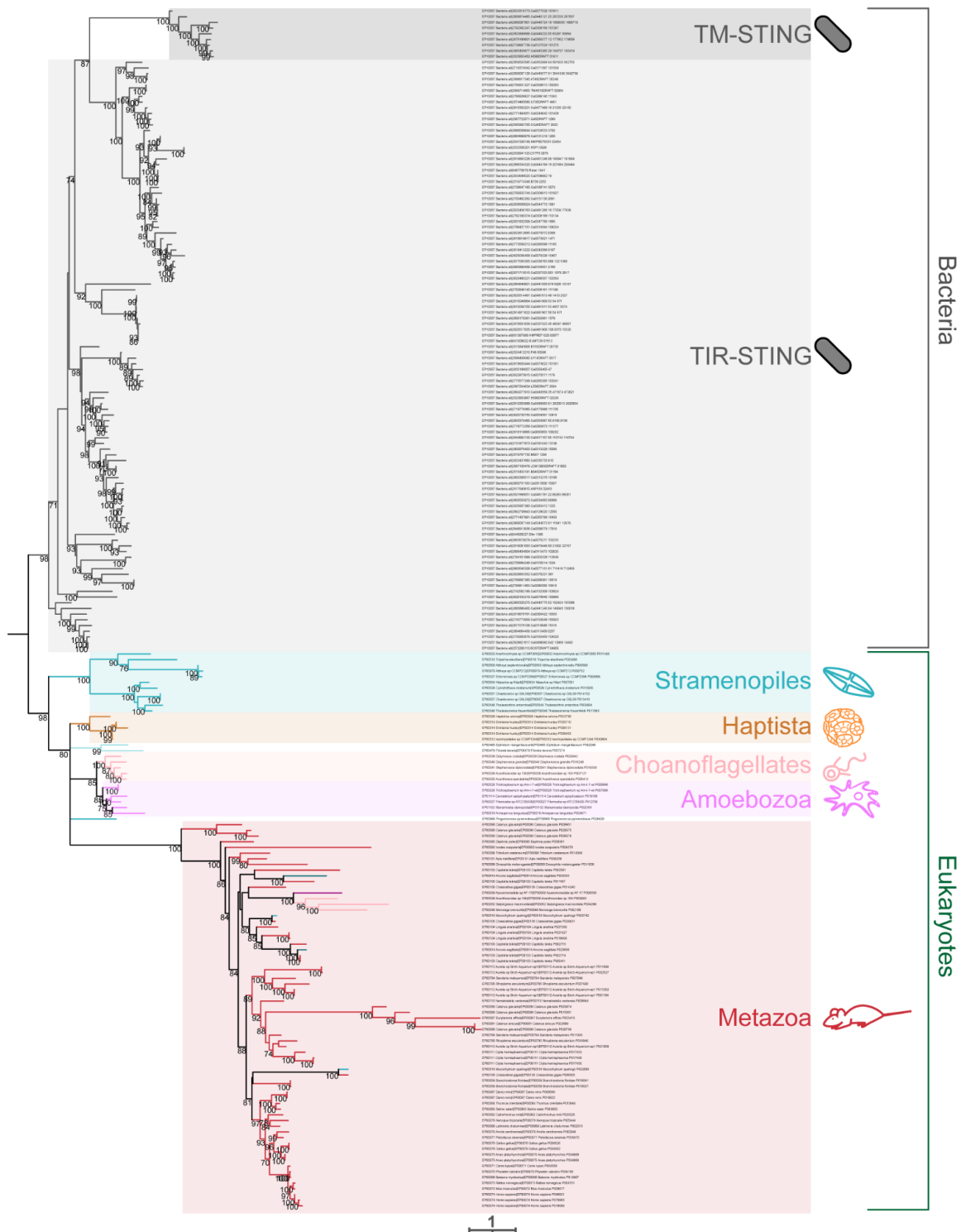


853

854 **Supplementary Figure 4: CD-NTase phylogenetic tree**

855 Maximum likelihood phylogenetic tree generated by IQtree of hits from iterative HMM searches
856 for diverse eukaryotic CD-NTases. Scale bar represents the number of amino acid substitutions

857 per position in the underlying MUSCLE alignment. Ultrafast bootstrap values calculated by
858 IQtree at all nodes with support >70 are shown. Branches with support values <70 were
859 collapsed to polytomies.
860

