1    **Heliorhodopsin evolution is driven by photosensory promiscuity in monoderms**

2    Paul-Adrian Bulzu[1], Vinicius Silva Kavagutti[1,2], Maria-Cecilia Chiriac[1], Charlotte
3    D. Vavourakis[3], Keiichi Inoue[4], Hideki Kandori[5,6], Adrian-Stefan Andrei[7], Rohit
4    Ghai[1]*

5    [1]Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre of
6    the Academy of Sciences of the Czech Republic, České Budějovice, Czech Republic.

7    [2]Department of Ecosystem Biology, Faculty of Science, University of South Bohemia,
8    Branišovská 1760, České Budějovice, Czech Republic.

9    [3]EUTOPS, Research Institute for Biomedical Aging Research, University of Innsbruck,
10   Austria.

11   [4]The Institute for Solid State Physics, The University of Tokyo, Kashiwa, Japan.

12   [5]Department of Life Science and Applied Chemistry, Nagoya Institute of Technology,
13   Showa, Nagoya 466-8555, Japan.

14   [6]OptoBioTechnology Research Center, Nagoya Institute of Technology, Showa, Nagoya
15   466-8555, Japan

16   [7]Limnological Station, Department of Plant and Microbial Biology, University of Zurich,
17   Kilchberg, Switzerland.


18   *Corresponding author: Rohit Ghai
19   Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre of
20   the Academy
21   of Sciences of the Czech Republic, Na Sádkách 7, 370 05, České Budějovice, Czech
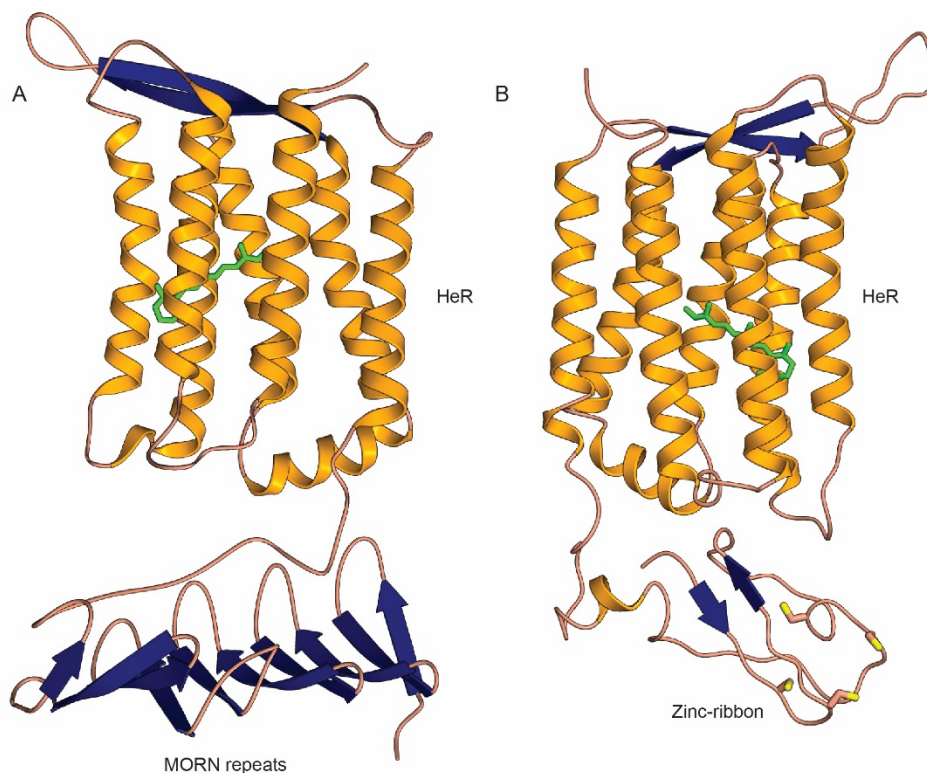22   Republic.
23   Phone: +420 387 775 881
24   Fax: +420 385 310 248
25   E-mail: ghai.rohit@gmail.com
26

27  The ability to harness Sun's electromagnetic radiation by channeling it into high-energy
28  phosphate bonds empowered microorganisms to tap into a cheap and inexhaustible
29  source of energy. Life`s billion-years history of metabolic innovations led to the
30  emergence of only two biological complexes capable of harvesting light: one based on
31  rhodopsins and the other on (bacterio)chlorophyll. Rhodopsins encompass the most
32  diverse and abundant photoactive proteins on Earth and were until recently canonically
33  split between type-1 (microbial rhodopsins) and type-2 (animal rhodopsins) families.
34  Unexpectedly, the long-lived type-1/type-2 dichotomy was recently amended through the
35  discovery of heliorhodopsins (HeRs) (Pushkarev et al. 2018), a novel and exotic family of
36  rhodopsins (i.e. type-3) that evaded recognition in our current homology-driven scrutiny
37  of life's genomic milieu. Here, we bring to resolution the debated monoderm/diderm
38  occurrence patterns by conclusively showing that HeR distribution is restricted to
39  monoderms. Furthermore, through investigating protein domain fusions, contextual
40  genomic information, and gene co-expression data we show that HeRs likely function as
41  generalised light-dependent switches involved in the mitigation of light-induced oxidative
42  stress and metabolic circuitry regulation. We reason that HeR's ability to function as
43  sensory rhodopsins is corroborated by their photocycle dynamics (Pushkarev et al. 2018)
44  and that their presence and function in monoderms is likely connected to the increased
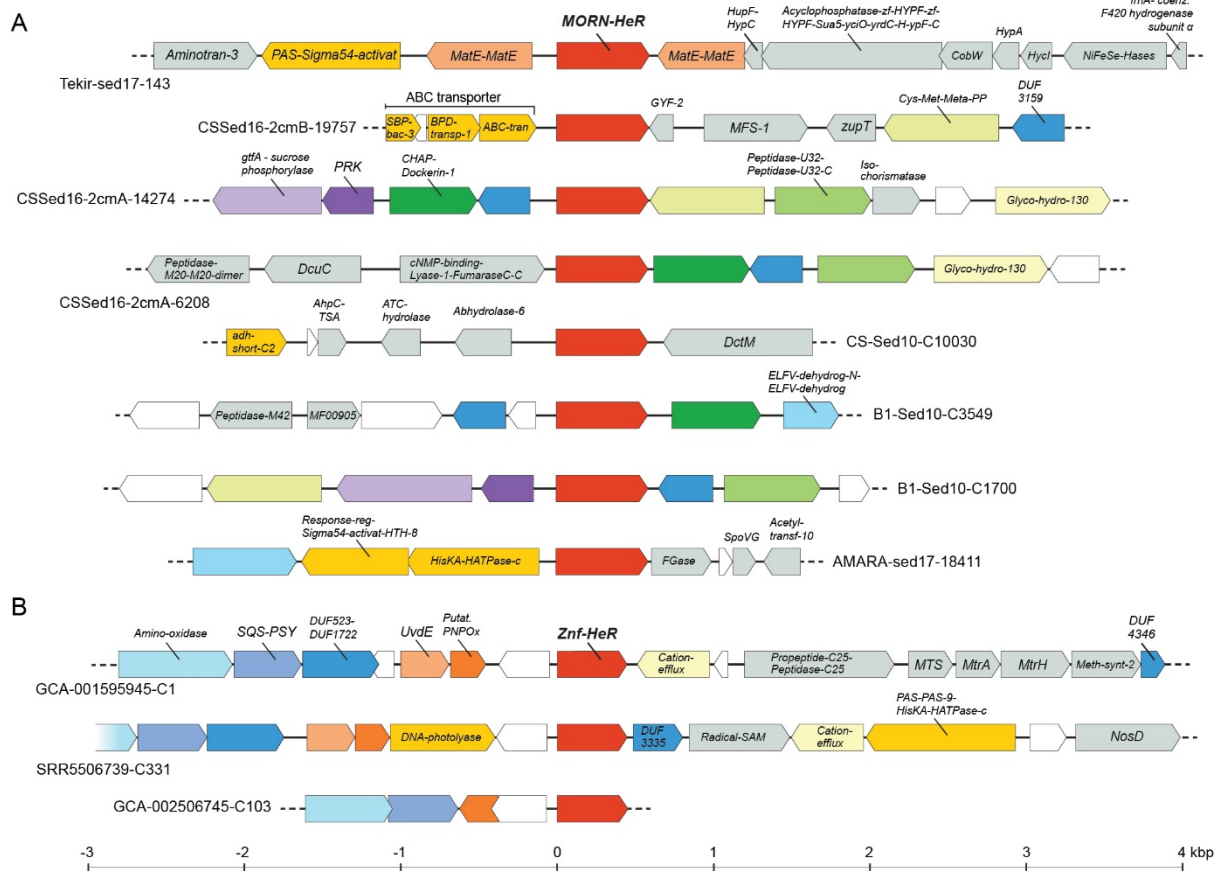45  sensitivity to light-induced damage of these organisms (Maclean et al. 2009).

46  Type-1 and -2 rhodopsins families share a similar topological conformation and little or no
47  sequence similarity among each other. Despite dissimilarities in function, structure and
48  phylogeny, type-1 and -2 rhodopsins have a similar membrane orientation with their N-
49  terminus being situated in the extracellular space. Identified during a functional
50  metagenomics screen and characterised by low sequence similarity when compared to
51  type-1 rhodopsins, HeRs attracted increasing research interest due to their peculiar
52  membrane orientation (i.e. N-terminus in the cytoplasm and the C-terminus in the
53  extracellular space)(Pushkarev et al. 2018), unusual protein structure (Kovalev et al. 2020)
54  and controversial taxonomic distribution (Flores-Uribe, Hevroni, and Ghai 2019). While
55  electrophysiological (Pushkarev et al. 2018), physicochemical (Tanaka et al. 2020) and
56  structural (Shihoya et al. 2019; Kovalev et al. 2020) studies achieved great progress in
57  elucidating a series of characteristics ranging from photocycle length (indicating no
58  pumping activity) to detailed protein organization, they provide no data regarding the
59  biological function of HeRs. Moreover, polarized opinions regarding the putative
60  ecological role and taxonomic distribution of HeR-encoding organisms (Flores-Uribe,
61  Hevroni, and Ghai 2019; Kovalev et al. 2020) call for the use of novel approaches in
62  establishing HeR functionality. This work draws its essence from the tenet that functionally
63  linked genes within prokaryotes are co-regulated, and thus occur close to each other
64  (Aravind 2000; Huynen et al. 2000). Within this framework, the functions of
65  uncharacterised genes (i.e. HeRs) can be inferred from their genomic surroundings. Here
66  we couple HeR's distributional patterns with contextual genomic information involving
67  protein domain fusions and operon organization, and gene expression data to shed light
68  on HeRs functionality.

69
70 **Figure 1.** Modelled three-dimensional (3D) structures of MORN-HeR and Znf-HeR protein
71 domain fusions. (A) 3D model of a heliorhdodopsin (HeR) containing three N-terminal
72 MORN domain repeats. (B) 3D model of a HeR containing an N-terminal Zn ribbon motif.
73 Both models are oriented with the extracellular side up and intracellular side down.
74 Retinal is coloured green and cysteine residues are depicted with yellow-topped orange
75 sticks.

76 Previous assessments of taxonomic distribution of HeRs reported conflicting data
77 regarding their presence in monoderm (Flores-Uribe, Hevroni, and Ghai 2019) and
78 diderm (Kovalev et al. 2020) prokaryotes. In order to accurately map HeR taxonomic
79 distribution we used the GTDB database (release 89), since it contains a wide-range of
80 high-quality genomes derived from isolated strains and environmental metagenome-
81 assembled genomes, classified within a robust phylogenomic framework (Parks et al.
82 2020). By scanning 24,706 genomes, we identified 450 *bona fide* HeR sequences
83 (topology: C-terminal inside and N-terminal outside, seven transmembrane helices and a
84 SxxxK motif in helix 7; Supplementary Table S1) spanned across 17 phyla (out of 151;
85 Supplementary Table S2). In order to assign HeR-containing genomes to either
86 monoderm or diderm categories, we employed a set of 27 manually curated protein
87 domain markers that are expected to be restricted to organisms possessing double-
88 membrane cellular envelopes (i.e. diderms) (Taib et al. 2020). While most analyses were
89 expected to be influenced by varying levels of genome completeness , we found that a
90 conservative criterion of presence of at least ten marker domains singled out all diderm
91 lineages (i.e. Negativicutes, Halanaerobiales and Limnochondria) (Taib et al. 2020;
92 Megrian et al. 2020) within the larger monoderm phylum Firmicutes, apart from correctly
93 identifying other well-known diderms. Except for three genomes (one each belonging to

94  Myxococcota, Spirochaetota and Dictyoglomota phyla), all other HeR occurrences were
95  restricted to monoderms (Supplementary Table S2). Examination of the HeR-encoding
96  Myxoccoccota contig by querying its predicted proteins against the RefSeq and GTDB
97  databases revealed it to be an actinobacterial contaminant. The *Spirochaeta* genome was
98  incomplete (60% estimated completeness) and only encoded for two outer membrane
99  marker genes, making any inferences regarding its affiliation to monoderm or diderm
100 bacteria impossible. However, we could not rule out that this genome could belong to a
101 member lacking lipopolysaccharides (LPS) (Taib et al. 2020). The Dictyoglomota genome
102 belongs to an isolate, and despite its high completeness, it encodes only five markers.
103 Combined with the notion that Dictyoglomota are known to have atypical membrane
104 architectures (Saiki et al. 1985), the presence of only five markers points towards the
105 absence of a classical diderm cell envelope. Apart from these exceptions, all other HeR-
106 encoding genomes are monoderm and, at least within this collection, we find no strong
107 evidence of HeRs being present in any organism that is conclusively diderm. We also
108 identified HeRs in several assembled metagenomes and metatranscriptomes (see
109 Methods). For improved resolution of taxonomic origin, we considered only contigs of at
110 least 5 Kb in length (n = 1,340 from metagenomes and n = 4 from metatranscriptomes).
111 Following a strict approach to taxonomy assignment (i.e. at least 60 % genes giving best-
112 hits to the same phylum and not just majority-rule), we could designate a phylum for most
113 HeRs. Without any exception, we found that all the contigs that received robust taxonomic
114 classification (n = 1,319) belonged to known monoderm phyla (Supplementary Table S3).
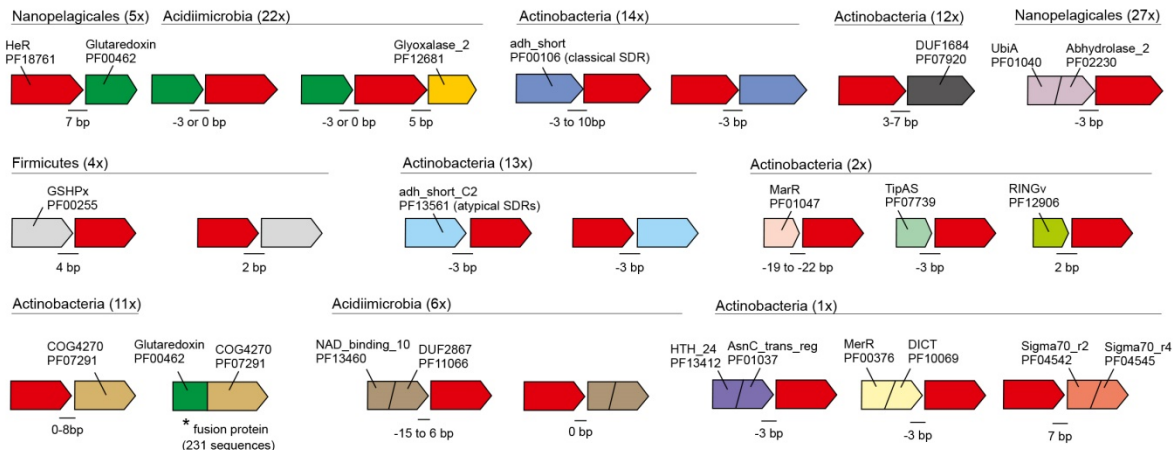


115

116 **Figure 2.** Genomic context of HeR-protein domain fusion genes. A) Representative
117 MORN-HeR encoding contigs identified in strictly anaerobic Firmicutes. B) Contigs
118 encoding Znf-HeR domain fusions. Neighbouring genes were depicted within an interval

119    spanning ~ 7 kb, centered on HeR. Genes occurring only once within the considered
120    intervals are coloured grey; genes encoding HisKA, PAS, regulatory domains, as well as
121    other discussed HeR neighbours are depicted bright yellow. Homologous genes
122    occurring multiple times found within each category of HeR-protein fusion contigs are
123    depicted using matching colours. Hypothetical genes are white.

