1    **Title:** Common ancestry of heterodimerizing TALE homeobox transcription factors across

2    Metazoa and Archaeplastida

3

4    **Short title**: Origins of TALE transcription factor networks

5

6    Sunjoo Joo[1,5], Ming Hsiu Wang[1,5], Gary Lui[1], Jenny Lee[1], Andrew Barnas[1], Eunsoo Kim[2],

7    Sebastian Sudek[3], Alexandra Z. Worden[3], and Jae-Hyeok Lee[1,4]¶

8

9    1. Department of Botany, University of British Columbia, 6270 University Blvd., Vancouver,

10   BC V6T 1Z4, Canada

11   2. Division of Invertebrate Zoology and Sackler Institute for Comparative Genomics,

12   American Museum of Natural History, 200 Central Park West, New York, NY 10024, United

13   States

14   3. Monterey Bay Aquarium Research Institute, 7700 Sandholdt Rd., Moss Landing, CA

15   95039, United States

16   4. Corresponding author

17   5. Equally contributed

18

19   Corresponding Author:

20   Jae-Hyeok Lee

21   #2327-6270 University Blvd.

22   1-604-827-5973

23   jae-hyeok.lee@botany.ubc.ca

24

25  **Abstract**

26

27  Homeobox transcription factors (TFs) in the TALE superclass are deeply embedded in the

28  gene regulatory networks that orchestrate embryogenesis. Knotted-like homeobox (KNOX)

29  TFs, homologous to animal MEIS, have been found to drive the haploid-to-diploid transition

30  in both unicellular green algae and land plants via heterodimerization with other TALE

31  superclass TFs, representing remarkable functional conservation of a developmental TF

32  across lineages that diverged one billion years ago. To delineate the ancestry of TALE-TALE

33  heterodimerization, we analyzed TALE endowment in the algal radiations of Archaeplastida,

34  ancestral to land plants. Homeodomain phylogeny and bioinformatics analysis partitioned

35  TALEs into two broad groups, KNOX and non-KNOX. Each group shares previously defined

36  heterodimerization domains, plant KNOX-homology in the KNOX group and animal PBC-

37  homology in the non-KNOX group, indicating their deep ancestry. Protein-protein interaction

38  experiments showed that the TALEs in the two groups all participated in heterodimerization.

39  These results indicate that the TF dyads consisting of KNOX/MEIS and PBC-containing

40  TALEs must have evolved early in eukaryotic evolution, a likely function being to accurately

41  execute the haploid-to-diploid transitions during sexual development.

42

43  **Author summary**

44  Complex multicellularity requires elaborate developmental mechanisms, often based on the

45  versatility of heterodimeric transcription factor (TF) interactions. Highly conserved TALE-

46  superclass homeobox TF networks in major eukaryotic lineages suggest deep ancestry of

47  developmental mechanisms. Our results support the hypothesis that in early eukaryotes, the

48  TALE heterodimeric configuration provided transcription-on switches via dimerization-

49  dependent subcellular localization, ensuring execution of the haploid-to-diploid transition

50  only when the gamete fusion is correctly executed between appropriate partner gametes, a

51  system that then diversified in the several lineages that engage in complex multicellular

52  organization.

53

54  **Keywords**: Archaeplastida evolution; developmental mechanism; KNOX transcription factor;

55  PBC-homology; TALE-class homeobox; transcription factor heterodimerization

56

57  **Introduction**

58

59  The homeobox transcription factors (TFs) are ubiquitous in eukaryotes, carrying a DNA-

60  binding homeodomain typically 60 amino acids, that folds into three $\alpha$-helices [1]. The

61  atypical or TALE (Three Amino acid Length Extension) superclass of homeobox TFs shares

62  a three-amino-acid insertion between helix 1 and 2 and plays essential roles during

63  embryonic development by participating in interactive TF networks. In animals, MEIS- and

64  PBC-class TALE proteins, such as Meis/Hth and Pbx/Exd, form heterodimers that in turn

65  form ternary complexes with HOX-class homeobox TFs, determining cellular fates along the

66  anterior-posterior axis of the developing embryo [2,3]. In plants, the interacting KNOX- and

67  BELL-class TFs in the TALE group play critical roles during organ formation and the

68  vegetative-to-reproductive transition in the undifferentiated cell mass known as the shoot

69  apical meristem [4,5].

70

71  The heterodimerization of TALE proteins serves as a trigger for precise execution of

72  developmental programs. Prior to heterodimerization, animal PBX proteins are localized in

73  the cytosol, and upon binding to MEIS, they translocate to the nucleus [6,7]. Similar

74  heterodimerization-dependent translocation is also observed for KNOX-BELL pairs in the

75  plant *Arabidopsis*, implying that this mechanism is a conserved regulatory feature of TALE

76  proteins [8]. In addition, TALE proteins differ in their DNA-binding specificity [9,10], which is

77  primarily determined by the homeodomain residues at positions 47, 50, and 54 [11], and

78  heterodimerization increases target affinity by bringing two such DNA-binding domains

79  together.

80

81  TALE-heterodimerization is mediated by class-specific homology domains located on the N-

82  terminal side adjacent to the homeodomain [12,13]. Animal MEIS and plant KNOX class

83  proteins share readily identifiable homology in their heterodimerization domain, leading to

84  the proposal of an ancestral TALE class named MEINOX [12]. In contrast, their partner

85  classes -- PBC and BELL -- exhibit no apparent homology in their heterodimerization

86  domains. Short shared sequence motifs and common secondary structures have been found

87  within the heterodimerization domains between MEINOX and PBC or BELL [14,15], but their

88  extent of the homology requires adequate taxon sampling to recover ancestral relationships.

89

90  An ancestral functions of TALE-TALE heterodimerization was revealed in studies of the

91  unicellular green alga *Chlamydomonas reinhardtii*: the KNOX ortholog GSM1 and a second

92  TALE protein GSP1 form heterodimers immediately after the fusion of sexual gametes, and

3

93    these drive the haploid-to-diploid transition by activating >200 diploid-specific genes and

94    inactivating >100 haploid-specific genes [10,16,17]. In subsequent studies, plant-type TALE-

95    TALE heterodimers between KNOX and BELL were shown to be required for the haploid-to-

96    diploid transition of the moss *Physcomitrella patens* [18,19]. Given the conserved role of

97    TALE heterodimerization as a developmental switch in the sexual life cycle of the plant

98    lineage, understanding its origins and diversification promises to shed light on the evolution

99    of developmental mechanisms during eukaryotic radiation and the emergence of land plants.

100

101    To delineate the ancestry of plant-type TALE heterodimerization, we performed a

102    phylogenetic and bioinformatics analysis of TALE TFs in the three algal radiations of the

103    Archaeplastida supergroup, the descendants of a single endosymbiosis event > one billion

104    years ago [20,21]. Our analysis showed that the TALEs were already diversified into two

105    groups at the origin of Archaeplastida, one sharing KNOX-homology and the other sharing

106    PBC-homology. Together with our protein-protein interaction data, we propose that all TALE

107    classes participate in heterodimerization networks via the KNOX- and PBC-homology

108    domains between the two ancestral groups.

109

110    **Results**

111

112    **TALEs in Archaeplastida are divided into two groups, KNOX and non-KNOX**

113    The Archaeplastida consists of three monophyletic phyla [22,23] (Fig 1). 1) Viridiplantae

114    include two divisions, Chlorophyta -- chlorophytes and prasinophytes (a paraphyletic group

115    of seven lineages [24]) -- and Streptophyta -- charophyte algae and land plants [25]. 2)

116    Rhodophyta (red algae) include diverse unicellular and multicellular organisms that diverge

117    into four major lineages [26] (S1 Spreadsheet). 3) Glaucophyta members include only four

118    cultured genera and possess plastids that carry ancestral features of the cyanobacterial

119    symbiont that gave rise to photosynthetic organelles in eukaryotes [27].