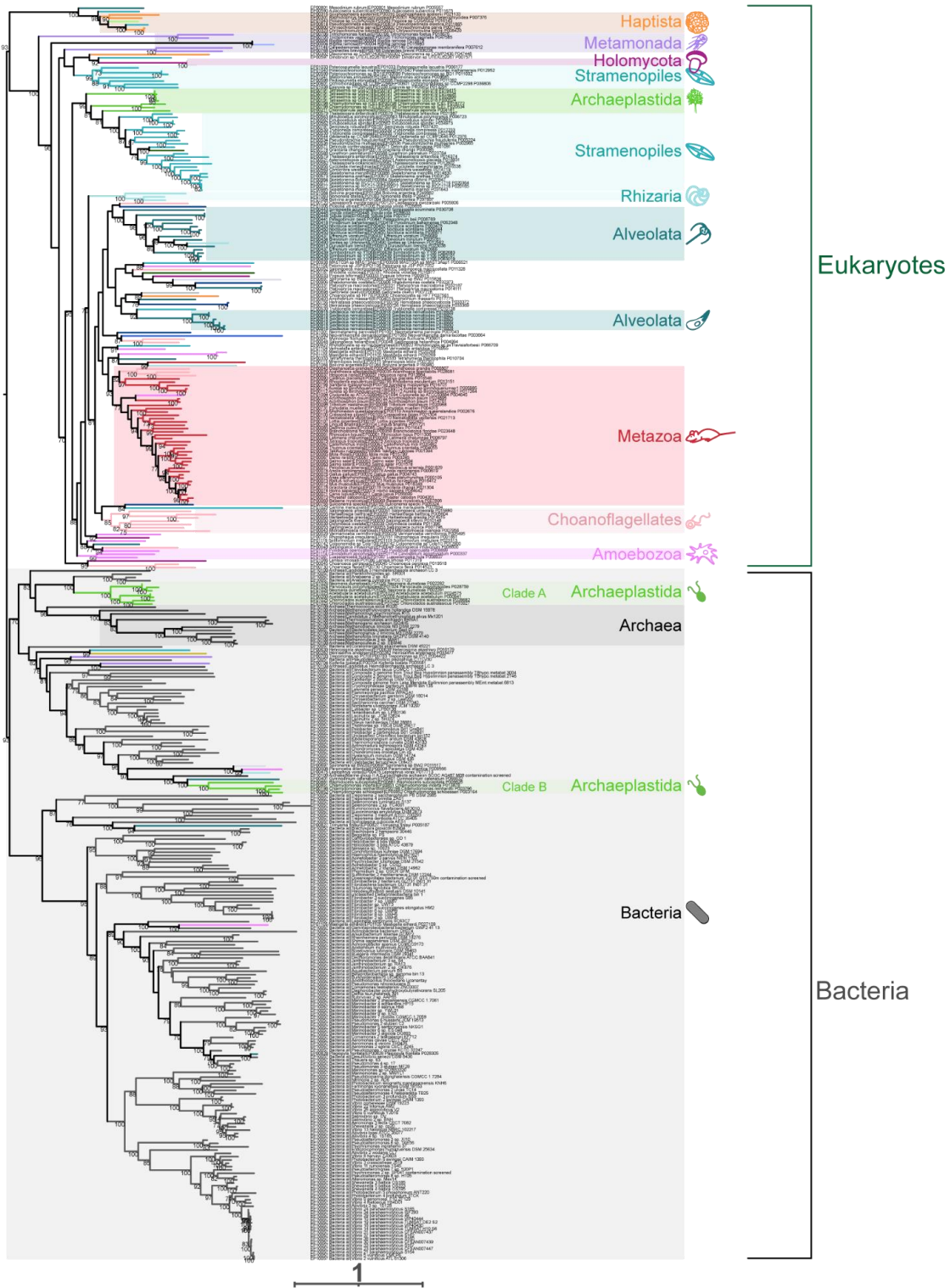


861
862
863
864

Supplementary Figure 5: STING phylogenetic tree

Maximum likelihood phylogenetic tree of hits from iterative HMM searches for diverse eukaryotic STING domains. Scale bar represents the number of amino acid substitutions per position in the