124    Domain fusions with rhodopsins are recently providing novel insights into the diverse
125    functional couplings that enhance the utility of a light sensor, e.g. the case of a
126    phosphodiesterase domain fused with a type-1 rhodopsin (Ikuta et al. 2020). As far as we
127    are aware, no domain fusions have been described for HeRs yet. In our search for such
128    domain fusions that may shed light on HeR functionality, the MORN repeat (Membrane
129    Occupation and Recognition Nexus, PF02493) was found in multiple copies (typically 3) at
130    the cytoplasmic N-terminus of some HeRs (n = 36). A tentative 3D model for a
131    representative MORN-HeR could be generated and is shown in Figure 1A. These MORN-
132    HeR sequences were phylogenetically restricted to two environmental branches of MAGs
133    recovered from haloalkaline sediments that affiliate to the family *Syntrophomonadaceae*
134    (phylum Firmicutes) (Timmers et al. 2018; Vavourakis et al. 2018, 2019) (Supplementary
135    Figure 1). The prototypic MORN repeat, consisting of 14 amino acids with the consensus
136    sequence YEGEWxNGKxHGYG, was first described in 2000 (Takeshima et al. 2000) from
137    junctophilins present in skeletal muscle and later recognized to be ubiquitous in both
138    eukaryotes and prokaryotes (El-Gebali et al. 2019). This conserved signature can be seen
139    in the alignment of MORN-repeats fused to HeRs (Supplementary Figure 2). MORN-
140    repeats have been shown to bind to phospholipids (Im et al. 2007; Ma et al. 2006),
141    promoting stable interactions with plasma membranes (Takeshima et al. 2000) and also
142    function as protein-protein interaction modules involved in di- and oligomerization (Sajko
143    et al. 2020). They are expected to be intracellular and provide a large putative interaction
144    surface (either with other MORN-HeRs or other proteins). A widespread adaptation of
145    bacteria to alkaline environmental conditions is the increased fluidity of their plasma
146    membranes achieved by the incorporation of branched-chain and unsaturated fatty acids
147    which ultimately influences the configuration and activity of membrane integral proteins
148    such as ATP synthases and various transporters (Kanno et al. 2015). Microbial rhodopsins
149    typically associate as oligomers in vivo, which is also the case with heliorhodopsins that
150    are known to form dimers (Shibata et al. 2018; Shihoya et al. 2019). The presence of
151    MORN-repeats in HeRs exclusively within extreme haloalkaliphilic bacteria (class
152    Dethiobacteria) may be accounted for via their potential role in stabilizing HeR dimers in
153    conditions of increased membrane fluidity (Supplementary Figure 4). Another possibility
154    would be the interaction of MORN-repeats with other MORN-repeat containing proteins
155    encoded in these MAGs. We could indeed identify multiple MORN-protein domain
156    fusions co-occurring in genomes of analysed Dethiobacteria (Supplementary Figures 1
157    and 3; Supplementary Table S15). Even though the nature of interactions amongst these
158    proteins with intracellular MORN-repeats is unclear, they raise the possibility that MORN-
159    repeats act as downstream transducers of conformational changes occuring in HeRs. Such
160    tandem repeat structures may function as versatile target recognition sites capable of
161    binding not only small molecules like nucleotides but also peptides and larger proteins
162    (Kajava 2012). If true, this would render HeRs as sensory rhodopsins. In support of this, we
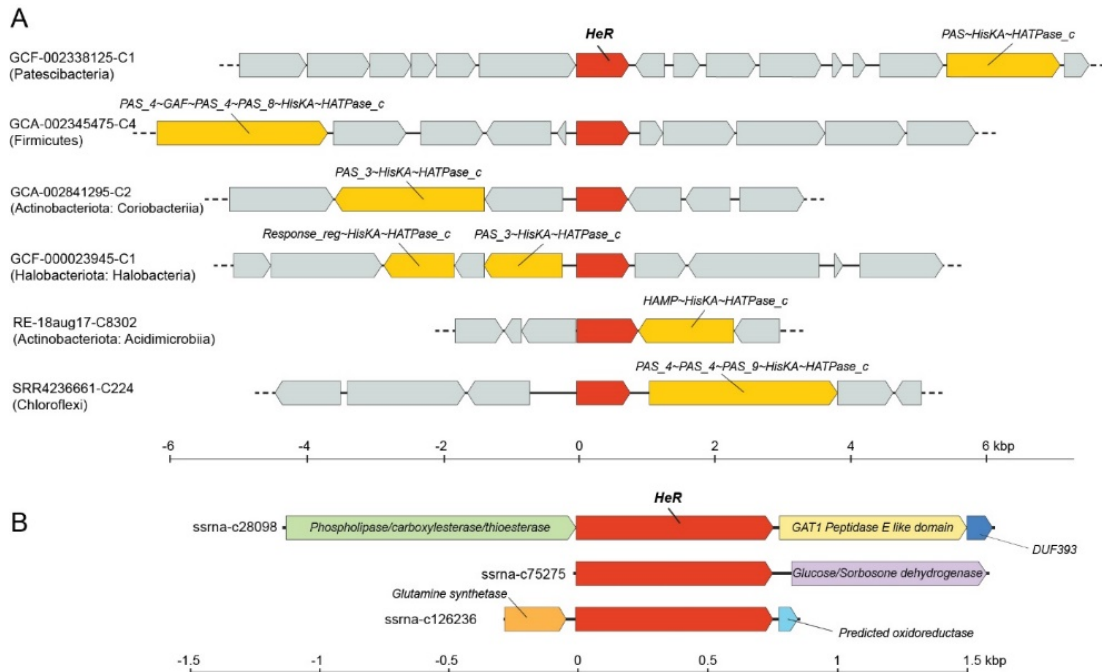163    found several genes in close proximity to MORN-HeRs encoding signature protein

164   domains (e.g. PAS, HisKA, HATPase_c) that are known to be involved in histidine kinase
165   signalling (Aravind, Iyer, and Anantharaman 2010) (Figure 2A).



166
167   **Figure 3**. Schematic representation of genes that may be transcriptionally linked to HeRs.
168   Taxonomic categories and number of occurrences are shown at the top of each putative
169   operon. Intergenic distances (in bp) are indicated at gene junctions. Negative distance
170   values indicate overlapping genes. Pfam or COG identifiers are used to represent domain
171   architectures. A star (*) indicates a fused gene (two domains: Glutaredoxin and COG4270)
172   found in at least 473 genomes from GTDB and 231 unique sequences in UniProt
173   suggesting a functional linkage of COG4270 with Glutaredoxin.

174   As no other obvious domains were found to be fused with HeRs using standard profile
175   searches, we examined all N and C-terminal extensions as well as loops longer than 50 aa
176   by performing more sensitive profile-profile searches using HHPred (Zimmermann et al.
177   2018). We found at least ten N-terminal extensions of HeRs (ntv1-ntv10), 22 variants of
178   ECL1 (extracellular loop 1), a single type of loop extension for ICL2 (intracellular loop 2)
179   and three variants of ICL3 (intracellular loop 3). A complete listing of all alignments and
180   summary results of HHpred can be found in Supplementary Table S8. Remarkably, we
181   found significant matches in a set of six sequences (all originating from
182   Thermoplasmatales archaea) to zinc ribbon proteins (Pfam domain zinc_ribbon_4) at the
183   N-terminus of some heliorhodopsins (these extensions are termed N-terminal variant 1 or
184   ntv1, Supplementary Table S8). Zinc ribbons belong to the larger family of Zinc-finger
185   domains (Krishna et al., 2003). A CxxC-17x-CxxC was found in this region that likely
186   coordinates a metal (e.g., zinc or iron). These CxxC_CxC type motifs are common to a
187   wider family of Zinc finger-like proteins that were initially found to bind to DNA and later
188   shown to be capable of binding to RNAs, proteins and small molecules (Krishna,
189   Majumdar, and Grishin 2003). Similar motifs are also seen in Rubredoxins and
190   Cys_rich_KTR domains. We term these fused ntv1 protein variants as Znf-HeRs (Zinc finger
191   Heliorhodopins). A modelled structure for a representative Znf-HeR is shown in Figure 1.
192   In one contig encoding a Znf-HeR we identified a histidine kinase that could be
193   functionally linked (Figure 2B). Notably, most identified Znf-HeRs are flanked by genes
194   known to be triggered by light exposure and play key roles in photoprotection (i.e.
195   carotenoid biosynthesis genes Lycopene cyclase, phytoene desaturase – Amino-oxidase,
196   squalene/phytoene synthase – SQS-PSY) and UV-induced DNA damage repair (DNA
197   photolyases, UV-DNA damage endonucleases – UvdE) (Rastogi et al. 2010; Yatsunami et
198   al. 2014). Recent research showed that HeRs from Thermoplasmatales archaea (TaHeR)

199 and uncultured freshwater Actinobacteria (48C12) (for which the structure is resolved and
200 lacks the ntv1 extension) might bind zinc (Hashimoto et al. 2020). As the zinc binding site
201 could not be precisely identified it was suggested that it could be located in the
202 cytoplasmic part, and responsible for modifying the function of HeR. Our discovery of Znf-
203 HeRs offers additional, more direct indications of the role of zinc in the possible
204 downstream signalling by HeRs.



205

206 **Figure 4.** A) Genes encoding HisKA domain signalling proteins identified in the proximity
207 of HeR genes from diverse phyla. All genes containing HisKA domains are coloured bright
208 yellow, HeRs are shown in red, and all other genes in grey. B) Transcripts obtained by
209 strand-specific metatranscriptomics from freshwater encoding genes co-expressed with
210 HeR.

211 Given the large number of long contigs encoding HeRs (from genomes and
212 metagenomes), we sought to identify candidate genes that could be transcribed together
213 with HeRs (in the same operon). We used the following strict criteria for obtaining such
214 genes 1) the intergenic distance between such a gene and the HeR must be less than 10
215 bp, and 2) the gene must be located on the same strand. A number of interesting
216 candidates emerged in this analysis with the most frequent ones being summarized in
217 Figure 3 (a complete table can be found in Supplementary Table S9). We identified
218 multiple instances in which genes with Glutaredoxin and GSHPx PFAM domains were
219 found adjacent to HeRs (n = 31). Glutaredoxins are small redox proteins with active
220 disulphide bonds that utilize reduced glutathione as an electron donor to catalyze thiol-
221 disulphide exchange reactions. They are involved in a wide variety of critical cellular
222 processes like the maintenance of cellular redox state, iron and redox-sensing, and
223 biosynthesis of iron-sulphur clusters (Lillig, Berndt, and Holmgren 2008; Rouhier et al.
224 2010). Glutathione is also used by glutathione peroxidase (GSHPx) to reduce hydrogen
225 peroxide and peroxide radicals i.e. as an anti-oxidative stress protection system (Bhabak
226 and Mugesh 2010). Additionally, there are also instances where Glutaredoxin and genes
227 containing Glyoxalase_2 domains may be co-transcribed with HeRs. Glyoxalases, in

228    concert with glutaredoxins, are critical for detoxification of methylglyoxal, a toxic
229    byproduct of glycolysis (Ferguson et al. 1998). Moreover, adjacent to HeRs we find at least
230    three instances where a catalase gene is also present (in Actinobacteria; see
231    Supplementary Figures S10-S11). Collectively, these observations suggest a role for HeRs
232    in oxidative stress mitigation. In one case, we found a gene encoding the DICT domain
233    (Figure 3) which is frequently associated to GGDEF, EAL, HD-GYP, STAS, and two-
234    component system histidine kinases. Notably, it has been predicted to have a role in light
235    response (Aravind, Iyer, and Anantharaman 2010).

236    Although we assembled contigs encoding HeRs from previously published
237    metatranscriptomes, the lack of strand-specific transcriptomes hampered any clear
238    conclusions on whether or not genes adjacent to HeRs are indeed co-transcribed, leaving
239    open the possibility that they might simply be artefacts of assembly (Zhao et al. 2015). In
240    order to gather more definitive evidence for co-transcription we performed strand-
241    specific metatranscriptome sequencing for a freshwater sample (see Methods). We
242    recovered six HeR-encoding transcripts that were > 1 kb in length. All these transcripts are
243    predicted to originate from highly abundant freshwater Actinobacteria with streamlined
244    genomes (four transcripts from "*Ca. Planktophila*" and two from "*Ca. Nanopelagicus*")
245    (Supplementary Table S12) (Neuenschwander et al. 2018). Overall, there are three types
246    of transcripts based upon gene content: class1 - encoding Glutamine synthetase catalytic
247    subunit and NAD+ synthetase; class2 - encoding a hydrolase, a peptidase and a DUF393
248    domain containing protein, and class3 - encoding glucose/sorbosone dehydrogenase
249    (GSDH) (Figure 4B and Supplementary Table S12). A common theme for glutamine
250    synthetase and NAD+ synthetase is that both utilize ammonia and ATP to produce
251    glutamine and NAD+ respectively. Moreover, some NAD+ synthetases may be glutamine
252    dependent (Resto, Yaffe, and Gerratana 2009). Glutamine synthetase in particular is a key
253    enzyme for nitrogen metabolism in prokaryotes at large (García-Domínguez, Reyes, and
254    Florencio 1999). For hydrolases and peptidases, the function prediction is somewhat
255    broad. Glucose/sorbosone dehydrogenase catalyses the production of gluconolactone
256    from glucose (Oubrie et al. 1999). Therefore, it appears that all six HeRs are generally co-
257    transcribed with genes involved in nitrogen assimilation and degradation/assimilation of
258    sugars and peptides. This would suggest that these processes are also influenced by light,
259    with such a link between light-dependent increase in sugar uptake and metabolic activity
260    being recently proposed in non-phototrophic Actinobacteria (Maresca et al. 2019). Light
261    also triggers photosynthetic activity, increasing availability of sugars and other nutrients
262    (e.g. glutamine and ammonia) for heterotrophs. In this vein, a link between a light sensing
263    mechanism, e.g. via heliorhodopsins, may lead to elevated metabolic activity.
264    In a previous study, histidine kinases were deemed absent in the vicinity of HeRs (Kovalev
265    et al. 2020). Given that our initial analyses predicted a sensory function, we examined
266    genomic regions spanning 10 kb up- and downstream of HeRs. Already in the case of
267    MORN-HeRs and Znf-HeRs we observed histidine kinase signalling components in close
268    proximity to them (Figure 2). In our search we detected multiple instances of histidine
269    kinases (HisKA) fused with PAS, GAF, MCP_Signal, HAMP or HATPase_c domains in the
270    gene neighbourhoods of HeRs in distinct phyla (e.g. Actinobacteria, Chloroflexi,
271    Patescibacteria, Firmicutes, Dictyoglomota, Thermoplasmatota) (Figure 4B; more details
272    in Supplementary Figures S5-S14). Moreover, in many cases multiple response regulator
273    genes were present in the same regions (Pfam domains Response_reg, Trans_reg_C).

274     Less frequently, GGDEF and EAL domains, usually associated with bacterial signalling
275     proteins, were also present. Using overrepresentation analysis (Shmakov et al. 2018), we
276     found that the occurrence of two-component system protein domains in the vicinity of
277     HeRs is statistically significant (see Methods and Supplementary Table S11). In addition to
278     these two-component system proteins, the same regions also appear enriched in redox
279     proteins (e.g. thioredoxin, peroxidase, catalase). The close association of two-component
280     systems, genes involved in oxidative stress mitigation and HeRs points towards a
281     functional interaction.

282     **Conclusions**

283     In conclusion, contextual genomic information shows that monoderm prokaryotes use
284     HeRs in multiple mechanisms for the activation of downstream metabolic pathways post
285     light sensing. Furthermore, we offer tantalizing clues regarding the involvement of HeRs
286     in multiple cellular processes and add new lines of inquiry for the primary role of HeRs in
287     light signal transduction. Additional support for the role of HeRs in light sensing is
288     inferred from the frequent association of HeRs with classical histidine kinases and
289     associated protein domains in multiple phyla. Furthermore, multiple types of N-terminal
290     domain fusions found in specific subfamilies of HeRs (i.e. MORN domains in
291     haloalkaliphilic Firmicutes and Zinc ribbon type domains in Thermoplasmatales archaea)
292     point to possible downstream signalling which may be effected by recruitment of
293     additional, as yet unknown, partner proteins.