120

121    To collect all the available homeobox protein sequences, we performed BLAST and Pfam-

122    motif searches against non-plant genomes and transcriptome assemblies throughout the

123    Archaeplastida (S1 Spreadsheet), identifying 327 proteins from 55 species as the

124    Archaeplastida homeobox collection (29 genomes and 18 transcriptomes; S2 Spreadsheet).

125    Of these, 102 possessed the defining feature of TALE proteins, a three-amino-acid insertion

126    between aa positions 23-24 in the homeodomain [28]. At least two TALE genes were

127    detected in most genomes except five genomes in the Trebouxiophyceae class of the

128    Chlorophyta (S1 Spreadsheet; see S1A Notes for further discussion of the absence of

129    TALEs in Trebouxiophyceae).

130

131    The collected TALE sequences were then classified by their homeodomain features using a

132    phylogenetic approach, with TALEs from animals, plants, and early-diverging eukaryotes

133    (Amoebozoa and Excavata) as outgroups (S1 Fig). The resultant TALE homeodomain

134    phylogeny distinguished two groups in all three phyla of Archaeplastida (Fig 2). 1) The

135    KNOX-group as a well-supported clade displayed a phylum-specific cladogram: two

136    Glaucophyta sequences at the base (as KNOX-Glauco) were separate from the next clade,

137    which combines Rhodophyta sequences (as KNOX-Red1) and a Viridiplantae-specific clade

138    with strong support (92/90/1.00). 2) The non-KNOX group, including the BELL and GSP1

139    homologs, contained clades of mixed taxonomic affiliations. These analyses showed that the

140    TALE proteins had already diverged into two groups before the evolution of the

141    Archaeplastida and that the KNOX-group is highly conserved throughout Archaeplastida.

142

143    **KNOX group sequences share the same heterodimerization domains throughout**

144    **Archaeplastida**

145    The next question was whether the plant KNOX class originated prior to the Viridiplantae

146    phylum. The plant KNOX proteins and the Chlorophyta GSM1 possess KNOX-homology

147    sequences, consisting of KN-A, KN-B and ELK domains, required for their

148    heterodimerization with other TALE proteins [10]; therefore, the presence of the KNOX

149    homology sequences suggest functional homology to the plant KNOX class.  To collect

150    homology domains without prior information, we performed ad-hoc homology domain

151    searches among the KNOX group sequences. Using the identified homology domains as

152    anchors, we carefully curated an alignment of the KNOX-group sequences combined with

153    any other TALE sequences with a KNOX-homology, (S2 Fig). From this KNOX alignment,

154    we defined KNOX-homologs as having amino acid similarity scores >50% for at least two of

155    the three domains comprising the KNOX-homology region (S3 Spreadsheet for calculated

156    domain homology). Using this criterion, all KNOX group sequences (excluding partial

157    sequences) possessed the KNOX homology (Fig 2, marked by red dots following their IDs),

158    indicating that the KNOX-homolog already existed before the evolution of eukaryotic

159    photosynthesis as represented by the Archaeplastida.

160

161    In addition to the KNOX-homology, the same search also revealed two novel domains at the

162    C-terminus of the homeodomain (S2 Fig): the first (KN-C1) was shared among the

163    Chlorophyta sequences, and the second (KN-C2) was shared among a group of KNOX

164    homologs in a clade outside the KNOX-group (KNOX-Red2).

165

166    **KNOX classes diverged independently among the algal phyla**

5

167    In Viridiplantae, we found a single KNOX homolog in most Chlorophyta species, whereas

168    KNOX1 and KNOX2 divergence was evident in the Streptophyta division, including the

169    charophyte *Klebsormidium flaccidum* and land plants (Fig 2). The newly discovered KN-C1

170    domain was specific to the Chlorophyta KNOX sequences and found in all but one species

171    (*Pyramimonas amylifera*). The absence of similarity between KN-C1 and the C-terminal

172    extensions of KNOX1/KNOX2 sequences suggests independent, lineage-specific KNOX

173    evolution in the Chlorophyta and Streptophyta (S2 Fig). We, therefore, refer to the

174    Chlorophyta KNOX classes as KNOX-Chloro in contrast to the KNOX1 and KNOX2 classes

175    in the Streptophyta.

176

177    The KNOX homologs in the Rhodophyta were divided into two classes: a paraphyletic group

178    close to the KNOX-Chloro clade, named KNOX-Red1, and a second group near the PBX-

179    Outgroup, named KNOX-Red2. KNOX-Red1 lacked a KN-A, whereas KNOX-Red2 lacked

180    an ELK and shared a KN-C2 domain (S2 Fig). We consider KNOX-Red1 as the ancestral

181    type, since the KNOX-Red1 sequences were found in all examined Rhodophyta taxa,

182    whereas the KNOX-Red2 sequences were restricted to two taxonomic classes

183    (Cyanidiophyceae and Florideophyceae). Interestingly, the KNOX-Red2 clade included two

184    green algal sequences, with strong statistical support (89/89/0.97; Fig 2); these possessed a

185    KN-C2 domain, suggesting their ancestry within the KNOX-Red2 class (S2 Fig; see S1B

186    Notes for further discussion about their possible origin via horizontal gene transfer).

187

188    Available TALE sequences were limited for the Glaucophyta. We found a single KNOX

189    homolog in two species, which possessed KN-A and KN-B domains but lacked an ELK

190    domain. We termed these KNOX-Glauco.

191

192    **Non-KNOX group TALEs possess animal type PBC-homology domain, suggesting a**

193    **shared ancestry between Archaeplastida and Metazoa**

194    Following the identification of KNOX homologs, the non-KNOX group in the Archaeplastida

195    was redefined as lacking KN-A and KN-B domains. Further classification of the non-KNOX

196    group was challenging due to its highly divergent homeodomain sequences. However, we

197    noticed that the number of non-KNOX genes per species was largely invariable: one in most

198    Rhodophyta and Glaucophyta genomes and two in the majority of Chlorophyta genomes,

199    suggesting their conservation within each radiation.

200

201    Our ad-hoc homology search provided critical information for non-KNOX classification,

202    identifying a homology domain shared among all Glaucophyta and Rhodophyta non-KNOX

203    sequences (Fig 3A and 3B). Since this domain showed a similarity to the second half of the

204    animal PBC-B domain (Pfam ID: PF03792) known as heterodimerization domain [12], we

205    named this domain PBL (PBC-B Like). Accordingly, we classified all the non-KNOX TALEs

206    in Glaucophyta and Rhodophyta as a single PBC-related homeobox class, PBX-Glauco or

207    PBX-Red. PBX-Glauco sequences also possessed the MEIONOX motif, conserved in the

208    animal PBC-B domain, indicating common ancestry of PBC-B and PBL domains (Fig 3A).

209

210    **GSP1 shares distant PBC-homology together with other non-KNOX group sequences**

211    **in Viridiplantae**

212    A remaining question was the evolution of the Chlorophyta non-KNOX sequences that

213    apparently lacked a PBC-homology. To uncover even a distant homology, we compared the

214    newly defined PBL domains with the Chlorophyta sequences by BLAST (cut-off E-value of

215    1E-1) and multiple sequence alignments. This query collected three prasinophyte and one

216    charophyte TALE sequences that possessed a MEINOX motif and a putative PBL-domain;

217    however, they showed very low sequence identity among themselves (Fig 3C). Further

218    query utilizing these four sequences identified 11 additional non-KNOX sequences. Nine of

219    these were made into two alignments, one including GSP1 homologs and the other

220    combining most prasinophyte sequences (S3A and S3B Fig). The two remaining sequences

221    (Picocystis_salinarum_04995 and Klebsormidium_flaccidum_00021_0250) showed a

222    homology to a PBX-Red sequence of *Chondrus cruentum* (ID:41034) in a ~ 200 aa-long

223    extension beyond the PBL domain, suggesting their PBX-Red ancestry (another potential

224    case of horizontal transfers; S4 Fig). All the Chlorophyta non-KNOX sequences that carry

225    the PBL-homology domains were classified as GLX (GSP1-like homeobox) in recognition of

226    the GSP1 protein of *Chlamydomonas* as the first characterized member of this class [29].

227

228    **Two non-KNOX paralogs of Chlorophyta heterodimerize with the KNOX homologs.**

229    Even with our sensitive iterative homology search, we could not identify a PBC/PBL-

230    homology in about half of the Chlorophyta non-KNOX sequences. Since most Chlorophyta

231    genomes possess one GLX homolog and one non-KNOX sequence without the PBL-

232    homology domain, we refer the latter collectively to Class-B (S5 Fig). Exceptions were found

233    in one prasinophyte clade (class Mamiellophyceae), whose six high-quality genomes all

234    contain two non-KNOX sequences lacking the PBL-homology. Nonetheless, these non-

235    KNOX sequences formed two groups, one more conserved and the other less conserved,

236    referred to the Mam-A and Mam-B classes, respectively (S6 and S7 Fig). Considering the

237    reductive genome evolution of the Mamiellophyceae [30], the conserved Mam-A class may

238    be derived from an ancestral GLX class.

239

7

240    Two divergent non-KNOX classes in Chlorophyta led to a critical question about their dyadic

241    networks. Previously studies had shown that TALE heterodimers required interaction

242    between MEIS and PBC domains in animals and between KNOX and PBL domains in

243    *Chlamydomonas* [6,10]. It was, therefore, predicted that all Glaucophyta and Rhodophyta

244    TALEs form heterodimers via their KNOX- and PBL-homology domains. On the other hand,

245    it remained to be tested whether the Chlorophyta TALEs lacking a PBL-domain can form

246    heterodimers with other TALEs.

247

248    To characterize interaction network of TALE class proteins in Chlorophyta, we selected three

249    prasinophyte species for protein-protein interaction assays: two species containing Mam-A

250    and Mam-B genes (*Micromonas commoda* and *Ostreococcus tauri*), and another species

251    (*Picocystis salinarum*), whose transcriptome contained one GLX and one Class-B sequence.