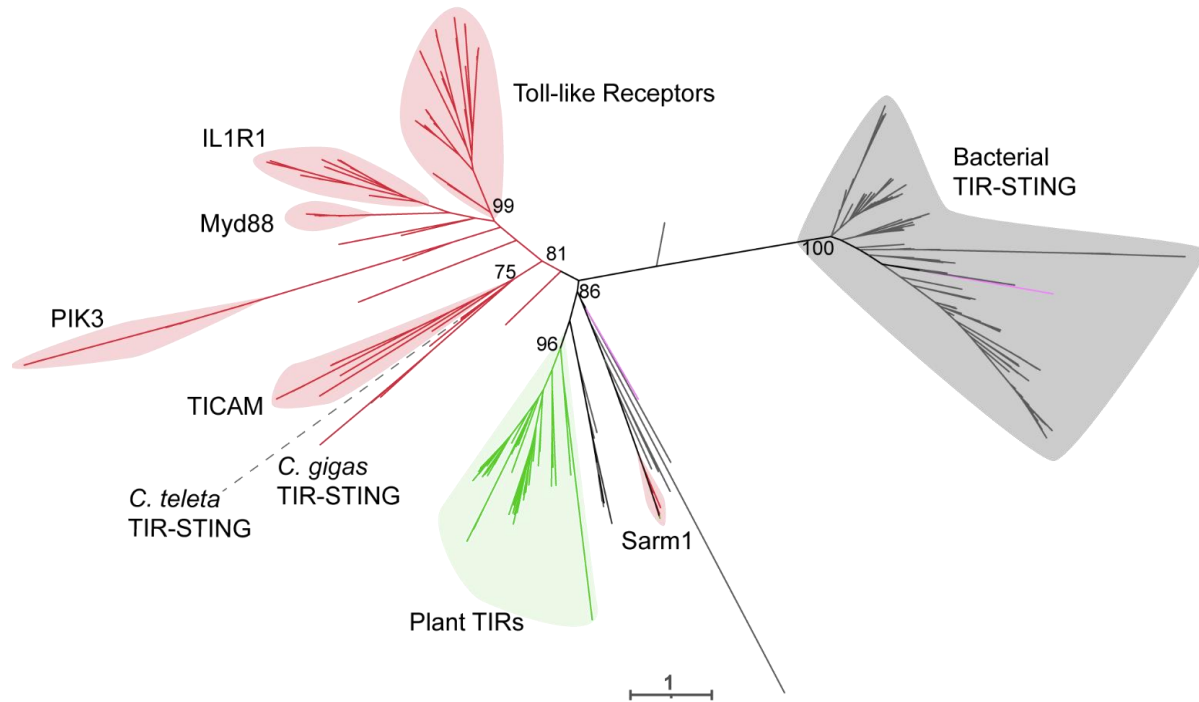
865 underlying MUSCLE alignment. Ultrafast bootstrap values calculated by IQtree at all nodes with
866 support >70 are shown. Branches with support values <70 were collapsed to polytomies.



867

868 **Supplementary Figure 6: Viperin phylogenetic tree**

869 Maximum likelihood phylogenetic tree generated by IQtree of hits from iterative HMM searches
870 for diverse eukaryotic Viperins. Scale bar represents the number of amino acid substitutions per
871 position in the underlying MUSCLE alignment. Ultrafast bootstrap values calculated by IQtree at
872 all nodes with support >70 are shown. Branches with support values <70 were collapsed to
873 polytomies.
874



875 **Supplementary Figure 7: TIR domain of *Crassostrea gigas*' TIR-STING is closely related**
876 **to metazoan TIR domains**
877

878 Unrooted maximum likelihood tree of diverse TIR domains. Scale bars on the phylogenetic tree
879 represent the number of amino acid substitutions per position in the underlying MUSCLE
880 alignment. Ultrafast bootstrap values calculated by IQtree at key nodes are shown.
881
882
883
884
885
886

887 **Supplementary Data**

888 **Supp. Data Fasta 1: CD-NTase**

889 Fasta file with all CD-NTase amino acid sequences analyzed.

890 **Supp. Data Fasta 2: STING**

891 Fasta file with all STING amino acid sequences analyzed.

892 **Supp. Data Fasta 3: Viperin**

893 Fasta file with all Viperin amino acid sequences analyzed.

894 **Supp. Data Table 1: Hmmscan excel file**

895 Hmmscan data for each CD-NTase, STING, and Viperin protein sequence.

896 **Supp. Data Tree 1: CD-NTase**

897 Newick file of maximum likelihood phylogenetic tree generated from a MUSCLE
898 alignment with IQtree. Newick file is used in Fig. 2 and Supp. Fig. 3 and 4. Node support
899 values calculated from ultrafast bootstraps.

900 **Supp. Data Tree 2: STING**

901 Newick file of maximum likelihood phylogenetic tree generated from a MUSCLE
902 alignment with IQtree. Newick file is used in Fig. 3 and Supp. Fig. 3 and 5. Node support values
903 calculated from ultrafast bootstraps.

904 **Supp. Data Tree 3: Viperin**

905 Newick file of maximum likelihood phylogenetic tree generated from a MUSCLE
906 alignment with IQtree. Newick file is used in Fig. 4 and Supp. Fig. 3 and 6. Node support values
907 calculated from ultrafast bootstraps.