294     We further propose a critical role for HeRs in protecting monoderm cells from light-
295     induced oxidative damage. In this sense, we observed a close association and probable
296     transcriptional linkage of HeRs to glyoxylases and glutaredoxins (sometimes seen as
297     overlapping genes). Given that light can induce the uptake and metabolism of sugars, as
298     previously discussed for certain Actinobacteria (Maresca et al. 2019), it is expected that
299     increased sugar availability resulting from photosynthesis leads to increased glycolytic
300     activity in heterotrophic bacteria. Glycolysis also produces small amounts of toxic
301     methylglyoxal that can be neutralized by the combined action of glyoxylases and
302     glutaredoxins. In this sense, it appears that at least in some Actinobacteria, glyoxylases
303     and glutaredoxins may be transcribed together with HeRs, but how the transcription is
304     controlled remains unclear. Additional evidence of transcriptional linkages of HeRs to
305     proteins like peroxiredoxin and catalase also imply a light-dependent activation, boosting
306     the cellular response to light induced oxidative damage which may be critical for both
307     aerobes and anaerobes. Evidence from strand-specific HeR transcripts originating from
308     freshwater Actinobacteria suggests the further involvement of HeRs in nitrogen and sugar
309     metabolism via glutamate synthase, NAD+ synthases and glucose/sorbosone
310     dehydrogenases in these organisms.

311     Overall, the picture that emerges (at least for some organisms) is one of HeR's role in
312     responding to light and transmitting the signal via histidine kinases. Downstream
313     processes that are ultimately regulated are diverse, including possible roles for HeRs in
314     the mitigation of light-induced oxidative damage and in the regulation of nitrogen
315     assimilation and carbohydrate metabolism, processes that may benefit from a light-
316     dependent activation through more efficient utilization of available resources.

317 Recent work has shown more support for the diderm-first ancestor (Coleman et al. 2020)
318 and given the far broader distribution of type-1 rhodopsins in both mono- and di-derm
319 organisms it appears likely that type-1 rhodopsins emerged prior to HeRs. The very
320 restricted distribution of HeRs to monoderms would support this view as well. Even so,
321 HeRs are not universally present in monoderms and when present, appear to be
322 associated with diverse genes involved in signal transduction, oxidative stress mitigation,
323 nitrogen and glucose metabolism. This would suggest they have been exapted as
324 generalized sensory switches that may allow light-dependent control of metabolic activity
325 in multiple lineages, somewhat similar to type-1 rhodopsins where minor modifications
326 have led to emergence of a wide variety of ion-pumps (Kandori 2020). The frequent
327 distribution of HeRs in aquatic environments (habitats characterised by increased light
328 penetration), where they commonly occur within phylum Actinobacteriota, helps us to
329 explain their monoderm-restricted presence. Abundant freshwater actinobacterial
330 lineages are generally typified by lower GC content (Ghai, McMahon, and Rodriguez-
331 Valera 2012) and increased vulnerability to oxidative stress damage (Kim et al. 2019). This
332 susceptibility is also illustrated by actinobacterial phages that exhibit positive selection
333 towards reactive oxygen species defense mechanisms (Kavagutti et al. 2019). Given the
334 fact that monoderms are generally more sensitive to light-induced damage and
335 corroborated with up-mentioned metabolic implications, we consider that HeRs evolved
336 as sensory switches capable of triggering a fast response against photo-oxidative stress in
337 prokaryotic lineages more sensitive to light.

338 **Methods**

339 **Metagenomes and metatranscriptomes.** We used previously published metagenomics
340 and metatranscriptomics data from freshwaters (Andrei et al. 2019; Kavagutti et al. 2019;
341 Mehrshad et al. 2018), haloalkaline brine and sediment (Vavourakis et al. 2018, 2019),
342 brackish sediments (Bulzu et al. 2019), GEOTRACES cruise (Biller et al. 2018) and TARA
343 expeditions (Salazar et al. 2019). In addition, we downloaded multiple environmental
344 metagenomes (sludge, marine, pond, estuary, etc.) from EBI MGnify
345 (https://www.ebi.ac.uk/metagenomics/) (Mitchell et al. 2020) and assembled them using
346 Megahit v1.2.9 (D. Li et al. 2016). All contigs in this work are named or retain existing
347 names that allow tracing them to their original datasets.

348 **Sequence search for *bona fide* rhodopsins**. Genes were predicted in metagenomics
349 contigs using Prodigal (Hyatt et al. 2010). Candidate rhodopsin sequences were scanned
350 with hmmsearch (Eddy 2011) using PFAM models (PF18761: heliorhodopsin, PF01036:
351 bac_rhodopsin) and only hits with significant e-values ( < 1e-3) were retained. Homologs
352 for these sequences were identified by comparison to a known set of rhodopsin
353 sequences (Bulzu et al. 2019) using MMSeqs2 (Hauser, Steinegger, and Söding 2016) and
354 alignments were made using MAFFT-linsi(Katoh and Standley 2013). These alignments
355 were used as input to Polyphobius (Käll, Krogh, and Sonnhammer 2005) for
356 transmembrane helix prediction. Only those sequences that had seven transmembrane
357 helices and either a SxxxK motif (for heliorhodopsins) or DxxxK motif (for
358 proteorhodopsins) in TM7 were retained. In addition, we also screened the entire
359 UniProtKB for HeRs. In total, we accumulated at least 4,108 (3,606+502) *bona fide* HeR
360 sequences.

**Taxonomic classification of assembled contigs**. Contigs were dereplicated using cd-hit (W. Li and Godzik 2006) (95% sequence identity and 95% coverage). Only contigs ≥ 10 kb were retained for this analysis. A custom protein database was created by predicting and translating genes in all GTDB genomes (release 89) (Parks et al. 2020) using Prodigal (Hyatt et al. 2010). These sequences were supplemented with viral and eukaryotic proteins from UniProtKB (UniProt Consortium 2019). Best-hits against predicted proteins in contigs were obtained using MMSeqs2 (Hauser, Steinegger, and Söding 2016). Taxonomy was assigned to a contig (minimum length 5 kb) only if ≥ 60% of genes in the contig gave best-hits to the same phylum. All contigs that appeared to originate from diderms were cross-checked against NCBI RefSeq (accessed online on 15[th] December 2020).

**Outer-envelope detection.** A set of protein domains found in genes encoding for the outer-envelope (Taib et al. 2020) was further reduced to include only those domains that were found mostly in known diderms. These domains were searched against the predicted proteins in all genomes in GTDB using hmmsearch (e-value < 1e-3). The results are shown in Supplementary Table S13.

**Protein function/structure predictions.** Predicted proteins were annotated using TIGRFAMs (Haft, Selengut, and White 2003) and COGs (Galperin et al. 2015). Domain predictions were carried out using the pfam_scan.pl script against the PFAM database (release 32) (El-Gebali et al. 2019). Profile-profile searches were carried out online using the HHPred server (Zimmermann et al. 2018). Additional annotations were added using Interproscan (Mitchell et al. 2019). Protein structure predictions were carried out using the Phyre2 server (Kelley et al. 2015) and structures were visualized with CueMol (http://www.cuemol.org/en/).

**Domains overrepresentation near heliorhodopsin.** A subset of high-quality MAGs (n = 240) containing HeR-encoding genes flanked both up- and downstream by a minimum of 10 genes were selected from GTDB (release 89) (Parks et al. 2020). For each genome, the probability of finding any particular domain by chance in a random subset of 20 genes was calculated using the hypergeometric distribution (without replacement) in R with the function *phyper* (*stats* package) (Johnson, Kemp, and Kotz 2005). In order to account for type I errors arising from multiple comparisons, hypergeometric test P-values were adjusted using the Benjamini-Hochberg procedure (Benjamini and Hochberg 1995). Further, we selected domains located in the proximity of HeR in at least 10% of genomes with low probability (FDR corrected P-value < 0.05). This procedure that was initially employed for the whole GTDB genome collection was repeated for individual phyla containing HeR-encoding genes within at least five genomes.

**Strand-specific freshwater transcriptome sequencing and assembly.** Sampling was performed on the 16[th] of August 2020 at 9:00 in Řimov reservoir, Czech Republic, (48°50'54.4"N, 14°29'16.7"E) using a hand-held vertical Friedinger (2 L) sampler. A total of 20 L of water were collected from a depth of 0.5 m and immediately transported to the laboratory. Serial filtration was carried out by passing water sample through a 20 μm pore size pre-filter mesh followed by a 5 μm pore size PES filter (Sterlitech) and a 0.22 μm pore size PES filter (Sterlitech, USA) using a Masterflex peristaltic pump (Cole-Palmer, USA).

404 Filtration was done at maximum speed for 15 minutes to limit cell lysis and RNA damage
405 as much as possible. A total volume of 3.7 L was filtered during this time. PES filters (5 µm
406 and 0.22 µm pore sizes) were loaded into cryo-vials pre-filled with 500 µl of DNA/RNA
407 Shield (Zymo Research, USA) and stored at -80°C. RNA was extracted from filters using the
408 Direct-zol RNA MicroPrep (Zymo Research, USA) after they had been previously thawed,
409 partitioned, and subjected to mechanical lysis by bead-beating in ZR BashingBead™ Lysis
410 tubes (with 0.1 and 0.5 mm spheres). DNase treatment was performed to remove
411 genomic DNA during RNA extraction as an "in-column" step described in the Direct-zol
412 protocol and was repeated after RNA elution, by using the Ambion Turbo DNA-freeTM Kit
413 (Life Technologies, USA). RNA was quantified using a NanoDrop® ND-1000 UV-Vis
414 spectrophotometer (Thermo Fisher Scientific, USA) and integrity was verified by agarose
415 gel (1%) electrophoresis. A total of 4.6 µg of RNA from the 0.22 µm pore size filter and 2.6
416 µg from the 5 µm pore size filter were sent for dUTP-marking based strand-specific
417 metatranscriptomic sequencing at Novogene (www.novogene.com). Following quality
418 control at Novogene, samples were mixed into one single reaction, then subjected to
419 rRNA depletion and used for stranded library preparation. Strand-specificity was achieved
420 by incorporation of dUTPs instead of dTTPs in the second-strand cDNA followed by
421 digestion of dUTPs by uracil-DNA glycosylase to prevent PCR amplification of this strand.
422 Paired-end (PE 150 bp) sequencing was carried out using a Novaseq 6000 platform. A
423 total of 166,213,184 raw sequencing reads, amounting to 24.9 Gb, were produced. *De*
424 *novo* assembly of metatranscriptomic data was performed using rnaSPAdes v.3.14.1
425 (Bushmanova et al. 2019) in reverse-forward strand-specific mode (--ss rf) with a custom k-
426 mers list 29, 39, 49, 59, 69, 79, 89, 99, 109, 119, 127. A total of 156,235 hard-filtered
427 transcripts of a minimum length of 1 kb were assembled. Protein coding sequences were
428 predicted *de novo* using Prodigal (Hyatt et al. 2010) in metagenomic mode (-p meta).
429 Protein domains were annotated by scanning with InterProScan(Mitchell et al. 2019) while
430 PFAM (Protein Families)(El-Gebali et al. 2019) domains were identified using the publicly
431 available perl script pfam_scan.pl (ftp://ftp.ebi.ac.uk/pub/databases/Pfam/Tools/).
432 Proteins were scanned locally using HMMER3 (Eddy 2011) against the COGs (Clusters of
433 Orthologous Groups) (Galperin et al. 2015) HMM database (e-value < 1e-5) and the
434 TIGRFAMs (TIGR Families) (Haft, Selengut, and White 2003) HMM collection with trusted
435 score cutoffs. BlastKOALA (Kanehisa, Sato, and Morishima 2016) was used to assign KO
436 identifiers (KO numbers). Annotations for representative transcripts encoding HeR are
437 summarised in Supplementary Table S12.

## Data availability

439 Sequence data generated in this study have been deposited in the European Nucleotide
440 Archive (ENA) at EMBL-EBI under project accession number PRJEB35770 (run
441 ERR5100021). The derived data that support the findings of this paper, including R code
442 used for statistical analyses, are available in FigShare
443 (https://figshare.com/s/7bb42426f2ad5e891fec). All other relevant data supporting the
444 findings of this study are available within the paper and its supplementary information
445 files.
446
447
448

## Acknowledgements

## Author contributions

R.G. and P.-A.B. designed the study. P.-A.B., A.-Ş.A. and R.G. wrote the manuscript. P.-A.B., R.G., V.S.K, M.-C.C., C.D.V and A.-Ş.A. analysed and interpreted the data. K.I. and H.K. performed rhodopsin structural analyses. All authors commented on and approved the manuscript.

## Competing interests

The authors declare no competing interests.

References

Andrei, Adrian-Ştefan, Michaela M. Salcher, Maliheh Mehrshad, Pavel Rychtecký, Petr Znachor, and Rohit Ghai. 2019. "Niche-Directed Evolution Modulates Genome Architecture in Freshwater Planctomycetes." *The ISME Journal* 13 (4): 1056–71.

Aravind, L. 2000. "Guilt by Association: Contextual Information in Genome Analysis." *Genome Research* 10 (8): 1074–77.

Aravind, L., L. M. Iyer, and V. Anantharaman. 2010. "Natural History of Sensor Domains in Bacterial Signaling Systems." In *Sensory Mechanisms in Bacteria: Molecular Aspects of Signal Recognition*, edited by Ray Dixon Stephen Spiro, 1–38. Caister Academic Press Norfolk, UK.

Benjamini, Y., and Y. Hochberg. 1995. "Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing." *Journal of the Royal Statistical Society*. https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/j.2517-6161.1995.tb02031.x.

Bhabak, Krishna P., and Govindasamy Mugesh. 2010. "Functional Mimics of Glutathione Peroxidase: Bioinspired Synthetic Antioxidants." *Accounts of Chemical Research* 43 (11): 1408–19.

Biller, Steven J., Paul M. Berube, Keven Dooley, Madeline Williams, Brandon M. Satinsky, Thomas Hackl, Shane L. Hogle, et al. 2018. "Marine Microbial Metagenomes Sampled across Space and Time." *Scientific Data* 5 (September): 180176.

Bulzu, Paul-Adrian, Adrian-Ştefan Andrei, Michaela M. Salcher, Maliheh Mehrshad, Keiichi Inoue, Hideki Kandori, Oded Beja, Rohit Ghai, and Horia L. Banciu. 2019. "Casting Light

489     on Asgardarchaeota Metabolism in a Sunlit Microoxic Niche." *Nature Microbiology* 4 (7):
490     1129–37.
491     Bushmanova, Elena, Dmitry Antipov, Alla Lapidus, and Andrey D. Prjibelski. 2019.
492     "RnaSPAdes: A de Novo Transcriptome Assembler and Its Application to RNA-Seq Data."
493     *GigaScience* 8 (9). https://doi.org/10.1093/gigascience/giz100.
494     Coleman, Gareth A., Adrián A. Davín, Tara Mahendrarajah, Anja Spang, Philip Hugenholtz,
495     Gergely J. Szöllősi, and Tom A. Williams. 2020. "A Rooted Phylogeny Resolves Early
496     Bacterial Evolution." *Cold Spring Harbor Laboratory*.
497     https://doi.org/10.1101/2020.07.15.205187.
498     Eddy, Sean R. 2011. "Accelerated Profile HMM Searches." *PLoS Computational Biology* 7
499     (10): e1002195.
500     El-Gebali, Sara, Jaina Mistry, Alex Bateman, Sean R. Eddy, Aurélien Luciani, Simon C.
501     Potter, Matloob Qureshi, et al. 2019. "The Pfam Protein Families Database in 2019."
502     *Nucleic Acids Research* 47 (D1): D427–32.
503     Ferguson, G. P., S. Tötemeyer, M. J. MacLean, and I. R. Booth. 1998. "Methylglyoxal
504     Production in Bacteria: Suicide or Survival?" *Archives of Microbiology* 170 (4): 209–18.
505     Flores-Uribe, J., G. Hevroni, and R. Ghai. 2019. "Heliorhodopsins Are Absent in Diderm
506     (Gram-negative) Bacteria: Some Thoughts and Possible Implications for Activity."
507     *Environmental Microbiology Reports*.
508     https://onlinelibrary.wiley.com/doi/abs/10.1111/1758-2229.12730.
509     Galperin, Michael Y., Kira S. Makarova, Yuri I. Wolf, and Eugene V. Koonin. 2015.
510     "Expanded Microbial Genome Coverage and Improved Protein Family Annotation in the
511     COG Database." *Nucleic Acids Research* 43 (Database issue): D261-9.
512     García-Domínguez, M., J. C. Reyes, and F. J. Florencio. 1999. "Glutamine Synthetase
513     Inactivation by Protein-Protein Interaction." *Proceedings of the National Academy of*
514     *Sciences of the United States of America* 96 (13): 7161–66.
515     Ghai, Rohit, Katherine D. McMahon, and Francisco Rodriguez-Valera. 2012. "Breaking a
516     Paradigm: Cosmopolitan and Abundant Freshwater Actinobacteria Are Low GC."
517     *Environmental Microbiology Reports* 4 (1): 29–35.
518     Haft, Daniel H., Jeremy D. Selengut, and Owen White. 2003. "The TIGRFAMs Database of
519     Protein Families." *Nucleic Acids Research* 31 (1): 371–73.
520     Hashimoto, Masanori, Kota Katayama, Yuji Furutani, and Hideki Kandori. 2020. "Zinc
521     Binding to Heliorhodopsin." *Journal of Physical Chemistry Letters* 11 (20): 8604–9.
522     Hauser, Maria, Martin Steinegger, and Johannes Söding. 2016. "MMseqs Software Suite
523     for Fast and Deep Clustering and Searching of Large Protein Sequence Sets."
524     *Bioinformatics* 32 (9): 1323–30.
525     Huynen, M., B. Snel, W. Lathe 3rd, and P. Bork. 2000. "Predicting Protein Function by
526     Genomic Context: Quantitative Evaluation and Qualitative Inferences." *Genome Research*
527     10 (8): 1204–10.
528     Hyatt, Doug, Gwo-Liang Chen, Philip F. Locascio, Miriam L. Land, Frank W. Larimer, and
529     Loren J. Hauser. 2010. "Prodigal: Prokaryotic Gene Recognition and Translation Initiation
530     Site Identification." *BMC Bioinformatics* 11 (March): 119.
531     Ikuta, Tatsuya, Wataru Shihoya, Masahiro Sugiura, Kazuho Yoshida, Masahito Watari,
532     Takaya Tokano, Keitaro Yamashita, et al. 2020. "Structural Insights into the Mechanism of
533     Rhodopsin Phosphodiesterase." *Nature Communications* 11 (1): 5605.