252    In all three species, we found that KNOX homologs interacted with all examined non-KNOX

253    proteins in Mam-A, Mam-B, Class-B, and GLX class (Fig 4A-4C). No interaction was

254    observed between the two non-KNOX proteins in any of the three species (Fig 4A-4C).

255    Similar to the GLX-KNOX heterodimerization, Mam-A and Mam-B also required additional

256    domains outside the homeodomain for their heterodimerization with the KNOX homologs

257    (S8 Fig). These results showed that the all divergent non-KNOX TALEs maintained their

258    original activity to form heterodimers with the KNOX homologs. Observed interacting

259    network among the TALE sequences is summarized in S9A Fig.

260

261    **TALE heterodimerization evolved early in eukaryotic history**

262    Our discovery of the PBC-homology in Archaeplastida suggests common ancestry of the

263    heterodimerizing TALES between Metazoa and Archaeplastida. It also predicted that other

264    eukaryotic lineages might possess TALEs with the PBC-homology. Outside animals, the

265    Pfam database contains only two PBC-B domain-harboring sequences, one from a

266    Cryptophyta species (*Guillardia theta,* ID_137502) and the other from an Amoebozoa

267    species (*Acanthamoeba castillian,* ID:XP_004342337)[31]. We further examined the

268    Excavata group, near to the posited root of eukaryotic phylogeny [22]. A search of two

269    genomes (*Naegleria gruberi* and *Bodo saltans*) collected 12 TALE homeobox sequences in

270    *N.gruberi*, and none in *B.saltans*, of which we found one with a PBC-homology domain

271    (ID:78561, Fig 3A) and one with a MEIS/KNOX-homology (ID:79931, S2 Fig). Our data

272    suggest that the heterodimerization domains -- the PBC-homology and MEIS/KNOX-

273    homology -- originated early in eukaryotic evolution and persisted throughout the major

274    eukaryotic radiations.

275

**276 Intron-retention supports the parallel evolution of the heterodimeric TALE classes**

**277 during eukaryotic radiations**

278 The ubiquitous presence of dyadic TALEs raised next question: Are all the dyadic TALEs

279 reported in this study the descendants of a single ancestral dyad, or do they result from

280 lineage-specific evolution from a single prototypical TALE (proto-TALE) that does not

281 engage in heterodimerization. To probe deep ancestry, we examined intron-retention, this

282 being regarded as a long-preserved character and less prone to occur by homoplasy (a

283 character displayed by a set of species but not present in their common ancestor) [32].  Five

284 intron positions were shared by at least two TALE classes, of which the 44/45 and 48[2/3]

285 introns qualified as the most ancestral since they were found throughout the Archaeplastida

286 and Metazoa (S10 Fig).

287

288 The 44/45 and 48[2/3] introns showed an intriguing exclusive distribution between the two

289 dyadic partners of each phylum: one possesses the 44/45 and the other the 48[2/3] intron

290 (S10 Fig). This mutually exclusive pattern suggested that two TALE genes with distinct

291 intron positions existed at the onset of the eukaryotic radiation. We consider the 44/45 intron

292 position as the most ancestral, given that it was conserved in most non-TALE homeobox

293 genes [12]. In this regard, we speculate that acquisition of the 48[2/3], and loss of the 44/45

294 intron, accompanied an early event wherein the proto-TALE with the 44/45 intron was

295 duplicated to generate a second TALE with the 48[2/3] intron. Since the two intron positions

296 were found within both the MEIS/KNOX and PBC/PBX/GLX groups, we propose that the

297 duplicated TALEs arose early and diversified to establish lineage-specific heterodimeric

298 configurations during eukaryotic radiations.

299

300 Given that the heterodimeric TALEs evolved in a lineage-specific manner, we asked what

301 the proto-TALE looked like at the time it underwent duplication. The following observations

302 suggest that the proto-TALE was a homodimerizing protein. First, the PBC-homology

303 domains of PBX/GLX class proteins identified in the Archaeplastida includes the MEINOX-

304 motif that was originally defined for its similarity to the MEIS/KNOX-homology domains (Fig

305 3) [14]. Second, PBX-Glauco sequences possess the ELK-homology within their PBL

306 domain (Fig 3), which align well to the ELK domains of KNOX class sequences in

307 Viridiplantae (S11 Fig). Therefore, the MEINOX-motif and ELK-homology across the

308 heterodimerizing KNOX and PBX groups supported the common origin of heterodimerizing

309 TALE groups from a single TALE by duplication followed by subfunctionalization.

310

**311 Discussion**

312

9

313 **TALE endowment in Archaeplastida**

314 Our study shows that all three Archaeplastida phyla possess TALEs, diverged into two

315 groups with distinct heterodimerization domains, the KNOX group with KN-A/KN-B domains

316 and the PBX (or GLX) group with PBL domains. The similarity between the KNOX/PBX and

317 the animal MEIS/PBC dyads led us to identify homologous heterodimerization domains in

318 the TALEs of other eukaryotic lineages including Excavata. Based on our findings, we

319 hypothesize that the TALE heterodimerization arose very early in eukaryotic evolution.

320

321 During > 1 BY of Archaeplastida history, TALE TF networks have undergone three

322 duplication events compared to the simple dyadic TALEs in Glaucophyta. In Viridiplantae,

323 the KNOX class persists as a single member throughout the mostly unicellular Chlorophyta,

324 whereas it duplicated into KNOX1 and KNOX2 in the multicellular Streptophyta [33]. In

325 Rhodophyta, two KNOX classes, KNOX-Red1 and KNOX-Red2 differ in KN-A and KN-B

326 domains, suggesting sub-functionalization. The third duplication event occurred in the non-

327 KNOX group of the Chlorophyta, whose sequences then underwent rapid divergence in their

328 homeodomain and heterodimerization domains, rendering their classification trickier than

329 other classes. Despite this divergence, proteins in one of the two radiations (Class-B and

330 Mam-B) were found to heterodimerize with KNOX homologs, suggesting that these non-

331 KNOX members serve as regulators of KNOX/GLX heterodimers. We summarize our finding

332 in Fig 1, S10B Fig.

333

334 **Is the plant BELL class homologous to the Chlorophyta GLX class?**

335 The BELL class is the only non-KNOX class in land plants, sharing a POX (Pre-homeobox)

336 domain (PF07526) [13] and lacking an identifiable PBL domain. The *K. flaccidum* genome,

337 the only genome available in the charophyte lineage from which land plant emerged,

338 contained three non-KNOX sequences, all possessing a PBL domain (Fig 3, S3,S4 Fig).

339 Therefore, the lack of PBL-homology in the plant BELL class appears to be due to

340 divergence or domain loss from an old charophyte class that had PBL-homology. We found

341 an intron at the 24[2/3] homeodomain position of a *K. flaccidum* GLX homolog, which was

342 previously identified as being specific to the plant BELL class (S8A Fig) [12], suggesting that

343 the plant BELL class evolved from an ancestral GLX gene. More taxon sampling in

344 charophytes is needed to confirm this inference.

345

346 **What would have been the critical drivers of TALE heterodimerization networks**

347 **emerging from ancestral homodimers?**

348 We found two conserved intron positions and shared sequence motifs between the KNOX-

349 and PBX-groups, generating our hypothesis that a proto-TALE protein initially engaged in

10

350   homodimerization and then duplicated and diversified into two heterodimerizing classes (Fig

351   1, S9A Fig). Heterodimerization-dependent subcellular localization [10,34], coupled with

352   numerous combinations of distinct DNA-binding modules that fine-tune target specificity,

353   then generated customized transcription-on switches.

354

355   During sexual development, it is critical to accurately detect the fusion of two cells before

356   initiating diploid development and to make sure that the mating combines correct partner

357   gametes. TF heterodimerization can implement both steps if one TF partner is contributed

358   by each gamete. In fact, TALE heterodimerization plays a central role as a developmental

359   switch for the haploid-to-diploid transition in green algae and land plants [10,19]. A similar

360   haploid-to-diploid transition triggered by TF heterodimerization has recently been

361   documented in *Dictyostelium* [35] and is well described in Basidiomycete fungi that utilize

362   non-TALE homeobox proteins such as bW and bE [36,37].

363

364   Discovery of new prokaryotic life forms, especially in the Archaea domain, suggests that

365   multiple symbiotic mergers of different life forms evolved into the proto-eukaryotes, possibly

366   first as a symbiotic community, which then evolved into the last eukaryotic common

367   ancestors (LECA) that rapidly diverged into the eukaryotic supergroups [38-40]. This

368   eukaryogenesis model predicts that the proto-eukaryotes → LECA transition required the

369   faithful transmission of traits between progenitor cells and their progeny to evolve as

370   individual lineages by Darwinian selection. Under this hypothesis, we anticipate that the

371   generation of the LECA may have been driven by the sexual mechanisms that distinguish a

372   cellular merger between the common descendants from a merger between unrelated

373   community members. Our proposal for the evolution of heterodimeric TALEs from the

374   homodimeric proto-TALE may provide one of the necessary mechanisms for the first sexual

375   mode of reproduction that might have driven the generation of the LECA from its proto-

376   eukaryotic ancestors.

377

378   **Does expansion of heterodimerizing TALE TFs relate to the emergence of**

379   **multicellular complexity?**

380   Plant studies have shown that the duplicated KNOX classes serve distinct functions: the

381   plant KNOX1 class regulates the differentiation of an undifferentiated cell mass into spores

382   in mosses or leafy organs in vascular plants, and the plant KNOX2 class regulates the

383   transition from haploid gametophytes to diploid sporophytes in mosses and controls

384   secondary cell wall development in vascular plants [18,41-43]. We propose that the

385   duplicated TALE heterodimers in the Streptophyta allowed independent regulation of cellular

386   differentiation and life cycle transitions, priming the emergence of land plants by expanding

11

387    the diploid phase of their life cycle from a dormant zygospore to a multicellular individual

388    bearing many meiotic spores. The repertoire of TALE heterodimers continued to expand

389    during land plant evolution, serving all the major organ differentiation programs in the diploid

390    phase of their life cycle.

391

392    Can a similar expansion of TALE heterodimers be found during Metazoan evolution? Our

393    search for TALE TFs in unicellular relatives of the Metazoa -- Spingoeca and Monosiga –

394    revealed a simple configuration with one MEIS- and one PBC-like TALE (S12, S13 Fig),

395    whereas at the Metazoan base one finds at least three MEIS-related classes and two PBC-

396    related classes [44]. These findings suggest the occurrence of a similar expansion of a

397    founding dyad during Metazoan evolution. Therefore, in both plants and animals, the TALE

398    TF network seems to be redeployed for complex multicellularity, departing from its posited

399    original function in sexual development.

400

401    Our results suggest that TALE TF networks represent early-evolving developmental

402    mechanisms. That said, the emergence of complex multicellularity doubtless required more

403    than TF networks. TF-based developmental cues need to be propagated via chromatin-level

404    regulatory mechanisms that establish the cellular memory during embryo development. The

405    extent to which chromatin-level regulatory mechanisms are involved in the development of

406    unicellular organisms is a critical question in elucidating the origins of complex

407    multicellularity.

408

409    **Materials and methods**

410    **Strains and culture conditions**

411    Axenic *Micromonas commoda* (RCC299) and *Ostreococcus tauri* (OTH95) were maintained

412    in Keller medium [45] in artificial seawater at room temperature. One hundred mL of a 14-

413    day-old culture was harvested for genomic DNA extraction. *Picocystis salinarum*

414    (CCMP1897) was obtained from the National Center for Marine Algae and Microbiota

415    (NCMA), maintained in L1 medium [46] in artificial sea water, and plated on 1.5% Bactoagar-

416    containing media for single-colony isolation. Genomic DNA of *P. salinarum* was then

417    obtained from a culture derived from one colony.

418

419    **Phylogenetic analysis and classification of homeobox genes**

420    Archaeplastida algal TALE homeodomains were collected from the available genomes and

421    transcriptomes listed in S1 Spreadsheet. Details of how TALE sequence was collected is

422    provided in S1A Methods. After excluding nearly identical sequences, a total of 96

423    sequences together with 18 reference TALE sequences were made into the final

12

424    homeodomain alignment with 70 unambiguously aligned positions with eight gapped and

425    one constant sites. Details of phylogenetic reconstruction is provided in S1B Methods.

426

427    **Bioinformatics analysis**

428    The entire TALE collection was divided into multiple groups representing major clades in the

429    homeodomain tree. Each group was individually analyzed by running MEME4.12 in the

430    motif-discovery mode with default option collecting up to 10 motifs at http://meme-suite.org/

431    [47]. The search provided multiple non-overlapping motifs, many of which were combined

432    according to previously identified domains such as bipartite KN-A/KN-B, ELK, and HD [14]

433    and independent domain searches against the INTERPRO database

434    (http://www.ebi.ac.uk/interpro/) [48]. All the collected TALE-associated homology domains

435    were aligned to generate HMM motifs, which we used to test if these homology domains are

436    specific to the TALE sequences. All the homology domain information was used to locate

437    any error in gene predictions, and gene models were updated if necessary (Details of the

438    gene model curation is provided in S1C Methods).

439

440    **Intron comparison**

441    Introns within the homeodomain were collected and labeled as site numbers of the

442    homeodomain (1-63). If an intron is between two codons it is denoted N/N+1, where N is the

443    last amino acid site number of the preceding exon; introns within a codon are denoted

444    N[n/n+1], where n is one or two for the codon nucleotide position relative to the splice-sites.

445

446    **Yeast-two-hybrid analysis**

447    *M. commoda* (affixed with Micco), *O. tauri* (affixed with Ostta), and *P. salinarum* (affixed with

448    Picsa) TALE protein coding sequences were cloned by PCR using primers designed herein

449    (S4 Spreadsheet) from genomic DNAs prepared by the phenol/chloroform extraction and

450    ethanol precipitation method. Micco_62153 and Picsa_04684 contained a single intron,

451    whereas all the other nine genes lacked an intron in the entire open reading frame. For

452    cloning of Micco_62153, we synthesized the middle fragment lacking the intron and ligated

453    them via *Xho*I and *Cla*I sites. For cloning details, see S1D Methods.

454

455    **Supporting Information**

456    **S1 Fig. Alignment of TALE homeodomain sequences of the Archae-algal collection.**

457    The 106 sequences were made into an alignment after excluding 20 near identical

458    sequences to reduce redundancy. Animal/amoeba/haptophyte outgroup sequences are

459    included as they share homology with Archaeplastida TALEs outside the homeodomain. The

460    three bars above the sequence numbers show predicted alpha helices. Discarded insertions

461    are noted in red arrowheads.

462    **S2 Fig. Homology domain alignment of KNOX homogogs.** MEIS class outgroup

463    sequences are included at the bottom. Class label is on the left. KN-A, KN-B, ELK,

464    HOMEOBOX, KN-C1, and KN-C2 domains are labeled on the top. Class groups are labeled

465    by colored bars on the left next to the gene names. Yellow, light green, and green shades in

466    sequences show more than 60%, 80%, or 100% similarity in each column. Gaps between

467    KN-A and KN-B and between KN-B and ELK have been eliminated.

468    **S3 Fig. GLX class is defined by PBL-Chloro domain.** (A-C) Three alignments are

469    adjusted with inserting gaps for direct comparison among different PBL-Chloro domains. (A)

470    GLX-Chloro class members. (B) GLX-Basal class members. (C) Three Viridiplantae

471    sequences with strong MEINOX homology domain. PBC-homology is shared among the

472    Chlorophyta non-KNOX sequences.

473    **S4 Fig. Extensive homology of Picsa_04995 and Klefl_00021_0250 to Chocr_41034**

474    **indicates their classification as PBX-Red.** MEINOX-homology and PBL-Red domains are

475    indicated by red bars below the alignment.

476    **S5 Fig. Alignment of Class-B TALE proteins in volvocales.** Short motifs are conserved

477    among all members in this class over the entire length of the sequence.

478    **S6 Fig. Alignment of Mam-A TALE proteins in mamiellophyceae.** Short motifs (Box1-4)

479    are conserved among all members in this class over the entire length of the sequence. Red

480    reverse triangle at 548-549 shows the truncation position of Micco_Mam-A-tr used in Yeast-

481    two-hybrid analysis.

482    **S7 Fig. Alignment of Mam-B TALE proteins in mamiellophyceae.** A conserved motif is

483    found between 180-197 amino acids in the alignment. Red reverse triangle at 100-101

484    shows the truncation position of Ostta_Mam-B-tr used in Yeast-two-hybrid analysis.

485    Homology is restricted to a single homology-A domain ouside the homeodomain.

486    **S8 Fig. Full-length proteins are necessary for mamiellophyceae non-KNOX TALE**

487    **proteins to form heterodimers.** Left and Right: Yeast-two-hybrid assays on Ade-/His-/Leu-

488    /Trp- medium. The construct information for the prey conjugated with the GAL4 DNA-binding

489    domain and for the bait conjugated with the GAL4 transcriptional activation domain is given

490    in the table below.

491    **S9 Fig. TALE interaction network defined by this study using yeast-two-hybrid**

492    **assays.** (A) Summary diagram for the TALE interaction network. (B) Yeast-two-hybrid

493    assays for the cross-species interaction of TALE proteins. Only one of the possible

494    reciprocal combinations of GAL4 domain conjugations is provided for simplicity. Large X

495    indicates no yeast in the sector. -LTHA: Leu-/Trp-/His-/Ade- medium; -LT: Leu-/Trp- medium.