534    Im, Yang Ju, Amanda J. Davis, Imara Y. Perera, Eva Johannes, Nina S. Allen, and Wendy F.
535    Boss. 2007. "The N-Terminal Membrane Occupation and Recognition Nexus Domain of
536    Arabidopsis Phosphatidylinositol Phosphate Kinase 1 Regulates Enzyme Activity." *The*
537    *Journal of Biological Chemistry* 282 (8): 5443–52.
538    Johnson, Norman L., Adrienne W. Kemp, and Samuel Kotz. 2005. *Univariate Discrete*
539    *Distributions*. John Wiley & Sons.
540    Kajava, Andrey V. 2012. "Tandem Repeats in Proteins: From Sequence to Structure."
541    *Journal of Structural Biology* 179 (3): 279–88.
542    Käll, Lukas, Anders Krogh, and Erik L. L. Sonnhammer. 2005. "An HMM Posterior Decoder
543    for Sequence Feature Prediction That Includes Homology Information." *Bioinformatics* 21
544    Suppl 1 (June): i251-7.
545    Kandori, Hideki. 2020. "Biophysics of Rhodopsins and Optogenetics." *Biophysical Reviews*
546    12 (2): 355–61.
547    Kanehisa, Minoru, Yoko Sato, and Kanae Morishima. 2016. "BlastKOALA and
548    GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome
549    Sequences." *Journal of Molecular Biology* 428 (4): 726–31.
550    Kanno, Manabu, Hideyuki Tamaki, Yasuo Mitani, Nobutada Kimura, Satoshi Hanada, and
551    Yoichi Kamagata. 2015. "PH-Induced Change in Cell Susceptibility to Butanol in a High
552    Butanol-Tolerant Bacterium, Enterococcus Faecalis Strain CM4A." *Biotechnology for*
553    *Biofuels* 8 (April): 69.
554    Katoh, Kazutaka, and Daron M. Standley. 2013. "MAFFT Multiple Sequence Alignment
555    Software Version 7: Improvements in Performance and Usability." *Molecular Biology and*
556    *Evolution* 30 (4): 772–80.
557    Kavagutti, Vinicius S., Adrian-Ştefan Andrei, Maliheh Mehrshad, Michaela M. Salcher, and
558    Rohit Ghai. 2019. "Phage-Centric Ecological Interactions in Aquatic Ecosystems Revealed
559    through Ultra-Deep Metagenomics." *Microbiome* 7 (1): 135.
560    Kelley, Lawrence A., Stefans Mezulis, Christopher M. Yates, Mark N. Wass, and Michael J.
561    E. Sternberg. 2015. "The Phyre2 Web Portal for Protein Modeling, Prediction and
562    Analysis." *Nature Protocols* 10 (6): 845–58.
563    Kim, Suhyun, Ilnam Kang, Ji-Hui Seo, and Jang-Cheon Cho. 2019. "Culturing the
564    Ubiquitous Freshwater Actinobacterial AcI Lineage by Supplying a Biochemical 'helper'
565    Catalase." *The ISME Journal* 13 (9): 2252–63.
566    Kovalev, K., D. Volkov, R. Astashkin, A. Alekseev, I. Gushchin, J. M. Haro-Moreno, I.
567    Chizhov, et al. 2020. "High-Resolution Structural Insights into the Heliorhodopsin Family."
568    *Proceedings of the National Academy of Sciences of the United States of America* 117 (8):
569    4131–41.
570    Krishna, S. Sri, Indraneel Majumdar, and Nick V. Grishin. 2003. "Structural Classification of
571    Zinc Fingers: Survey and Summary." *Nucleic Acids Research* 31 (2): 532–50.
572    Li, Dinghua, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-Fung Ting, Kunihiko
573    Sadakane, Hiroshi Yamashita, and Tak-Wah Lam. 2016. "MEGAHIT v1.0: A Fast and
574    Scalable Metagenome Assembler Driven by Advanced Methodologies and Community
575    Practices." *Methods* 102 (June): 3–11.
576    Li, Weizhong, and Adam Godzik. 2006. "Cd-Hit: A Fast Program for Clustering and
577    Comparing Large Sets of Protein or Nucleotide Sequences." *Bioinformatics* 22 (13): 1658–
578    59.

579   Lillig, Christopher Horst, Carsten Berndt, and Arne Holmgren. 2008. "Glutaredoxin
580   Systems." *Biochimica et Biophysica Acta* 1780 (11): 1304–17.
581   Ma, Hui, Ying Lou, Wen Hui Lin, and Hong Wei Xue. 2006. "MORN Motifs in Plant PIPKs
582   Are Involved in the Regulation of Subcellular Localization and Phospholipid Binding." *Cell
583   Research* 16 (5): 466–78.
584   Maclean, Michelle, Scott J. MacGregor, John G. Anderson, and Gerry Woolsey. 2009.
585   "Inactivation of Bacterial Pathogens Following Exposure to Light from a 405-Nanometer
586   Light-Emitting Diode Array." *Applied and Environmental Microbiology* 75 (7): 1932–37.
587   Maresca, Julia A., Jessica L. Keffer, Priscilla P. Hempel, Shawn W. Polson, Olga
588   Shevchenko, Jaysheel Bhavsar, Deborah Powell, Kelsey J. Miller, Archana Singh, and
589   Martin W. Hahn. 2019. "Light Modulates the Physiology of Nonphototrophic
590   Actinobacteria." *Journal of Bacteriology* 201 (10). https://doi.org/10.1128/JB.00740-18.
591   Megrian, Daniela, Najwa Taib, Jerzy Witwinowski, Christophe Beloin, and Simonetta
592   Gribaldo. 2020. "One or Two Membranes? Diderm Firmicutes Challenge the Gram-
593   Positive/Gram-Negative Divide." *Molecular Microbiology* 113 (3): 659–71.
594   Mehrshad, Maliheh, Michaela M. Salcher, Yusuke Okazaki, Shin-Ichi Nakano, Karel Šimek,
595   Adrian-Stefan Andrei, and Rohit Ghai. 2018. "Hidden in Plain Sight-Highly Abundant and
596   Diverse Planktonic Freshwater Chloroflexi." *Microbiome* 6 (1): 176.
597   Mitchell, Alex L., Alexandre Almeida, Martin Beracochea, Miguel Boland, Josephine
598   Burgin, Guy Cochrane, Michael R. Crusoe, et al. 2020. "MGnify: The Microbiome Analysis
599   Resource in 2020." *Nucleic Acids Research* 48 (D1): D570–78.
600   Mitchell, Alex L., Teresa K. Attwood, Patricia C. Babbitt, Matthias Blum, Peer Bork, Alan
601   Bridge, Shoshana D. Brown, et al. 2019. "InterPro in 2019: Improving Coverage,
602   Classification and Access to Protein Sequence Annotations." *Nucleic Acids Research* 47
603   (D1): D351–60.
604   Neuenschwander, Stefan M., Rohit Ghai, Jakob Pernthaler, and Michaela M. Salcher.
605   2018. "Microdiversification in Genome-Streamlined Ubiquitous Freshwater
606   Actinobacteria." *The ISME Journal* 12 (1): 185–98.
607   Oubrie, A., H. J. Rozeboom, K. H. Kalk, A. J. Olsthoorn, J. A. Duine, and B. W. Dijkstra.
608   1999. "Structure and Mechanism of Soluble Quinoprotein Glucose Dehydrogenase." *The
609   EMBO Journal* 18 (19): 5187–94.
610   Parks, Donovan H., Maria Chuvochina, Pierre-Alain Chaumeil, Christian Rinke, Aaron J.
611   Mussig, and Philip Hugenholtz. 2020. "A Complete Domain-to-Species Taxonomy for
612   Bacteria and Archaea." *Nature Biotechnology* 38 (9): 1079–86.
613   Pushkarev, Alina, Keiichi Inoue, Shirley Larom, José Flores-Uribe, Manish Singh, Masae
614   Konno, Sahoko Tomida, et al. 2018. "A Distinct Abundant Group of Microbial Rhodopsins
615   Discovered Using Functional Metagenomics." *Nature* 558 (7711): 595–99.
616   Rastogi, Rajesh P., Richa, Ashok Kumar, Madhu B. Tyagi, and Rajeshwar P. Sinha. 2010.
617   "Molecular Mechanisms of Ultraviolet Radiation-Induced DNA Damage and Repair."
618   *Journal of Nucleic Acids* 2010 (December): 592980.
619   Resto, Melissa, Jason Yaffe, and Barbara Gerratana. 2009. "An Ancestral Glutamine-
620   Dependent NAD(+) Synthetase Revealed by Poor Kinetic Synergism." *Biochimica et
621   Biophysica Acta* 1794 (11): 1648–53.
622   Rouhier, Nicolas, Jérémy Couturier, Michael K. Johnson, and Jean-Pierre Jacquot. 2010.
623   "Glutaredoxins: Roles in Iron Homeostasis." *Trends in Biochemical Sciences* 35 (1): 43–52.

624  Saiki, Takashi, Yasuhiko Kobayashi, Kiyotaka Kawagoe, and Teruhiko Beppu. 1985.
625  "Dictyoglomus Thermophilum Gen. Nov., Sp. Nov., a Chemoorganotrophic, Anaerobic,
626  Thermophilic Bacterium." *International Journal of Systematic and Evolutionary*
627  *Microbiology* 35 (3): 253–59.
628  Sajko, S., I. Grishkovskaya, J. Kostan, and M. Graewert. 2020. "Structures of Three MORN
629  Repeat Proteins and a Re-Evaluation of the Proposed Lipid-Binding Properties of MORN
630  Repeats." *BioRxiv*. https://www.biorxiv.org/content/10.1101/826180v2.abstract.
631  Salazar, Guillem, Lucas Paoli, Adriana Alberti, Jaime Huerta-Cepas, Hans-Joachim
632  Ruscheweyh, Miguelangel Cuenca, Christopher M. Field, et al. 2019. "Gene Expression
633  Changes and Community Turnover Differentially Shape the Global Ocean
634  Metatranscriptome." *Cell* 179 (5): 1068-1083.e21.
635  Shibata, Mikihiro, Keiichi Inoue, Kento Ikeda, Masae Konno, Manish Singh, Chihiro
636  Kataoka, Rei Abe-Yoshizumi, Hideki Kandori, and Takayuki Uchihashi. 2018. "Oligomeric
637  States of Microbial Rhodopsins Determined by High-Speed Atomic Force Microscopy and
638  Circular Dichroic Spectroscopy." *Scientific Reports* 8 (1): 8262.
639  Shihoya, Wataru, Keiichi Inoue, Manish Singh, Masae Konno, Shoko Hososhima, Keitaro
640  Yamashita, Kento Ikeda, et al. 2019. "Crystal Structure of Heliorhodopsin." *Nature* 574
641  (7776): 132–36.
642  Shmakov, Sergey A., Kira S. Makarova, Yuri I. Wolf, Konstantin V. Severinov, and Eugene V.
643  Koonin. 2018. "Systematic Prediction of Genes Functionally Linked to CRISPR-Cas Systems
644  by Gene Neighborhood Analysis." *Proceedings of the National Academy of Sciences of*
645  *the United States of America* 115 (23): E5307–16.
646  Taib, Najwa, Daniela Megrian, Jerzy Witwinowski, Panagiotis Adam, Daniel Poppleton,
647  Guillaume Borrel, Christophe Beloin, and Simonetta Gribaldo. 2020. "Genome-Wide
648  Analysis of the Firmicutes Illuminates the Diderm/Monoderm Transition." *Nature Ecology*
649  *& Evolution*, October. https://doi.org/10.1038/s41559-020-01299-7.
650  Takeshima, H., S. Komazaki, M. Nishi, M. Iino, and K. Kangawa. 2000. "Junctophilins: A
651  Novel Family of Junctional Membrane Complex Proteins." *Molecular Cell* 6 (1): 11–22.
652  Tanaka, Tatsuki, Manish Singh, Wataru Shihoya, Keitaro Yamashita, Hideki Kandori, and
653  Osamu Nureki. 2020. "Structural Basis for Unique Color Tuning Mechanism in
654  Heliorhodopsin." *Biochemical and Biophysical Research Communications* 533 (3): 262–67.
655  Timmers, Peer H. A., Charlotte D. Vavourakis, Robbert Kleerebezem, Jaap S. Sinninghe
656  Damsté, Gerard Muyzer, Alfons J. M. Stams, Dimity Y. Sorokin, and Caroline M. Plugge.
657  2018. "Metabolism and Occurrence of Methanogenic and Sulfate-Reducing Syntrophic
658  Acetate Oxidizing Communities in Haloalkaline Environments." *Frontiers in Microbiology* 9
659  (December): 3039.
660  UniProt Consortium. 2019. "UniProt: A Worldwide Hub of Protein Knowledge." *Nucleic*
661  *Acids Research* 47 (D1): D506–15.
662  Vavourakis, Charlotte D., Adrian-Stefan Andrei, Maliheh Mehrshad, Rohit Ghai, Dimity Y.
663  Sorokin, and Gerard Muyzer. 2018. "A Metagenomics Roadmap to the Uncultured
664  Genome Diversity in Hypersaline Soda Lake Sediments." *Microbiome* 6 (1): 168.
665  Vavourakis, Charlotte D., Maliheh Mehrshad, Cherel Balkema, Rutger van Hall, Adrian-
666  Ştefan Andrei, Rohit Ghai, Dimity Y. Sorokin, and Gerard Muyzer. 2019. "Metagenomes
667  and Metatranscriptomes Shed New Light on the Microbial-Mediated Sulfur Cycle in a
668  Siberian Soda Lake." *BMC Biology* 17 (1): 69.

669    Yatsunami, Rie, Ai Ando, Ying Yang, Shinichi Takaichi, Masahiro Kohno, Yuriko
670    Matsumura, Hiroshi Ikeda, et al. 2014. "Identification of Carotenoids from the Extremely
671    Halophilic Archaeon Haloarcula Japonica." *Frontiers in Microbiology* 5 (March): 100.
672    Zhao, Shanrong, Ying Zhang, William Gordon, Jie Quan, Hualin Xi, Sarah Du, David von
673    Schack, and Baohong Zhang. 2015. "Comparison of Stranded and Non-Stranded RNA-
674    Seq Transcriptome Profiling and Investigation of Gene Overlap." *BMC Genomics* 16
675    (September): 675.
676    Zimmermann, Lukas, Andrew Stephens, Seung-Zin Nam, David Rau, Jonas Kübler, Marko
677    Lozajic, Felix Gabler, Johannes Söding, Andrei N. Lupas, and Vikram Alva. 2018. "A
678    Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server at Its
679    Core." *Journal of Molecular Biology* 430 (15): 2237–43.