496    **S10 Fig. Intron-retention pattern suggests parallel evolution of KNOX and non-KNOX**

497    **group classes from common duplicated TALE ancestors.** (A) Intron locations collected

14

498    from 12 TALE classes are shown with arrows. Half arrows indicate cases where not all the

499    class members share the position. White arrows indicate shared positions in at least two

500    different classes, and black arrows indicate class-specific positions. The numbers above the

501    consensus sequence show 60 amino acid positions; the three-amino-acid extension is

502    denoted as 'abc.' Row color depicts two alternative domain configurations: purple for

503    MEIS/KNOX types, and navy for PBX/GLX types. Class names are colored according to

504    their phylogenetic groups: green for Viridiplantae, red for Rhodophyta, blue for Glaucophyta,

505    and black for outgroups. The numbers following the class names show how many genes

506    provided the intron information. Of the shared positions, purple triangles on the top mark

507    those shared between MEIS/KNOX and PBX/GLX classes, blue triangles mark those shared

508    between GLX and BELL classes, and red triangles mark those shared between KNOX

509    classes. A notable exception is the KNOX-Red1 class, for which three Rhodophyta clades

510    show different intron locations (44/45, 48[2/3] or 53[2/3]), indicating that the 44/45 intron

511    position can indeed be displaced to 48[2/3] or elsewhere, albeit infrequently. The unique

512    46/47 intron in the PBX-Glauco (Cyapa_20927) would presumably have resulted from a

513    similar displacement in intron position. (B) Distribution of conserved introns among the TALE

514    homeobox classes. Identified TALE classes are mapped on the Arachaeplastida phylogeny.

515    The 44/45 intron is marked by blue outline and the 48[2/3] intron is marked by red outline.

516    Underlines of the class names indicate the presence of a PBC-homology domain. The

517    Archaeplastida phylogeny is modified from figure 1 of Jackson et al. (2015).

518    **S11 Fig. ELK-domain alignment.**

519    **S12 Fig. Identification of MEIS homologs in choanoflagellates.**

520    **S13 Fig. Identification of PBX homologs in choanoflagellates.**

521

522    **S1 Spreadsheet Genomic resources used in this study.** A total of 374 homeobox protein

523    sequences are compiled for this analysis, of which 113 TALE protein sequences are

524    collected. The number of total homeobox proteins and TALE superclass were estimated

525    largely from our homeodomain search described in the materials and methods section.

526    Under Genome annotation, 'Draft' indicates a genome without annotation, 'Trans' indicates a

527    transcriptome assembly.

528    **S2 Spreadsheet Archaeplastidal homeobox collection of TALE protein analyzed in**

529    **this study.** For outgroups, only TALE members that are analyzed in this study are included.

530    **S3 Spreadsheet KNOX domain homology among KNOX classes**

531    **S4 Spreadsheet Primers used in this study**

532    **S5 Spreadsheet Yeast-two-hybrid constructs used in this study**

533    **S6 Spreadsheet Homeobox profile in Trebouxiophyceae**

534    **S1 Methods**

15

535 A. Collecting TALE homeobox protein sequences. B. Phylogenetic reconstruction. C.

536 Homology motif/domain search. D. Intron comparison. E. Cloning of Yeast-two-hybrid

537 constructs.

538 **S1 Notes**

539 A. Lack of TALE TFs in Trebouxiophyceae. B. Horizontal transfer may explain the presence

540 of Rhodophyta TALE heterodimers in *Picocystis and Klebsormidium* of Viridiplantae.

541

## Acknowledgements

550

## Author Contributions

552 **Conceptualization**: Sunjoo Joo, Alexandra Z. Worden, Jae-Hyeok Lee

553 **Data curation**: Sunjoo Joo, Ming Hsiu Wang, Gary Lui, Jenny Lee, Sebastian Sudek, Jae-

554 Hyeok Lee

555 **Formal analysis**: Sunjoo Joo, Ming Hsiu Wang, Jae-Hyeok Lee

556 **Funding acquisition**: Sunjoo Joo, Alexandra Z. Worden, Jae-Hyeok Lee

557 **Investigation**: Sunjoo Joo, Ming Hsiu Wang, Jenny Lee, Andrew Barnas, Jae-Hyeok Lee

558 **Project administration**: Alexandra Z. Worden, Jae-Hyeok Lee

559 **Resources**: Sunjoo Joo, Ming Hsiu Wang, Eunsoo Kim, Sebastian Sudek

560 **Supervision**: Sunjoo Joo, Jae-Hyeok Lee

561 **Writing – original draft**: Sunjoo Joo, Ming Hsiu Wang, Jae-Hyeok Lee

562 **Writing – review & editing**: Sunjoo Joo, Eunsoo Kim, Alexandra Z. Worden, Jae-Hyeok

563 Lee

564

## References

566 1. Billeter M, Qian YQ, Otting G, Muller M, Gehring W, Wuthrich K. Determination of

567 the nuclear magnetic resonance solution structure of an Antennapedia

568 homeodomain-DNA complex. J Mol Biol. 1993;234: 1084–1093.

569 doi:10.1006/jmbi.1993.1661

570    2.      Azpiazu N, Morata G. Functional and regulatory interactions between Hox and
571            extradenticle genes. Genes & Development. 1998;12: 261–273.

572    3.      Hudry B, Thomas-Chollier M, Volovik Y, Duffraisse M, Dard A, Frank D, et al.
573            Molecular insights into the origin of the Hox-TALE patterning system. eLife. 2014;3:
574            e01939. doi:10.7554/eLife.01939

575    4.      Hake S, Smith HMS, Holtan H, Magnani E, Mele G, Ramirez J. The role of KNOX
576            genes in plant development. Annu Rev Cell Dev Biol. Annual Reviews; 2004;20:
577            125–151. doi:10.1146/annurev.cellbio.20.031803.093824

578    5.      Hay A, Tsiantis M. KNOX genes: versatile regulators of plant development and
579            diversity. Development. 2010;137: 3153–3165. doi:10.1242/dev.030049

580    6.      Berthelsen J, Kilstrup-Nielsen C, Blasi F, Mavilio F, Zappavigna V. The sub cellular
581            localization of PBX1 and EXD proteins depends on nuclear import and export
582            signals and is modulated by association with PREP1 and HTH. Genes &
583            Development. 1999;13: 946–953.

584    7.      Stevens KE, Mann RS. A balance between two nuclear localization sequences and
585            a nuclear export sequence governs extradenticle subcellular localization. Genetics.
586            Genetics; 2007;175: 1625–1636. doi:10.1534/genetics.106.066449

587    8.      Bhatt AM, Etchells JP, Canales C, Lagodienko A, Dickinson H. VAAMANA--a
588            BEL1-like homeodomain protein, interacts with KNOX proteins BP and STM and
589            regulates inflorescence stem growth in Arabidopsis. Gene. 2004;328: 103–111.
590            doi:10.1016/j.gene.2003.12.033

591    9.      Smith HMS, Hake S. The Interaction of Two Homeobox Genes,
592            BREVIPEDICELLUS and PENNYWISE, Regulates Internode Patterning in the
593            Arabidopsis Inflorescence. Plant Cell. 2003;15: 1717–1727.

594    10.     Lee J-H, Lin H, Joo S, Goodenough U. Early sexual origins of homeoprotein
595            heterodimerization and evolution of the plant KNOX/BELL family. Cell. 2008;133:
596            829–840. doi:10.1016/j.cell.2008.04.028

597    11.     Noyes MB, Christensen RG, Wakabayashi A, Stormo GD, Brodsky MH, Wolfe SA.
598            Analysis of homeodomain specificities allows the family-wide prediction of preferred
599            recognition sites. Cell. 2008;133: 1277–1289. doi:10.1016/j.cell.2008.05.023

600    12.     Bürglin TR. Analysis of TALE superclass homeobox genes (MEIS, PBC, KNOX,
601            Iroquois, TGIF) reveals a novel domain conserved between plants and animals.
602            Nucleic Acids Research. 1997;25: 4173–4180.

603    13.     Bellaoui M, Pidkowich MS, Samach A, Kushalappa K, Kohalmi SE, Modrusan Z, et
604            al. The Arabidopsis BELL1 and KNOX TALE homeodomain proteins interact
605            through a domain conserved between plants and animals. Plant Cell. 2001;13:
606            2455–2470. doi:10.1105/tpc.010161

607  14.  Bürglin TR. The PBC domain contains a MEINOX domain: coevolution of Hox and
608       TALE homeobox genes? Dev Genes Evol. 1998;208: 113–116.

609  15.  Mukherjee K, Brocchieri L, Bürglin TR. A comprehensive classification and
610       evolutionary analysis of plant homeobox genes. Molecular Biology and Evolution.
611       2009;26: 2775–2794. doi:10.1093/molbev/msp201

612  16.  Nishimura Y, Shikanai T, Nakamura S, Kawai-Yamada M, Uchimiya H. Gsp1
613       triggers the sexual developmental program including inheritance of chloroplast
614       DNA and mitochondrial DNA in Chlamydomonas reinhardtii. Plant Cell. 2012;24:
615       2401–2414. doi:10.1105/tpc.112.097865

616  17.  Joo S, Nishimura Y, Cronmiller E, Hong RH, Kariyawasam T, Wang MH, et al.
617       Gene regulatory networks for the haploid-to-diploid transition of Chlamydomonas
618       reinhardtii. Plant Physiol. 2017;175: 314–332. doi:10.1104/pp.17.00731

619  18.  Sakakibara K, Ando S, Yip HK, Tamada Y, Hiwatashi Y, Murata T, et al. KNOX2
620       Genes Regulate the Haploid-to-Diploid Morphological Transition in Land Plants.
621       Science. 2013;339: 1067–1070. doi:10.1126/science.1230082

622  19.  Horst NA, Katz A, Pereman I, Decker EL, Ohad N, Reski R. A single homeobox
623       gene triggers phase transition, embryogenesis and asexual reproduction. Nature
624       Plants. 2016;2: 15209. doi:10.1038/nplants.2015.209

625  20.  Worden AZ, Lee JH, Mock T, Rouze P, Simmons MP, Aerts AL, et al. Green
626       Evolution and Dynamic Adaptations Revealed by Genomes of the Marine
627       Picoeukaryotes Micromonas. Science. 2009;324: 268–272.
628       doi:10.1126/science.1167222

629  21.  Archibald JM. Genomic perspectives on the birth and spread of plastids.
630       Proceedings of the National Academy of Sciences. National Academy of Sciences;
631       2015;112: 10147–10153. doi:10.1073/pnas.1421374112

632  22.  Adl SM, Simpson AGB, Lane CE, Lukeš J, Bass D, Bowser SS, et al. The revised
633       classification of eukaryotes. J Eukaryot Microbiol. 2012;59: 429–493.
634       doi:10.1111/j.1550-7408.2012.00644.x

635  23.  Worden AZ, Follows MJ, Giovannoni SJ, Wilken S, Zimmerman AE, Keeling PJ.
636       Rethinking the marine carbon cycle: Factoring in the multifarious lifestyles of
637       microbes. Science. 2015;347: 1257594–1257594. doi:10.1126/science.1257594

638  24.  Guillou, Eikrem, Massana, Romari, Vaulot. Diversity of Picoplanktonic
639       Prasinophytes Assessed by Direct Nuclear SSU rDNA Sequencing of
640       Environmental Samples and Novel Isolates Retrieved from Oceanic and Coastal
641       Marine Ecosystems. Annals of Anatomy. 2004;155: 22–22.
642       doi:10.1078/143446104774199592

643    25.    Lewis LA, McCourt RM. Green algae and the origin of land plants. American
644            Journal of Botany. 2004;91: 1535–1556. doi:10.3732/ajb.91.10.1535

645    26.    Yoon HS, Muller KM, Sheath RG, Ott FD, Bhattacharya D. Defining the major
646            lineages of red algae (RHODOPHYTA). Journal of Phycology. 2006;42: 482–492.
647            doi:10.1111/j.1529-8817.2006.00210.x

648    27.    Jackson C, Clayden S, Reyes-Prieto A. The Glaucophyta: the blue-green plants in
649            a nutshell. Acta Societatis Botanicorum Poloniae. 2015;84: 149–165.
650            doi:10.5586/asbp.2015.020

651    28.    Bertolino E, Reimund B, Wildt-Perinic D, Clerc RG. A novel homeobox protein
652            which recognizes a TGT core and functionally interferes with a retinoid-responsive
653            motif. J Biol Chem. 1995;270: 31178–31188.

654    29.    Kurvari V, Grishin NV, Snell WJ. A gamete-specific, sex-limited homeodomain
655            protein in Chlamydomonas. J Cell Biol. 1998;143: 1971–1980.

656    30.    Piganeau G, Grimsley N, Moreau H. Genome diversity in the smallest marine
657            photosynthetic eukaryotes. Res Microbiol. 2011;162: 570–577.
658            doi:10.1016/j.resmic.2011.04.005

659    31.    Clarke M, Lohan AJ, Liu B, Lagkouvardos I, Roy S, Zafar N, et al. Genome of
660            Acanthamoeba castellanii highlights extensive lateral gene transfer and early
661            evolution of tyrosine kinase signaling. Genome Biology. 2013;14: R11.
662            doi:10.1186/gb-2013-14-2-r11

663    32.    Sverdlov AV, Rogozin IB, Babenko VN, Koonin EV. Conservation versus parallel
664            gains in intron evolution. Nucleic Acids Research. 2005;33: 1741–1748.
665            doi:10.1093/nar/gki316

666    33.    Frangedakis E, Saint-Marcoux D, Moody LA, Rabbinowitsch E, Langdale JA.
667            Nonreciprocal complementation of KNOX gene function in land plants. The New
668            phytologist. 2016;341: 95–604. doi:10.1111/nph.14318

669    34.    Longobardi E, Penkov D, Mateos D, De Florian G, Torres M, Blasi F. Biochemistry
670            of the tale transcription factors PREP, MEIS, and PBX in vertebrates. Wellik D,
671            Torres M, Ros M, editors. Dev Dyn. 2014;243: 59–75. doi:10.1002/dvdy.24016

672    35.    Hedgethorne K, Eustermann S, Yang J-C, Ogden TEH, Neuhaus D, Bloomfield G.
673            Homeodomain-like DNA binding proteins control the haploid-to-diploid transition
674            inDictyostelium. Science Advances. 2017;3: e1602937.
675            doi:10.1126/sciadv.1602937

676    36.    Kües U, Asante-Owusu RN, Mutasa ES, Tymon AM, Pardo EH, O'Shea SF, et al.
677            Two classes of homeodomain proteins specify the multiple a mating types of the
678            mushroom Coprinus cinereus. Plant Cell. 1994;6: 1467–1475.
679            doi:10.1105/tpc.6.10.1467

680   37.   Spit A, Hyland RH, Mellor EJC, Casselton LA. A role for heterodimerization in
681         nuclear localization of a homeodomain protein. Proc Natl Acad Sci USA. 1998;95:
682         6228–6233. doi:10.1073/pnas.95.11.6228

683   38.   O'Malley MA. Endosymbiosis and its implications for evolutionary theory. Proc Natl
684         Acad Sci USA. 2015;112: 10270–10277. doi:10.1073/pnas.1421389112

685   39.   López-García P, Eme L, Moreira D. Symbiosis in eukaryotic evolution. Journal of
686         Theoretical Biology. 2017;434: 20–33. doi:10.1016/j.jtbi.2017.02.031

687   40.   Zaremba-Niedzwiedzka K, Caceres EF, Saw JH, Bäckström D, Juzokaite L,
688         Vancaester E, et al. Asgard archaea illuminate the origin of eukaryotic cellular
689         complexity. Nature. 2017;541: 353–358. doi:10.1038/nature21031

690   41.   Barton MK, Poethig RS. Formation of the shoot apical meristem in Arabidopsis
691         thaliana: an analysis of development in the wild type and in the shoot meristemless
692         mutant. Development. 1993;119: 823–831.

693   42.   Li E, Bhargava A, Qiang W, Friedmann MC, Forneris N, Savidge RA, et al. The
694         Class II KNOX gene KNAT7 negatively regulates secondary wall formation in
695         Arabidopsis and is functionally conserved in Populus. The New phytologist.
696         2012;194: 102–115. doi:10.1111/j.1469-8137.2011.04016.x

697   43.   Furumizu C, Alvarez JP, Sakakibara K, Bowman JL. Antagonistic roles for KNOX1
698         and KNOX2 genes in patterning the land plant body plan following an ancient gene
699         duplication. Qu L-J, editor. PLoS Genet. 2015;11: e1004980.
700         doi:10.1371/journal.pgen.1004980

701   44.   Mukherjee K, Bürglin TR. Comprehensive analysis of animal TALE homeobox
702         genes: new conserved motifs and cases of accelerated evolution. J Mol Evol.
703         2007;65: 137–153. doi:10.1007/s00239-006-0023-0

704   45.   Keller MD, Selvin RC, Claus W, Guillard RRL. MEDIA FOR THE CULTURE OF
705         OCEANIC ULTRAPHYTOPLANKTON1,2. Journal of Phycology. 2007;23: 633–
706         638. doi:10.1111/j.1529-8817.1987.tb04217.x

707   46.   Guillard RRL, Hargraves PE. Stichochrysis immobilis is a diatom, not a
708         chrysophyte. Phycologia. 1993;32: 234–236. doi:10.2216/i0031-8884-32-3-234.1

709   47.   Bailey TL, Johnson J, Grant CE, Noble WS. The MEME Suite. Nucleic Acids
710         Research. 2015;43: W39–49. doi:10.1093/nar/gkv416

711   48.   Finn RD, Attwood TK, Babbitt PC, Bateman A, Bork P, Bridge AJ, et al. InterPro in
712         2017-beyond protein family and domain annotations. Nucleic Acids Research.
713         2017;45: D190–D199. doi:10.1093/nar/gkw1107

714

715   **Figure Legends**

716

717 **Fig 1. Common origin of heterodimerizing TALE homeobox TFs.** Hypothesized

718 homodimerizing proto-TALE protein (top) duplicated before the eukaryotic radiations into

719 animals/fungi/amoebae vs. algae/plants. Lineage-specific diversification soon followed,

720 generating heterodimeric configurations distinct at the phylum-level. (Left) Each lineage

721 possesses one or two classes of potential heterodimeric TALEs, which are summarized onto

722 the eukaryotic phylogeny. A representative species name is given for each analyzed lineage.