1    **Supplementary Information**

2    **Heliorhodopsin evolution is driven by photosensory promiscuity in monoderms**

3    Paul-Adrian Bulzu[1], Vinicius Silva Kavagutti[1,2], Maria-Cecilia Chiriac[1], Charlotte

4    D. Vavourakis[3], Keiichi Inoue[4], Hideki Kandori[5,6], Adrian-Stefan Andrei[7], Rohit

5    Ghai[1]*

6    [1]Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre of

7    the Academy of Sciences of the Czech Republic, České Budějovice, Czech Republic.

8    [2]Department of Ecosystem Biology, Faculty of Science, University of South Bohemia,

9    Branišovská 1760, České Budějovice, Czech Republic.

10   [3]EUTOPS, Research Institute for Biomedical Aging Research, University of Innsbruck,

11   Austria.

12   [4]The Institute for Solid State Physics, The University of Tokyo, Kashiwa, Japan.

13   [5]Department of Life Science and Applied Chemistry, Nagoya Institute of Technology,

14   Showa, Nagoya 466-8555, Japan.

15   [6]OptoBioTechnology Research Center, Nagoya Institute of Technology, Showa, Nagoya

16   466-8555, Japan

17   [7]Limnological Station, Department of Plant and Microbial Biology, University of Zurich,

18   Kilchberg, Switzerland.


19   *Corresponding author: Rohit Ghai

20   Department of Aquatic Microbial Ecology, Institute of Hydrobiology, Biology Centre of the

21   Academy

22   of Sciences of the Czech Republic, Na Sádkách 7, 370 05, České Budějovice, Czech

23   Republic.

24   Phone: +420 387 775 881

25   Fax: +420 385 310 248

26   E-mail: ghai.rohit@gmail.com

27  We examined the distribution and co-occurrence of HeRs and type-1 rhodopsins
28  (Supplementary Tables S4-S7) in the GTDB database (release 89), as it has been previously
29  suggested that these proteins tend to coexist within the same organisms (Kovalev et al.
30  2020). From all 24,706 scanned genomes we identified and retrieved 1,455 *bona fide*
31  type-1 rhodopsin-containing genomes from which 69 (4.74 %) proved to also harbour
32  HeRs. Since 15.33 % from all identified HeRs (n = 450) co-occur with type-1 rhodopsins we
33  consider it as being more an exception than a norm and find no data to sustain a
34  physiological dependency between these two rhodopsin families.

## MORN-protein domain fusions

36  The bacteria encoding MORN-HeRs were previously predicted to be strict anaerobes
37  (Timmers et al. 2018). Although mostly recovered from sediments, these MAGs also
38  encode other proteins directly or indirectly associated with the presence of light (i.e.
39  bacteriophytochrome COG4251, DNA repair photolyase COG1533,
40  Deoxyribodipyrimidine photolyase COG0415). Therefore, the co-occurrence of strict
41  anaerobiosis and light-dependent components indicates the top sediment layer as their
42  likely habitat. The available 3D structure of MORN-repeats shows a single repeat to consist
43  of short beta-pleated regions folded back upon themselves, creating a flat surface area
44  that expands when the repeats are present in multiple tandem copies (Wilson et al., 2002).
45  The presence of periplasmic proteins with MORN-Big_2 or MORN-Big2-PASTA domains
46  may indicate the extracellular MORN repeats as adaptors between typical sensor domains
47  (i.e. PASTA/Big_2) and the transducer protein kinases acting in the cytoplasm
48  (Supplementary Figure 3). However, as the MORN-repeats that are fused to HeRs in these
49  organisms are intracellular, they could only interact with MORN-repeats on the
50  cytoplasmic side. In one such case, intracellular MORN-repeats were fused to ATPase
51  component of an ABC-type efflux pump (likely involved in drug resistance or Cu2+/Na2+
52  ion efflux). The other candidate found was a PknB protein kinase-FHA-MORN repeat fusion
53  that was predicted to be in an atypical membrane orientation. In such proteins,
54  dimerization domains (e.g. PASTA) are extracellular (Supplementary Figure 3C) while
55  cytoplasmic protein kinase domains function as part of signal transduction pathways in a
56  wide range of gram-positive bacteria (Kang et al. 2005). The dimerization of PknB is
57  essential for autophosphorylation and activation of the kinases through an allosteric
58  mechanism (Lombana et al. 2010). The predicted reverse orientation (if correct) of this
59  protein renders MORN-mediated interactions with HeR unlikely in these organisms.
60  However, the presence of MORN-repeats in PknB type proteins for which dimerization is
61  essential for function (most likely via MORN-repeats) strengthens the possibility that the
62  MORN-repeats aid in MORN-HeR dimerization as well.

## Various Cysteine-rich motif-containing heliorhodopsins

64  Another extension found in four sequences (N-terminal variant 2 or ntv2 in Supplementary
65  Table S8) with many more cysteines (n = 10) that did not give any significant hits to known
66  proteins was also identified. The presence of multiple cysteine residues and a conserved
67  tryptophan residue is reminiscent of RINGv finger domains that coordinate two metals.
68  However, RING finger domains also have a highly conserved histidine residue that was not
69  detected. Indeed, we also found at least two instances in freshwater Actinobacteria (*Ca.*
70  *Planktophila*) where a RINGv domain-containing protein is located right upstream of the
71  heliorhodopsin gene and in the same orientation (see Figure 3).

72 A third variant (ntv3) with seven cysteines was found (in Thermoplasmatales and in
73 Euryarchaeota) but without conserved histidines or tryptophan. However, it shows broad
74 similarity with ntv3 in presence of conserved prolines and arginines before the cysteine
75 motifs. Apart from the n-terminal extensions motifs, at least three sequences of the
76 intracellular loop (ICL3), were also found to be rich in cysteines (n = 8) which also
77 presented a conserved tryptophan similar to ntv2 described above (intracellular loop 3
78 variant 1, icl3v1).

79 Thus, at least 17 sequences presented cysteine-rich motifs either at the N-terminus (ntv1,
80 ntv2 and ntv3) or in the intracellular loop (icl3v1) at the cytoplasmic side of
81 heliorhodopsins suggesting the possibility that these might be transducers of the
82 conformational change in heliorhodopsins upon light excitation. The conserved cysteines
83 in these proteins could bind either iron or zinc and are likely redox-active. While we did
84 find many other types of N-terminal extensions, we were unable to find any significant hits
85 to these even by sensitive sequence searches (Supplementary Table S8).

## Sources of retinal for HeR function

87 We also found several NAD-dependent short-chain dehydrogenases that might encode
88 for retinol dehydrogenases in the vicinity of HeRs (adh_short in Figure 3, Supplementary
89 Table S9). It has been mentioned before that as HeRs are able to efficiently capture retinal
90 from exogenous sources and that HeR-encoding microbes do not have a retinal
91 biosynthesis pathway (Shihoya et al. 2019).

92 *De novo* retinal biosynthesis requires five genes that if supplied in-trans to a non-retinal
93 producing microbe may result in functional rhodopsins (Sabehi et al. 2005). The final step
94 of the *de novo* pathway uses a beta-carotene monooxygenase that converts beta-carotene
95 to retinal. However, retinal may also be converted from retinol by the action of retinol
96 dehydrogenases. We used the curated GTDB database to further probe the co-
97 occurrence of HeR and type-1 rhodopsins along with genes for retinal biosynthesis. Of the
98 total of 381 genomes that encoded only HeR, we find that only a single genome encoded
99 all genes necessary for the *de novo* production of retinal, but 213 (55%) also encoded the
100 beta-carotene monooxygenase and 241 genomes (63%) encoded at least one retinol
101 dehydrogenase (Supplementary Table S10). Considering genomes that encoded only
102 type-1 rhodopsins (n = 1,386), 596 (43%) encoded the complete pathway for retinal
103 biosynthesis and additionally 995 (71%) also encoded at least one retinol dehydrogenase.
104 It appears that microbes encoding only HeR mostly lack the complete pathway for *de novo*
105 retinal biosynthesis and that apart from exogenous capture of retinal, conversion from
106 beta-carotene (via beta-carotene monooxygenase) or from retinol (via retinol
107 dehydrogenases) may be at work.

## HeR genomic context

109 We performed gene context analysis of HeRs by combining the maximum-likelihood
110 phylogenetic tree generated for representative HeR sequences (n = 872) with HeR gene
111 neighbourhood information (Supp. Figure 4; iTOL:
112 https://itol.embl.de/tree/14723125092152021608050562). The resulting tree places most
113 HeR sequences (n = 835) within 19 conspicuous phylogenetic clusters which we further
114 denominate as C1-C19 (see Supp. Figure 15). Among them, Actinobacteriota-encoded
115 HeRs are by far the most numerous (n = 533) accounting for eight well-defined clusters

116  (i.e. C1-7 and C11) and a small sub-cluster within Patescibacteria (CPR)-dominated C17.
117  Regarding Actinobacteriota, C1 and C3 are represented by order Nanopelagicales, C2
118  includes chiefly members of Microtrichales (class Acidimicrobiia), C4 comprises
119  Actinomycetales HeRs and C5 includes classes Coriobacteriia and Thermoleophilia.
120  Notably, C5 brings together HeRs recovered mainly from lesser studied sediment habitats
121  including nine Chloroflexota (class Dehalococcoidia) sequences. Predominantly marine
122  C6 includes Acidimicrobiia HeRs from order Microtrichales and other poorly classified
123  representatives from within this class while C7 has both marine and freshwater
124  Microtrichales together with Propionibacteriales genus *Nocardioides*.

125  Cluster C11 stands out in this analysis due to the high level of evolutionary conservation of
126  both HeRs and their neighbouring genes, bringing together exclusively members of
127  marine Actinobacteria from order *Ca.* Actinomarinales (TMED189). Importantly, in C11 we
128  notice that synteny is only conserved among genes sharing the same orientation as HeR
129  while gene "gains" and/or "losses" occur only in the opposite orientation. Despite very
130  high phylogenetic relatedness within C11 and therefore the unsurprisingly similar gene
131  context amongst its members, the differences between (+) and (-) strand feature
132  conservation (relative to HeR orientation) indicate a potentially relevant transcriptional unit
133  comprised of genes: afuA - iron(III) transport system substrate-binding protein (K02012),
134  *afuB* - iron(III) transport system permease protein (K02011), *afuC* - iron(III) transport system
135  ATP-binding protein (K02010), *HeR* – Heliorhodopsin (PF18761), an 11-subunit respiratory
136  complex I operon (*nuoA, B, C, D, H, I, J, K, L, M, N*), *NDUFAF7* - NADH dehydrogenase
137  [ubiquinone] 1 alpha subcomplex assembly factor 7 (K18164/PF02636), *htpX* - heat shock
138  protein (K03799), *pspE* - phage shock protein E (K03972) containing a *Rhodanese*
139  (PF00581) domain and *adenosine_kinase* (cd01168) (Note: full-length annotated contigs
140  deposited in Figshare). In summary, HeRs from C11 are always preceded by genes
141  encoding a complete ABC-type ferric iron uptake system and followed by an operon
142  encoding a 11-subunit, "ancestral"-type (Moparthi and Hägerhäll 2011) respiratory
143  complex I and by accessory components required for correct assembly and function of this
144  complex (*NDUFAF7, htpX, pspE*) (Zurita Rendón et al. 2014; Pagani and Galante 1983;
145  Alexander and Volini 1987; Sakoh, Ito, and Akiyama 2005). The last conserved gene
146  encodes a pfkB family adenosine kinase (cd01168), a key purine salvage enzyme that
147  phosphorylates adenosine to generate adenosine monophosphate (AMP) (Long, Escuyer,
148  and Parker 2003).

149  The presence of HeRs within the same transcriptional unit as above mentioned energy
150  metabolism components could indicate them as modulators or even light-induced sensory
151  "switches" of such processes, a mechanism perhaps similar to cryptochrome-driven
152  metabolic synchronization with substrate availability described in other Actinobacteria
153  (Maresca et al. 2019). Notably, beside the 11-subunit complex I, that lacks the NADH
154  dehydrogenase module (subunits nuoE, nuoF, nuoG) (Moparthi and Hägerhäll 2011), *Ca.*
155  *Actinomarina* (TMED189) genomes also encode the full-sized 14-subunit variant of
156  respiratory complex I in close proximity to the first (for example in GCA-902516125.1).
157  Curiously, the association of HeRs with complex I genes is reminiscent of that between the
158  transmembrane, sensory, EAL-domain containing protein seen in *Bacillus cereus* located
159  upstream of a similar 11-subunit complex I operon (Moparthi and Hägerhäll 2011).

160     The last cluster featuring a significant number of Actinobacteria HeRs (n = 12) is C17. In
161     this CPR-dominated cluster, Actinobacteria HeRs form a well-defined group sharing a
162     common ancestor with a small CPR sub-cluster. While gene context appears conserved
163     within these Actinobacteria, this does not apply to the putative "sister" CPR sub-cluster.

164     Clusters C8-C10 share a common ancestor with C11 and include HeRs encoded largely in
165     strict or facultative anaerobic prokaryotes recovered from sediments (including activated
166     sludge). Notably, the phylogenetic tree (iTOL link above) shows two Asgardarchaeota
167     (class Heimdallarchaeia) HeRs branching with very high support (SH-test/UFBoot =
168     97.2/98) as a sister clade to all C8-C10 members, after the split with the common ancestor
169     shared with C11. Despite the high support for the Asgardarchaeota HeR split, defining a
170     credible cluster will require including additional sequences once more genomes become
171     available. The basal, C10 cluster, is comprised of Archaea-derived HeRs from phyla
172     Crenarchaeota (class Bathyarchaeia) and Euryarchaeota (class Methanobacteria). Clusters
173     C8 and mainly C9 include the MORN-HeR encoding Firmicutes as well as a few
174     Chloroflexota (class Anaerolineae) HeRs.

175     C12-C16 form a separate super-cluster showing moderate-to-low support for internal
176     branching patterns and include mostly, although not exclusively, anaerobic Archaea (C12
177     and C14, with the notable exception of aerobic Halobacterota within C12) and anaerobic
178     Chloroflexota (classes Anaerolineae - C13, C15 and Dehalococcoidia – C16). Notably, C16
179     includes one Thermoplasmatota HeR (encoded in contig SRR5506739-C331) with a zinc-
180     finger extension (Znf-HeR) at the N-terminus. C18 is the most basal cluster with confidently
181     assigned taxonomy. It includes exclusively aerobic Chloroflexota members of the
182     Ellin6529 lineage.

183     Although no consensus taxonomy could be determined for members of C19 – the first
184     branching group after the split with proteorhodopsins, the abundance of "eukaryotic
185     signature proteins" (ESP) (e.g. Arf, Roc, Rab, etc.) points towards either a eukaryotic origin
186     or unidentified, ESP-rich archaea (Hartman and Fedorov 2002; Dong, Wen, and Tian
187     2007).

188     An extended phylogenetic tree of HeRs, including additional dereplicated sequences
189     identified and retrieved from UniProtKB and GTDB, is available in FigShare (see Supp.
190     Methods - Phylogenetic tree of HeRs).


191     **Supplementary Methods**

192     **Re-assembling of HeR-encoding Spirochaeta.** The unexpected detection of HeR in a
193     previously published *Spirochaeta* (diderm organism) genome prompted further
194     investigation. The original Illumina short-read dataset SRX2623364 was downloaded from
195     NCBI SRA (Sequence Read Archive) and preprocessed by using a combination of tools
196     provided by the BBMap project (https://sourceforge.net/projects/bbmap/). This involved
197     removing poor-quality reads with bbduk.sh (qtrim = rl, trimq = 18), identifying phiX and p-
198     Fosil2 control reads (k = 21) and removing Illumina sequencing adapters (k = 21). Further,
199     *de novo* assembly of preprocessed paired-end reads was done by Megahit v1.2.9 (D. Li et
200     al. 2016) with k-mer list: 29, 39, 49, 59, 69, 79, 89, 99, 109, 119, 127, and with default
201     parameters. A total of 886 contigs with an average length of 4.98 kbp were produced.

202  Protein-coding genes were predicted by Prodigal (Hyatt et al. 2010) and taxonomically
203  classified by scanning with MMSeqs2 against the GTDB database. A *Spirochaeta* contig
204  (length = 304,032 bp) was identified and scanned for the presence of HeR against the
205  PFAM database. Both taxonomy (*Spirochaeta*) and HeR presence were consistent with
206  previously published results.

207  **Phylogenetic tree of HeRs**. An extensive collection of predicted HeR amino acid
208  sequences (n = 4,108) was generated from: 1) all HeR sequences available in UniProtKB (n
209  = 502), 2) HeR identified within dereplicated contigs (using cd-hit-est -c 0.95 -aS 0.95)
210  assembled from publicly available metagenomes and metatranscriptomes (n = 3,145), 3)
211  HeR identified within dereplicated high-quality MAGs included in the GTDB database (n =
212  455) and 4) HeR assembled from the strand-specific metatranscriptomic dataset
213  generated in this study (n = 6). This collection was simplified by keeping only
214  representative sequences (n = 1,669) chosen following clustering with MMSeqs2
215  (Steinegger and Söding 2017) at 90% sequence identity and 90% coverage (mode: easy-
216  cluster; -c 0.90; --min-seq-id 0.90). Representative HeR sequences were aligned together
217  with 30 selected proteorhodopsins serving as outgroup by using PASTA (Mirarab et al.
218  2015) (resulting alignment with 1,699 sequences, 2,486 columns, 2,264 distinct patterns,
219  1,091 parsimony-informative sites, 698 singleton sites, 697 constant sites). A Maximum
220  Likelihood (ML) phylogenetic tree was constructed with IQ-TREE2 (Minh et al. 2020) (1,000
221  iterations for ultrafast bootstrapping (Hoang et al. 2018) and SH testing, respectively; best
222  model chosen by ModelFinder (Kalyaanamoorthy et al. 2017): LG+G4; additional
223  parameters recommended for short sequence alignments -nstop 500 -pers 0.2). The
224  generated tree was annotated to include labels containing: HeR-encoding contig name,
225  habitat of origin and consensus GTDB taxonomic classification (if available). Data including
226  alignment and the annotated phylogenetic tree are deposited in FigShare
227  (https://figshare.com/s/7bb42426f2ad5e891fec).

228  **Phylogenetic tree of HeRs with gene context**. A simplified depiction of HeR genomic
229  context across representative taxonomic groups harbouring such genes was constructed
230  by merging HeR phylogenetic information with available HeR gene neighbourhood data
231  (Supplementary Figure 4). For this purpose, we established a collection of representative,
232  dereplicated HeR-encoding contigs (n = 872) of at least 5 kb and with clear consensus
233  taxonomy from two sources: 1) HeR-encoding contigs assembled from publicly available
234  metagenomes and metatranscriptomes (n = 3,145) and 2) HeR contigs assembled from
235  the strand-specific metatranscriptomic dataset generated in this study (n = 6).
236  Dereplication of contigs was previously achieved using cd-hit-est (W. Li and Godzik 2006)
237  with identity cutoffs of 95% and coverage of 95% (-c 0.95 -aS 0.95). **Phylogenetic tree
238  building:** Curated HeR sequences recovered from selected contigs were aligned together
239  with 30 proteorhodopsins serving as outgroup by using PASTA (Mirarab et al. 2015)
240  (resulting alignment with 902 sequences with 957 columns, 895 distinct patterns, 550
241  parsimony-informative sites, 198 singleton sites, 209 constant sites). A Maximum
242  Likelihood (ML) phylogenetic tree was constructed with IQ-TREE2 (Minh et al. 2020) (1,000
243  iterations for ultrafast bootstrapping (Hoang et al. 2018) and SH testing, respectively; best
244  model chosen by ModelFinder (Kalyaanamoorthy et al. 2017): LG+I+G4; additional
245  parameters recommended for short sequence alignments -nstop 500 -pers 0.2).

246 **Reconstruction of HeR gene neighborhoods:** Coding sequences were predicted *de*
247 *novo* using Prodigal (Hyatt et al. 2010) in metagenomic (-p meta) mode. Protein domains
248 were annotated by scanning predicted coding sequences against the PFAM (Protein
249 Families) database using the publicly available perl script pfam_scan.pl. Predicted protein
250 sequences from all contigs were clustered together using the MMSeqs2 easy-cluster
251 workflow with 50% identity and 80% coverage cutoffs (--min-seq-id 0.5 -c 0.8 -e 1e-3 --
252 cluster-reassign) and a minimum of 2 sequences per cluster. Clusters were sorted
253 according to their size (i.e. number of sequences) and colour codes were assigned to the
254 top largest 63. All clusters containing HeRs were assigned matching colours. HeR gene
255 neighbourhood data was combined with the reconstructed HeR phylogenetic tree in iTOL
256 (Letunic and Bork 2016) (https://itol.embl.de/). To facilitate visual interpretation, a number
257 of adjustments and rules were applied: 1) all contigs were oriented according to the sense
258 (+) of encoded HeRs, 2) contigs were centered on HeR genes with a maximum of 10
259 neighbouring genes depicted up- and downstream, 3) information regarding gene
260 lengths was not included, all of them being shown as equally sized rectangles, 4)
261 homologous genes (i.e. members of the same MMSeqs2 defined cluster) share matching
262 colours within each phylogenetically defined cluster, 5) all HeRs are coloured the same
263 across all phylogenetic clusters, 6) grey rectangles indicate genes with few homologues
264 and/or singletons, 7) taxonomy, habitat information and phylogenetic clusters are colour
265 coded on independent strips.

266 **MORN-HeR phylogenomic tree**. A maximum-likelihood (ML) phylogenomic tree was
267 constructed for Firmicutes MAGs encoding MORN-HeR protein domain fusions along with
268 other representatives of this phylum that assumably lack such genes. The established
269 collection of (n = 68) MAGs and reference genomes (Supplementary Table S15) was
270 scanned by hmmsearch against a previously published list of (n = 120) conserved protein
271 marker HMMs (Parks et al. 2018). Four divergent markers (TIGR00442, TIGR00539,
272 TIGR00643, TIGR00717) were identified by scanning with CD-search (e-value < 1e-2) and
273 removed. Curated amino acid sequences for the selected 116 phylogenetic markers were
274 aligned with PRANK (Löytynoja 2014) and resulting alignments were trimmed by BMGE
275 (parameters: -t AA -g 0.5 -b 3 -m BLOSUM30) (Criscuolo and Gribaldo 2010). Individually
276 trimmed alignments were concatenated resulting in a block of 68 sequences with 40,714
277 columns, 36,088 distinct patterns, 28,978 parsimony informative sites, 3,366 singleton
278 sites and 8,370 constant sites. The concatenated alignment was used with IQ-TREE
279 (v.1.6.12) to construct the ML phylogenomic tree (parameters: 1,000 iterations of ultrafast
280 bootstrapping (Hoang et al. 2018) and SH testing (Minh et al. 2020), respectively; best
281 model (LG+F+R5) chosen by ModelFinder (Kalyaanamoorthy et al. 2017).
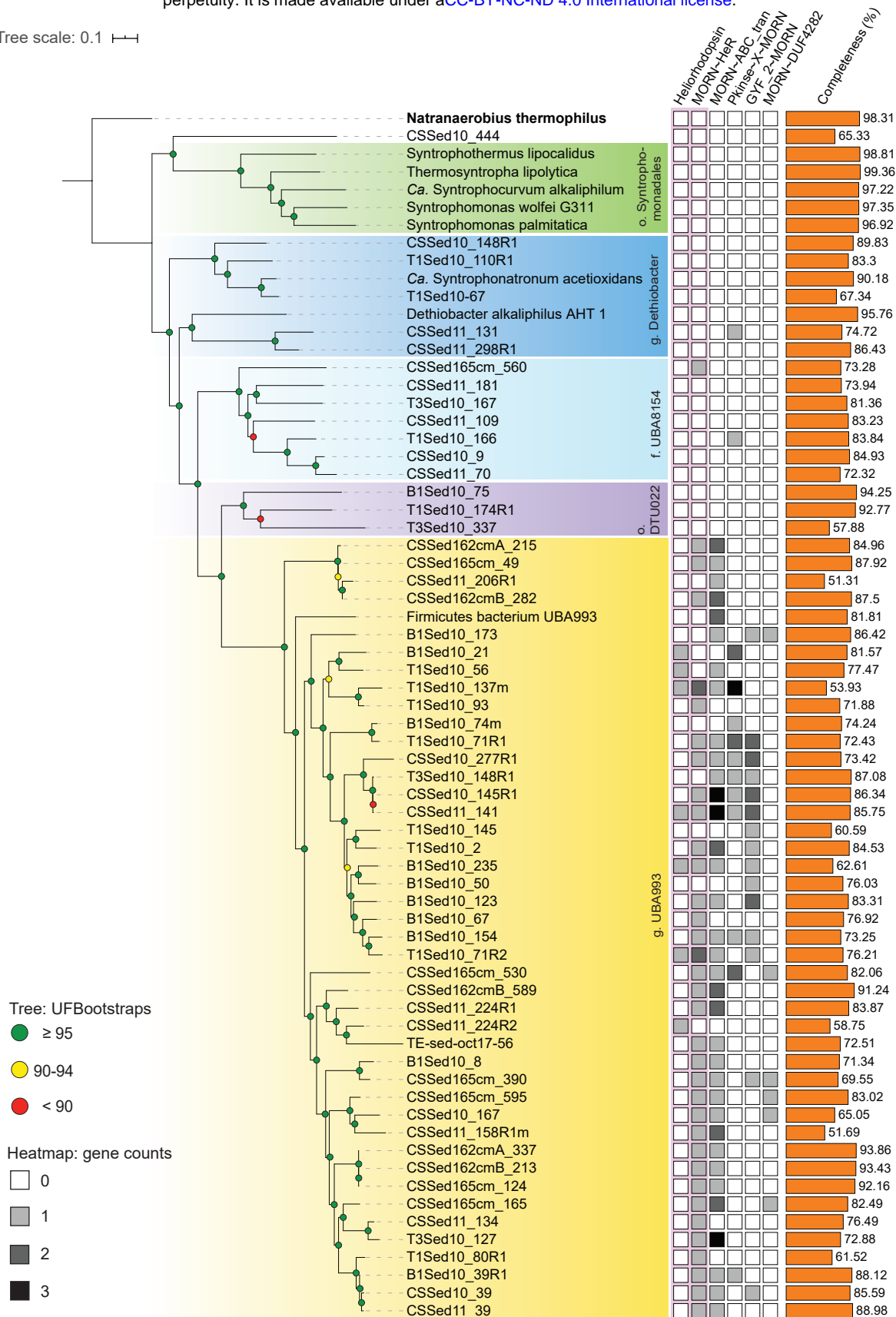
282 **Multiple sequence alignment (MSA) of MORN-HeR proteins.** Predicted amino acid
283 sequences containing full-length MORN-MORN-MORN-HeR protein domain fusions (n =
284 36) were retrieved from Firmicutes MAGs (n = 35) previously used to construct the
285 phylogenomic tree presented in Supplementary Figure 1. All sequences were aligned
286 together using the PSI-Coffee alignment method (Chang et al. 2012) provided by the T-
287 Coffee online server (http://tcoffee.crg.cat) with default parameters. The resulting MSA is
288 presented with annotations in Supplementary Figure 2 while the original alignment file
289 generated by PSI-Coffee is available in FigShare
290 (https://figshare.com/s/7bb42426f2ad5e891fec).

## Supplementary References
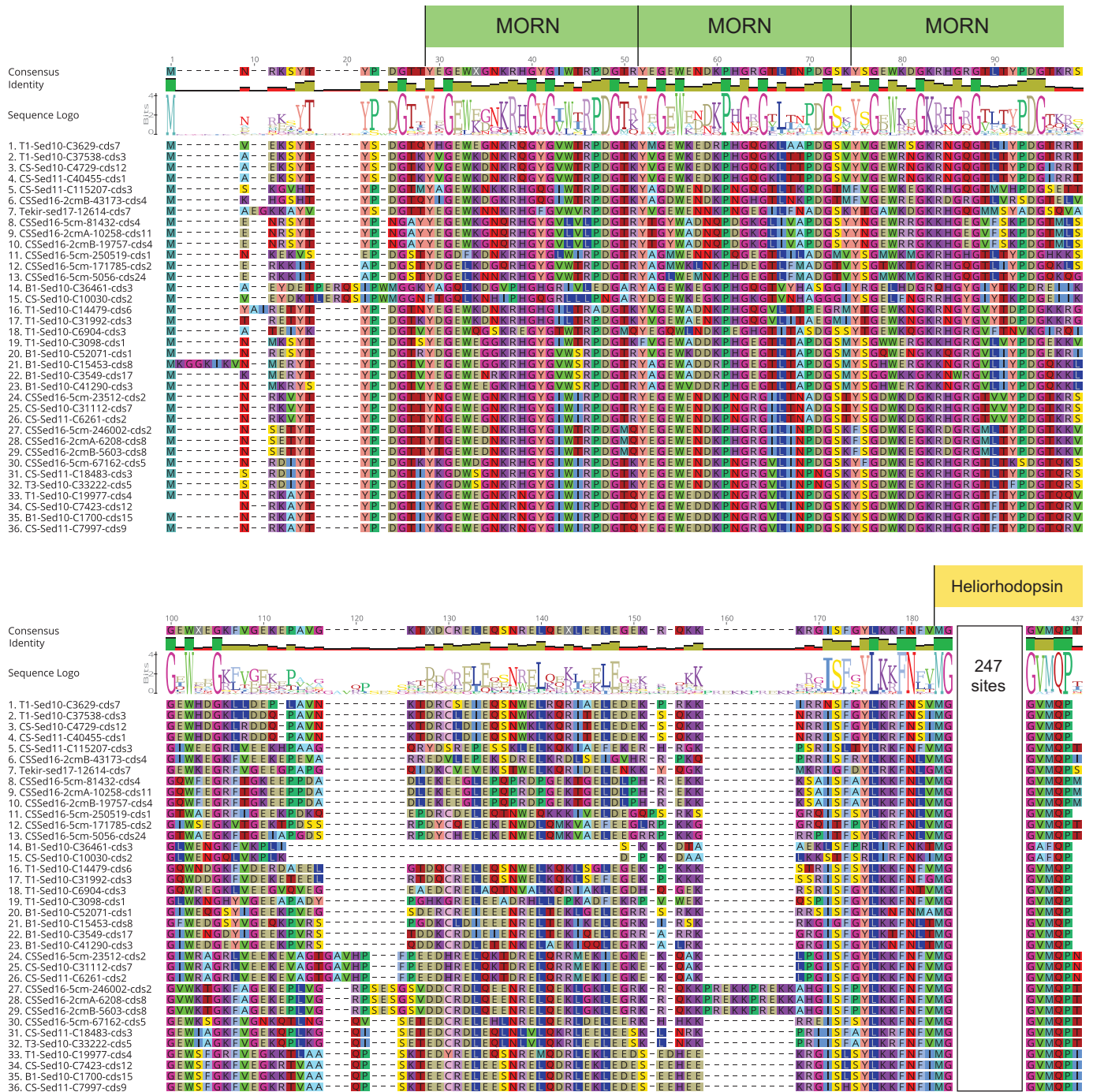
Alexander, K., and M. Volini. 1987. "Properties of an Escherichia Coli Rhodanese." *The Journal of Biological Chemistry* 262 (14): 6595–6604.

Chang, Jia-Ming, Paolo Di Tommaso, Jean-François Taly, and Cedric Notredame. 2012. "Accurate Multiple Sequence Alignment of Transmembrane Proteins with PSI-Coffee." *BMC Bioinformatics* 13 Suppl 4 (March): S1.

Criscuolo, Alexis, and Simonetta Gribaldo. 2010. "BMGE (Block Mapping and Gathering with Entropy): A New Software for Selection of Phylogenetic Informative Regions from Multiple Sequence Alignments." *BMC Evolutionary Biology* 10 (July): 210.

Dong, Jiu-Hong, Jian-Fan Wen, and Hai-Feng Tian. 2007. "Homologs of Eukaryotic Ras Superfamily Proteins in Prokaryotes and Their Novel Phylogenetic Correlation with Their Eukaryotic Analogs." *Gene* 396 (1): 116–24.

Hartman, Hyman, and Alexei Fedorov. 2002. "The Origin of the Eukaryotic Cell: A Genomic Investigation." *Proceedings of the National Academy of Sciences of the United States of America* 99 (3): 1420–25.

Hoang, Diep Thi, Olga Chernomor, Arndt von Haeseler, Bui Quang Minh, and Le Sy Vinh. 2018. "UFBoot2: Improving the Ultrafast Bootstrap Approximation." *Molecular Biology and Evolution* 35 (2): 518–22.