723 (Right) Summary of TALE configurations, coupling members of the PBC/PBX/GLX group

724 that shares PBC-homology domains and of the MEIS/KNOX group that shows homology in

725 the KN-A/B domains N-terminal to the homeodomain. Lightly shaded boxes depict homology

726 domains, whose names are provided above. Open areas in the domain boxes indicate the

727 absence of MEINOX-motif for PBX-Red, KN-A for KNOX-Red1 and ELK for KNOX-Red2.

728 Colored vertical lines in the HD indicate two shared introns at 44/45 (orange over 'H' In HD)

729 and 48[2/3] (blue over 'D' in HD), whose alternating existence between the two groups

730 suggests independent diversification of TALE heterodimerization. HD: Homeodomain; PBL-

731 C: PBL-Chloro; PBL-R: PBL-Red.

732

733 **Fig 2.  Maximum likelihood (ML) phylogeny of the TALE superclass homeodomain in**

734 **Archaeplastida supports ancient division between KNOX- and non-KNOX TALE**

735 **groups.** The consensus tree out of 1000 bootstrap trees is shown. The three numbers at

736 critical nodes show %bootstrap, %SH, and Bayesian posterior probability in support of

737 clades. The tree contains two outgroup clades marked by black squares at nodes, and two

738 Archaeplastida clades, one combining most KNOX sequences marked by the red square

739 and the other combining all non-KNOX sequences marked by the blue square. Vertical bars

740 on the right depict the distribution of outgroup in black, KNOX in red, and non-KNOX

741 sequences in blue. Red dots by the sequence names indicate the presence of KN-A or KN-B

742 domains, and blue dots indicate the presence of a PBC-homology domain. Truncated

743 sequences not available for homology domain analysis are marked with open black boxes.

744 Filled black boxes indicate the absence of a KN-A/B or PBC-homology domain. Proposed

745 classification is indicated by the vertical lines. Dotted vertical lines indicate suggested class

746 members placed outside the main clade for the class in the phylogeny. PBX-Red sequences

747 are found in four separate clades, marked by purple shades on the blue section of the

748 vertical bars. Colors of the sequence names indicate their phylogenetic group: Blue for

749 Glaucophyta, purple for Rhodophyta, green for prasinophytes, light blue for the

750 chlorophytes, orange for Streptophyta, and black for outgroups. The ruler shows genetic

751 distance. Details of the sequences analyzed by this phylogeny are provided in S2

752 Spreadsheet.

753

754    **Fig 3. Archaeplastida non-KNOX group TALEs possess a PBC-like domain (PBL)**

755    **consisting of N-terminal MEINOX homology and C-terminal PBC-B homology.** Amino

756    acid letters in black with gray shades, in white with light shades, and in white with black

757    shades show more than 60%, 80%, or 100% similarity in each column. Inverse red triangles

758    indicate the discarded sequences in un-aligned insertions. (A) PBL-Glauco domain

759    alignment, including two Glaucophyta sequences sharing homology in both MEINOX

760    homology and C-terminal half of the PBC-B domain with non-Archaeplastida TALE

761    sequences. Red box indicates the ELK domain. (B) PBL-Red domain alignment. All

762    Rhodophyta non-KNOX sequences possess a PBL domain with poor MEINOX homology.

763    (C) PBL-Chloro domain alignment. Cyanophora_paradox_20927.63 is included for

764    comparison. Picocystis_salinarum_02499 is a founding member of GLX class with a PBL-

765    Chloro domain. (D) Comparison among PBL domains. The top row shows the consensus

766    made from the alignment of (A), (B), and (C) combined and the lower consensus sequences

767    are collected from the individual alignments presented in (A), (B), and (C).

768

769    **Fig 4. TALE TFs engage in heterodimerization networks between KNOX and non-**

770    **KNOX groups.** The bait constructs conjugated to the GAL4 DNA-binding domain (DBD) and

771    the prey constructs conjugated to the GAL4 transcriptional activation domain (AD) are listed

772    in the table. Construct combinations, numbered 1-8, are arranged in wedges clock-wise,

773    starting at 9 o'clock as labeled in the -LT panels. Confirmed interacting pairs are shown in

774    bold faces in the table. The laminin and T-Antigen (T-Ag) pair, known to be interacting

775    partners, was plated in the 8th sector as a positive control. (A) Assays using *M. commoda*

776    TALEs. (B) Assays using *O. tauri* TALEs. (C) Assays using *P. salinarum* TALEs. KNOX-tr

777    refers to the N-terminal truncated KNOX construct for preventing self-activation. (D)

778    Detailed construct information is provided in S5 Spreadsheet.

779

22

1 **Figures for Joo et al.**

2 Article title: Common ancestry of heterodimerizing TALE homeobox transcription factors across
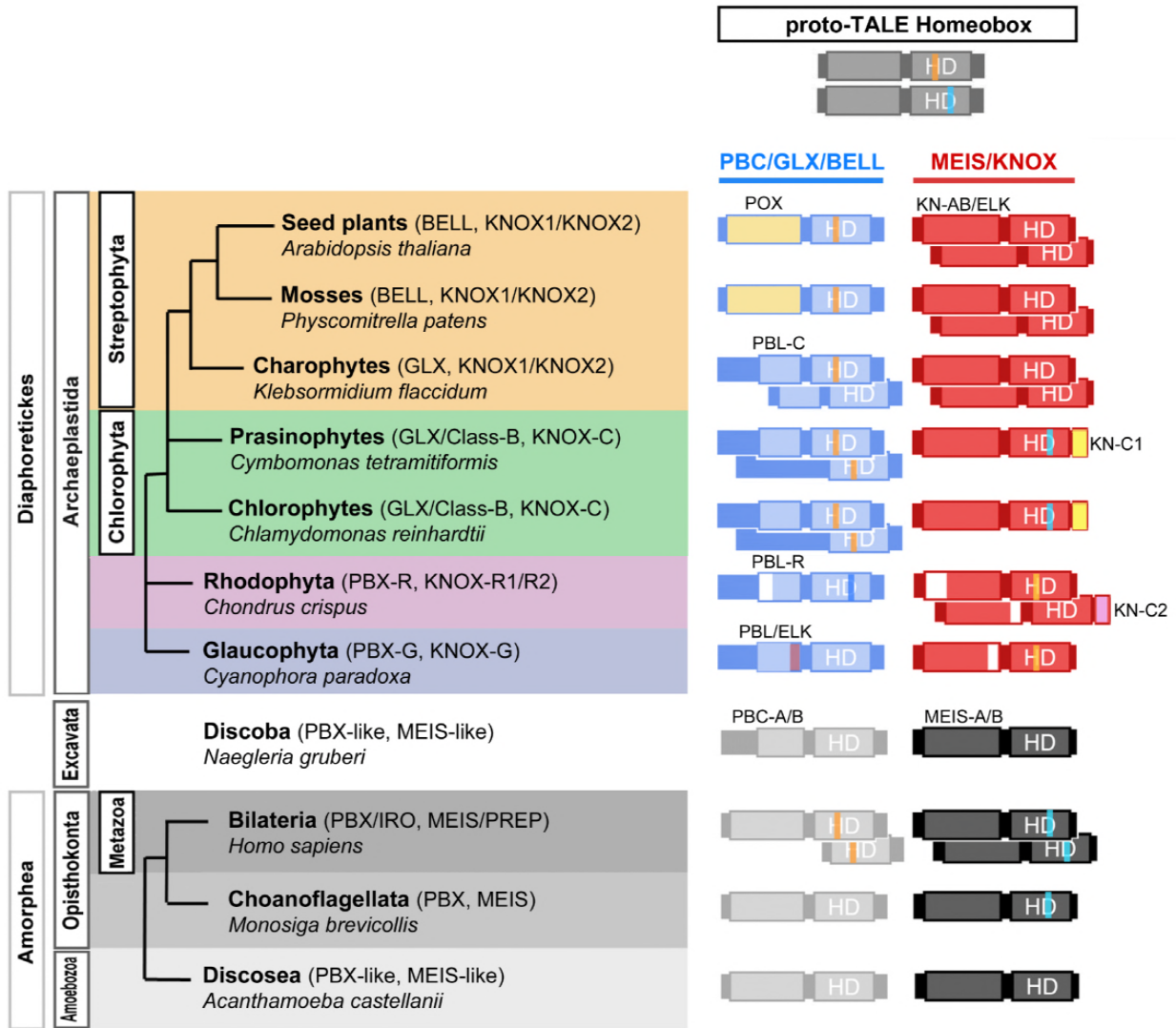
3 Metazoa and Archaeplastida.