Hyatt, Doug, Gwo-Liang Chen, Philip F. Locascio, Miriam L. Land, Frank W. Larimer, and Loren J. Hauser. 2010. "Prodigal: Prokaryotic Gene Recognition and Translation Initiation Site Identification." *BMC Bioinformatics* 11 (March): 119.

Kalyaanamoorthy, Subha, Bui Quang Minh, Thomas K. F. Wong, Arndt von Haeseler, and Lars S. Jermiin. 2017. "ModelFinder: Fast Model Selection for Accurate Phylogenetic Estimates." *Nature Methods* 14 (6): 587–89.

Kang, Choong-Min, Derek W. Abbott, Sang Tae Park, Christopher C. Dascher, Lewis C. Cantley, and Robert N. Husson. 2005. "The Mycobacterium Tuberculosis Serine/Threonine Kinases PknA and PknB: Substrate Identification and Regulation of Cell Shape." *Genes & Development* 19 (14): 1692–1704.

Kovalev, K., D. Volkov, R. Astashkin, A. Alekseev, I. Gushchin, J. M. Haro-Moreno, I. Chizhov, et al. 2020. "High-Resolution Structural Insights into the Heliorhodopsin Family." *Proceedings of the National Academy of Sciences of the United States of America* 117 (8): 4131–41.

Letunic, Ivica, and Peer Bork. 2016. "Interactive Tree of Life (ITOL) v3: An Online Tool for the Display and Annotation of Phylogenetic and Other Trees." *Nucleic Acids Research* 44 (W1): W242-5.

Li, Dinghua, Ruibang Luo, Chi-Man Liu, Chi-Ming Leung, Hing-Fung Ting, Kunihiko Sadakane, Hiroshi Yamashita, and Tak-Wah Lam. 2016. "MEGAHIT v1.0: A Fast and Scalable Metagenome Assembler Driven by Advanced Methodologies and Community Practices." *Methods* 102 (June): 3–11.

Li, Weizhong, and Adam Godzik. 2006. "Cd-Hit: A Fast Program for Clustering and Comparing Large Sets of Protein or Nucleotide Sequences." *Bioinformatics* 22 (13): 1658–59.

Lombana, T. Noelle, Nathaniel Echols, Matthew C. Good, Nathan D. Thomsen, Ho-Leung Ng, Andrew E. Greenstein, Arnold M. Falick, David S. King, and Tom Alber. 2010. "Allosteric Activation Mechanism of the Mycobacterium Tuberculosis Receptor Ser/Thr Protein Kinase, PknB." *Structure* 18 (12): 1667–77.

336  Long, Mary C., Vincent Escuyer, and William B. Parker. 2003. "Identification and
337        Characterization of a Unique Adenosine Kinase from Mycobacterium Tuberculosis."
338        *Journal of Bacteriology* 185 (22): 6548–55.
339  Löytynoja, Ari. 2014. "Phylogeny-Aware Alignment with PRANK." *Methods in Molecular*
340        *Biology*  1079: 155–70.
341  Maresca, Julia A., Jessica L. Keffer, Priscilla P. Hempel, Shawn W. Polson, Olga Shevchenko,
342        Jaysheel Bhavsar, Deborah Powell, Kelsey J. Miller, Archana Singh, and Martin W.
343        Hahn. 2019. "Light Modulates the Physiology of Nonphototrophic Actinobacteria."
344        *Journal of Bacteriology* 201 (10). https://doi.org/10.1128/JB.00740-18.
345  Minh, Bui Quang, Heiko A. Schmidt, Olga Chernomor, Dominik Schrempf, Michael D.
346        Woodhams, Arndt von Haeseler, and Robert Lanfear. 2020. "IQ-TREE 2: New Models
347        and Efficient Methods for Phylogenetic Inference in the Genomic Era." *Molecular*
348        *Biology and Evolution* 37 (5): 1530–34.
349  Mirarab, Siavash, Nam Nguyen, Sheng Guo, Li-San Wang, Junhyong Kim, and Tandy Warnow.
350        2015. "PASTA: Ultra-Large Multiple Sequence Alignment for Nucleotide and Amino-
351        Acid Sequences." *Journal of Computational Biology: A Journal of Computational*
352        *Molecular Cell Biology* 22 (5): 377–86.
353  Moparthi, Vamsi K., and Cecilia Hägerhäll. 2011. "The Evolution of Respiratory Chain
354        Complex I from a Smaller Last Common Ancestor Consisting of 11 Protein Subunits."
355        *Journal of Molecular Evolution* 72 (5–6): 484–97.
356  Pagani, S., and Y. M. Galante. 1983. "Interaction of Rhodanese with Mitochondrial NADH
357        Dehydrogenase." *Biochimica et Biophysica Acta* 742 (2): 278–84.
358  Parks, Donovan H., Maria Chuvochina, David W. Waite, Christian Rinke, Adam Skarshewski,
359        Pierre-Alain Chaumeil, and Philip Hugenholtz. 2018. "A Standardized Bacterial
360        Taxonomy Based on Genome Phylogeny Substantially Revises the Tree of Life."
361        *Nature Biotechnology* 36 (10): 996–1004.
362  Sabehi, Gazalah, Alexander Loy, Kwang-Hwan Jung, Ranga Partha, John L. Spudich, Tal
363        Isaacson, Joseph Hirschberg, Michael Wagner, and Oded Béjà. 2005. "New Insights
364        into Metabolic Properties of Marine Bacteria Encoding Proteorhodopsins." *PLoS*
365        *Biology*. https://doi.org/10.1371/journal.pbio.0030273.
366  Sakoh, Machiko, Koreaki Ito, and Yoshinori Akiyama. 2005. "Proteolytic Activity of HtpX, a
367        Membrane-Bound and Stress-Controlled Protease from Escherichia Coli." *The Journal*
368        *of Biological Chemistry* 280 (39): 33305–10.
369  Shihoya, Wataru, Keiichi Inoue, Manish Singh, Masae Konno, Shoko Hososhima, Keitaro
370        Yamashita, Kento Ikeda, et al. 2019. "Crystal Structure of Heliorhodopsin." *Nature*
371        574 (7776): 132–36.
372  Steinegger, Martin, and Johannes Söding. 2017. "MMseqs2 Enables Sensitive Protein
373        Sequence Searching for the Analysis of Massive Data Sets." *Nature Biotechnology* 35
374        (11): 1026–28.
375  Timmers, Peer H. A., Charlotte D. Vavourakis, Robbert Kleerebezem, Jaap S. Sinninghe
376        Damsté, Gerard Muyzer, Alfons J. M. Stams, Dimity Y. Sorokin, and Caroline M.
377        Plugge. 2018. "Metabolism and Occurrence of Methanogenic and Sulfate-Reducing
378        Syntrophic Acetate Oxidizing Communities in Haloalkaline Environments." *Frontiers*
379        *in Microbiology* 9 (December): 3039.
380  Zurita Rendón, Olga, Lissiene Silva Neiva, Florin Sasarman, and Eric A. Shoubridge. 2014.
381        "The Arginine Methyltransferase NDUFAF7 Is Essential for Complex I Assembly and
382        Early Vertebrate Embryogenesis." *Human Molecular Genetics* 23 (19): 5159–70.
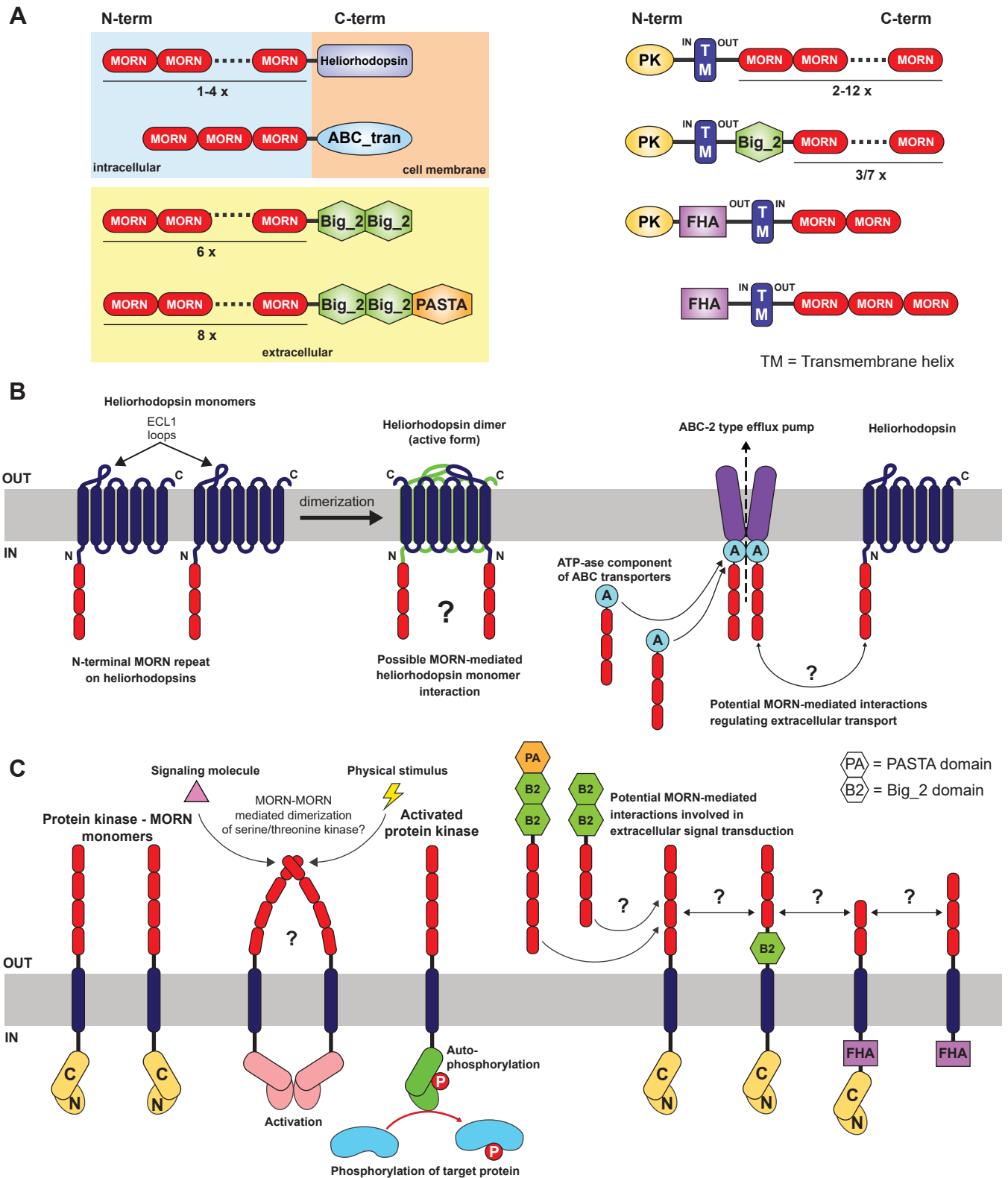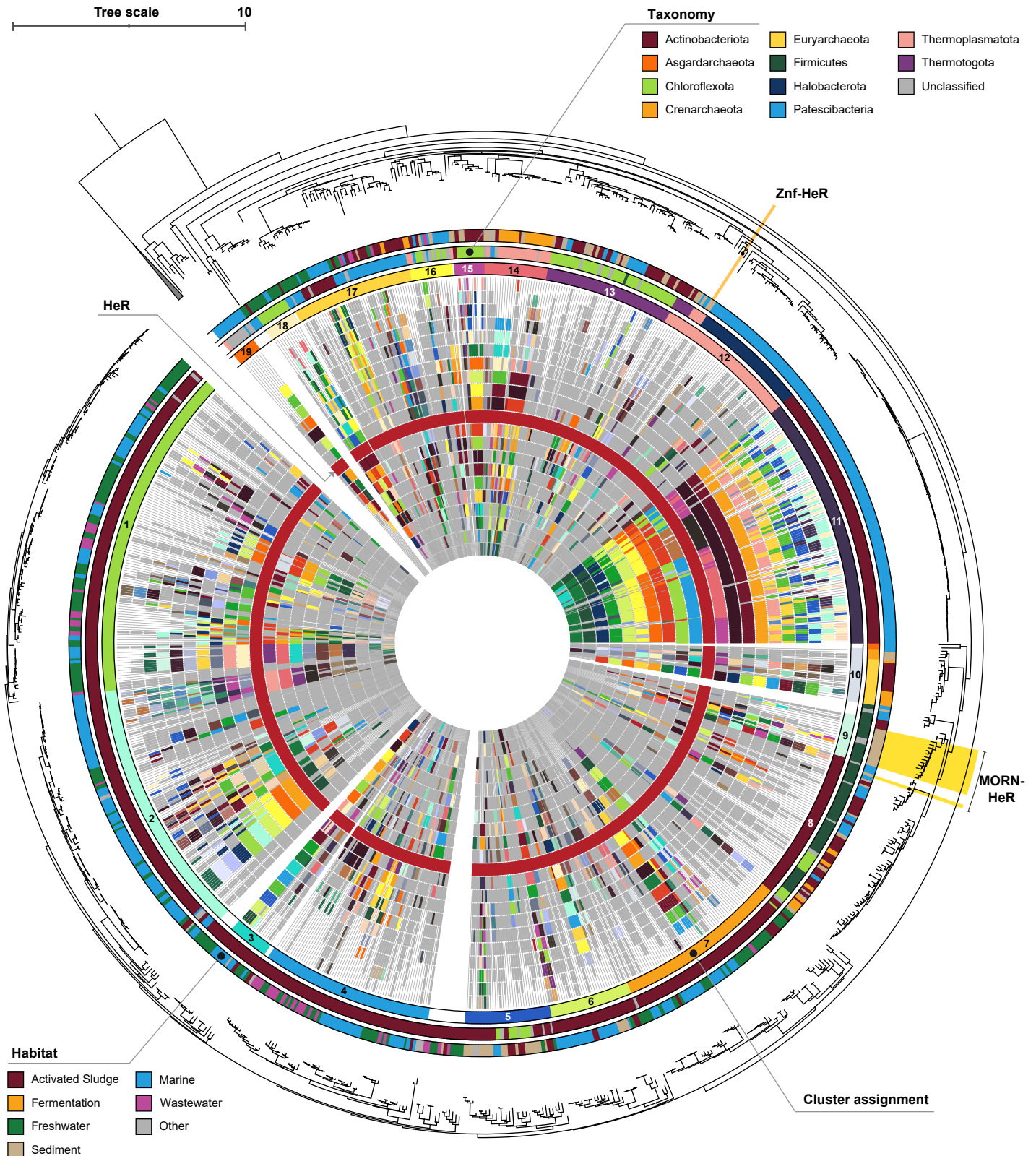
**Supplementary Figure 1.** Phylogenomic tree of MORN-Heliorhodopsin (MORN-HeR) encoding MAGs. Green circles indicate high confidence UFBootstrap values (≥ 95). Occurences of genes encoding heliorhodopsins, MORN-heliorhodopsin as well as other MORN-domain fusions are depicted in the adjacent matrix. Genome completeness values are depicted as a histogram (estimated by CheckM). All genomes are members of the *Firmicutes* phylum, with different taxonomic subdivisions highlighted on the tree (taxonomy by GTDBtk). *Natranaerobius thermophilus* was used to root the tree. The majority of genomes included here have been previously used for phylogenomic analyses by Timmers et. al, 2018. Among reference genomes, *Firmicutes bacterium UBA993* and *Ca. Syntrophocurvum alkaliphilum* are new additions.

**Supplementary Figure 2.** Multiple sequence alignment of MORN-HeR protein domain fusions predicted in Firmicutes MAGs. Each sequence shows 3 consecutive MORN domains (indicated by green rectangles). Aligned HeR domains display a high level of conservation and are truncated for illustration purposes (indicated by yellow rectangle). The original full-length alignment was deposited in FigShare.

**Supplementary Figure 3**. Summary of MORN-repeat proteins predicted from metagenome-assembled genomes (MAGs) of anaerobic Firmicutes. **A**). Schematics of frequently co-occuring MORN-repeat containing proteins in Firmicutes MAGs including MORN-HeR, MORN-ABC transporters - likely involved in drug resistance or $Cu^{2+}/Na^{2+}$ ion efflux, periplasmic proteins containing bacterial immunoglobulin-like folds (Big_2, PASTA), MORN-protein kinase fusions where MORN repeats and p-kinase domains are commonly separated by transmembrane α-helices on opposite sides of the cellular membrane and MORN-forkhead associated (FHA) domains. **B**) Potential interactions between MORN-HeR monomers, MORN-ATPase components of ABC transporters and MORN-HeR and ABC transporters mediated or stabilized by the presence of MORN-repeat fusions. **C**) Potential interactions of MORN-Protein-kinases associated with functions such as extracellular signal transduction.

**Supplementary Figure 4.** Genomic context of heliorhodopsin (HeR) genes across representative taxonomic groups. The phylogenetic tree was constructed using 872 HeR amino acid sequences and 30 proteorhodopsins used as an outgroup for rooting. Gene neighbourhoods (10 genes up- and downstream) for each HeR were depicted schematically. Abundant homologues are coloured within each defined phylogenetic cluster while less abundant genes and/or singletons are depicted in gray. All contigs were centered and oriented according to encoded HeR (dark red circle). Information regarding relative gene lengths was not included. Particular HeR-protein domain fusions are indicated separately: Znf-HeR and MORN-HeR.

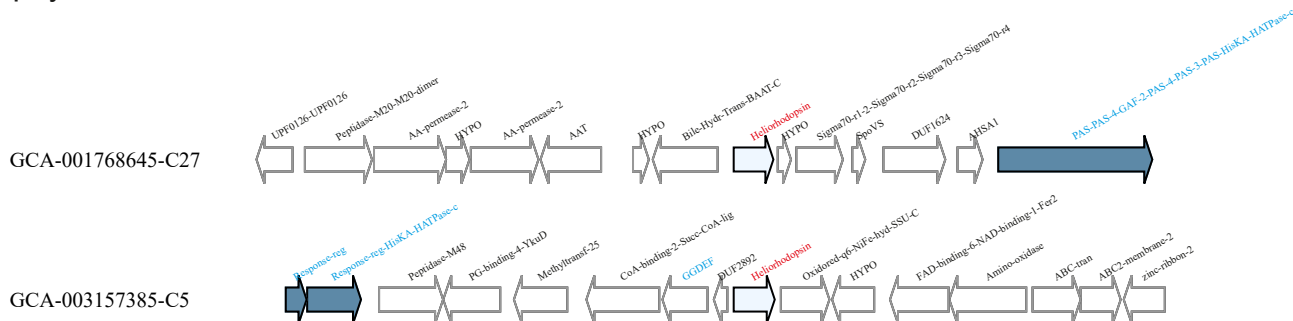**Supplementary Figure 5.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Chloroflexi.
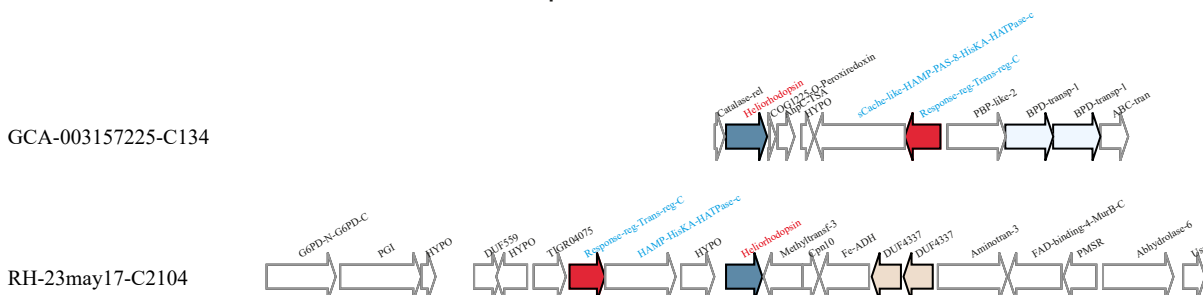
**Supplementary Figure 6.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Firmicutes.

**Supplementary Figure 7.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Patescibacteria.

**Supplementary Figure 8.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Actinobacteriota, class Acidimicrobiia.

**Supplementary Figure 9:** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Actinobacteriota, class Actinobacteria.

**Supplementary Figure 10.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Actinobacteriota, class Coriobacteriia.

**Supplementary Figure 11.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Actin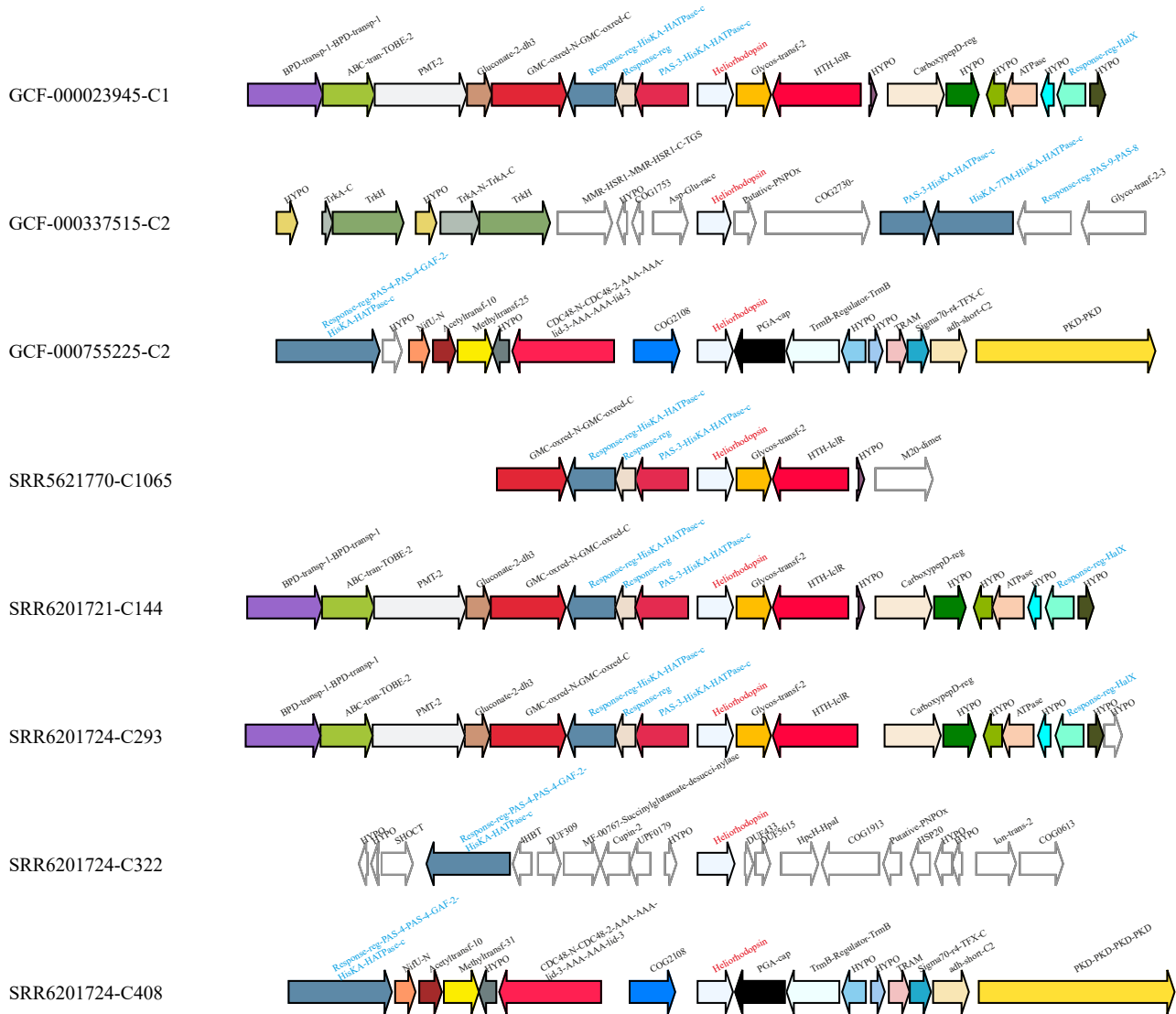obacteriota (classes RBG-13-55-18 and Thermoleophilia), phylum Dictyoglomota (class Dictyoglomia) and phylum Thermoplasmatota (class E2).

**Supplementary Figure 12.** Histidine kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of Heliorhodopsins (red labels) originating from the phylum Halobacterota, class Halobacteria.
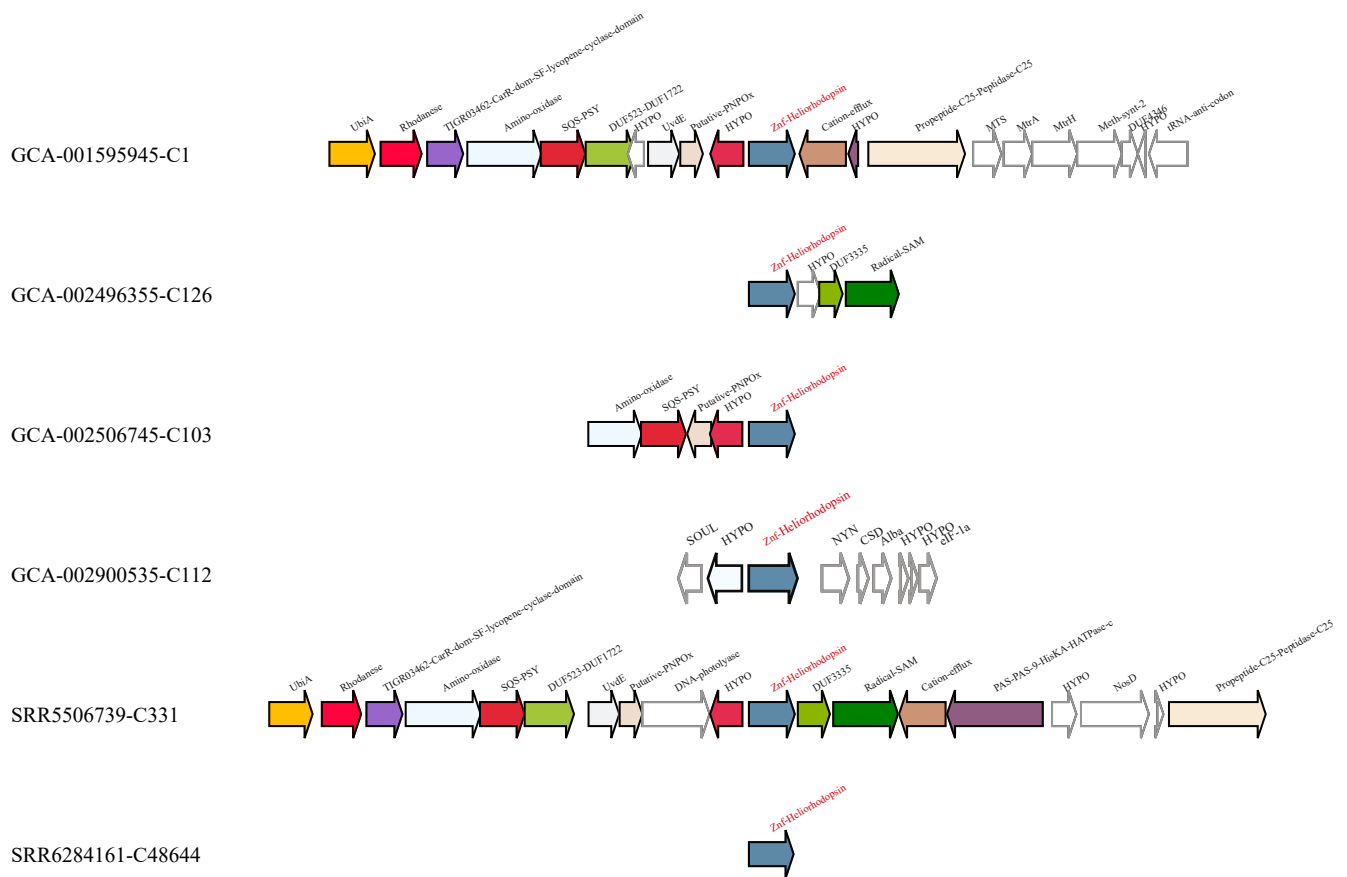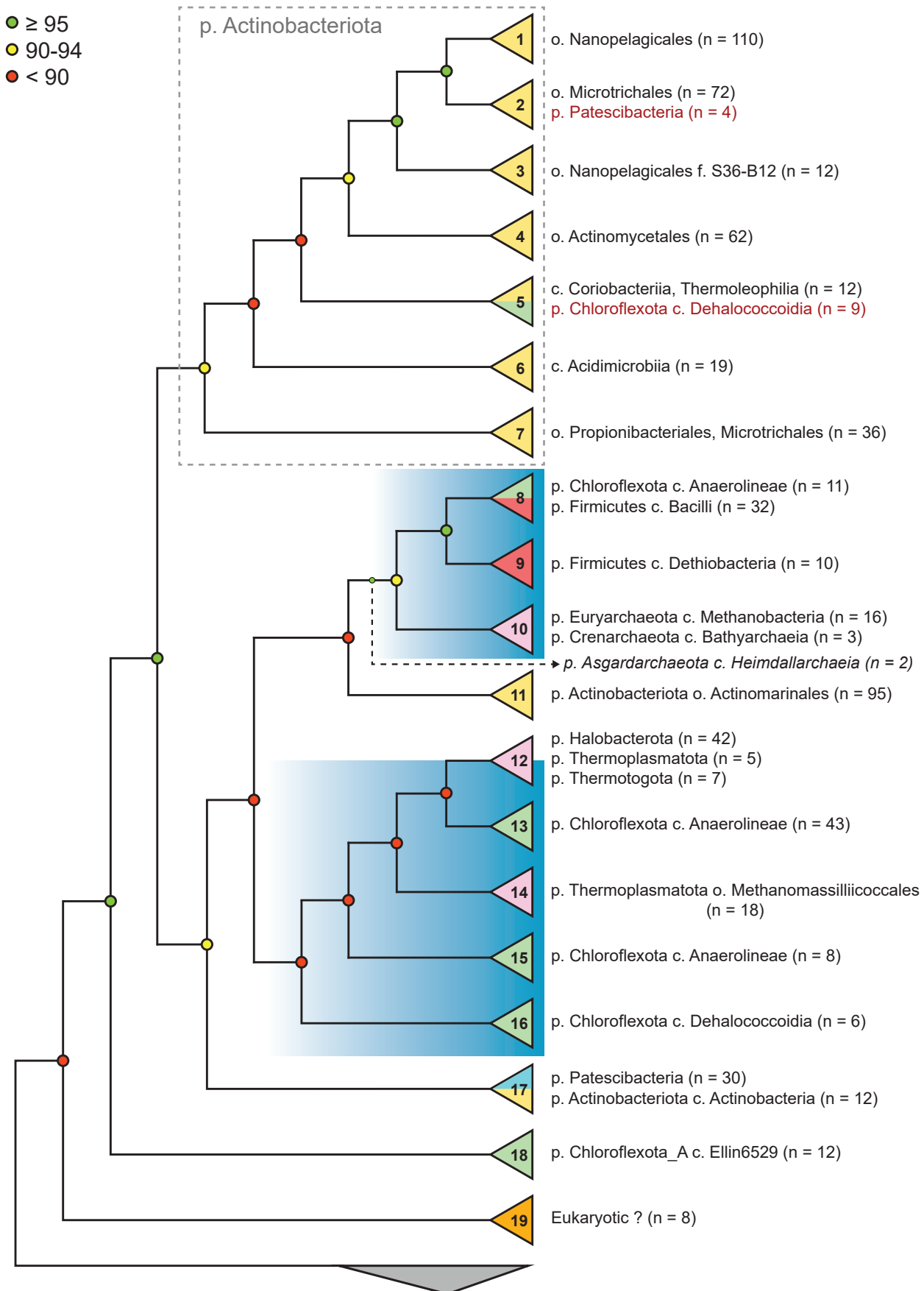
**Supplementary Figure 13.** Histidine-kinases and related genes (blue labels) in the genomic neighborhood (within 10 kb) of MORN-Heliorhodopsins (red labels) originating from Firmicutes (phylum) MAGs (sediment and brine metagenomes).

**Supplementary Figure 14.** Genomic neighborhood (within 10 kb) of Znf-Heliorhodopsins (red labels) originating from archaeal contigs from the phylum Thermoplasmatales, class E2.

**Supplementary Figure 15.** Simplified representation (cladogram) of the HeR phylogenetic tree used for gene context analysis. Cluster numbers (defined in Supp. Fig. 4) are indicated on triangles at the tip of each branch. Actinobacteriota clusters are coloured yellow, Archaea - purple, Chloroflexota - green, Eukaryota - orange, Firmicutes - red, Patescibacteria - blue. Taxonomy and sequence counts are shown only for representatives of each cluster. The blue rectangles highlight clusters where all or most members are anaerobic organisms. The outgroup (proteorhodopsins; n = 30) is depicted as a gray triangle at the bottom.