4 Authors: Sunjoo Joo, Ming Hsiu Wang, Gary Lui, Jenny Lee, Andrew Barnas, Eunsoo Kim,

5 Sebastian Sudek, Alexandra Z. Worden, and Jae-Hyeok Lee.
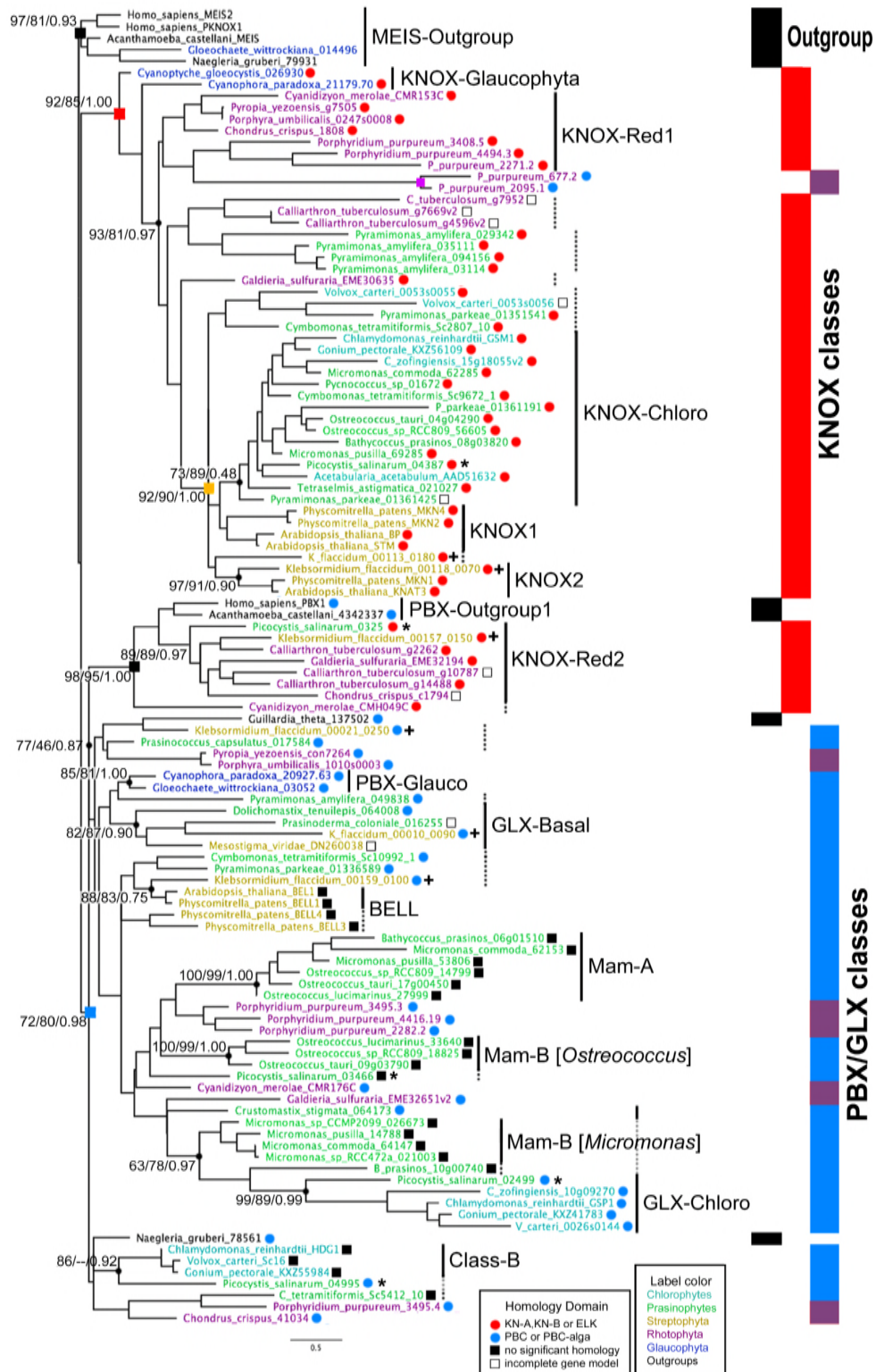
6

7

8

9    **Fig 1**.



10

**Fig 1**. **Common origin of heterodimerizing TALE homeobox TFs.** We propose that a homodimerizing
proto-TALE protein duplicated (top) prior to the major bifurcation resulting in animals/fungi/amoebae vs.
algae/plants. Lineage-specific diversification soon followed, generating heterodimeric configurations
distinct at the phylum-level. These configurations usually couple members of the PBC/PBX/GLX group
that shares PBC-homology domains and MEIS/KNOX group that shows homology in the KN-A/B domains
N-terminal to the homeodomain. Each lineage possesses one or two classes of potential heterodimeric
partners. Major TALE classes are mapped onto the eukaryotic phylogeny. A representative species name
is given for each analyzed lineage. Open boxes in the domain diagrams indicate the absence of MEINOX
for PBX-Red, KN-A for KNOX-Red1 and ELK for KNOX-Red2. Colored vertical lines in the HD indicate
two proposed ancestral introns at 44/45 (orange over 'H' In HD) and 48[2/3] (blue over 'D' in HD), whose
alternating existence between the two groups suggests independent diversification of TALE
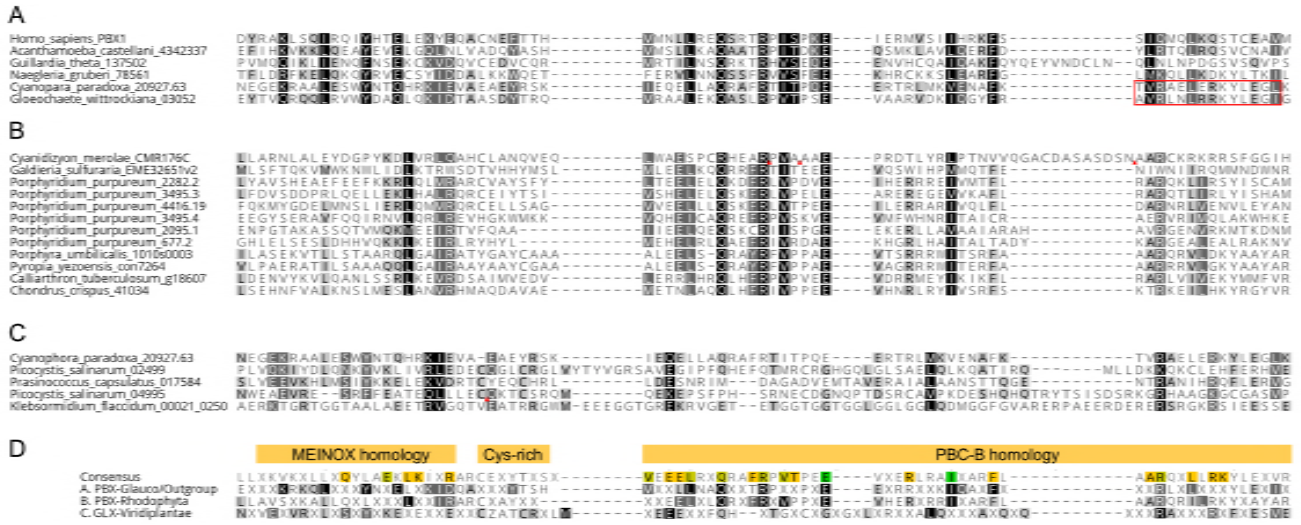heterodimerization.

23

**Fig 2. Maximum likelihood (ML) phylogeny of the TALE superclass homeodomain in**

**Archaeplastida supports ancient division between KNOX- and non-KNOX TALE groups.**

**Fig 2. Maximum likelihood (ML) phylogeny of the TALE superclass HD in Archaeplastida supports ancient division between KNOX- and non-KNOX TALE groups.** The consensus tree out of 1000 bootstrap trees is shown. The three numbers shown at nodes are %bootstrap, %SH, and Bayesian posterior probability in support of clades. The tree contains two outgroup clades marked by black squares, and two Archaeplastida clades, one combining most KNOX sequences marked by the red square and the other combining all non-KNOX sequences marked by the blue square. Vertical bars on the right depict the distribution of outgroup, KNOX, and non-KNOX sequences. KNOX sequences are marked with red dots indicating the presence of KN-A or KN-B domains. GLX/PBX sequences are marked with blue dots indicating the presence of a PBC-homology domain. Truncated sequences not available for homology domain analysis are marked with open black boxes. Filled box indicates the absence of a KN-A/B or PBC-homology domain. Proposed classification is shown by black vertical lines. Dotted lines indicate sequences related to a class but placed outside the main clade for the class. PBX-Red sequences are found in four paraphyletic clades, marked by purple shades on the blue vertical bar. Sequence IDs containing the species name are colored by their phylogeny: Blue for Glaucophyta, purple for Rhodophyta, green for prasinophytes, light blue for the chlorophytes, orange for Streptophyta, and black for outgroups. The ruler shows genetic distance. All the sequences and their phylogenetic information are found in S2 Spreadsheet.

*Gloeochaete_wittrockiana_014496 is considered as a sequence from a bannelid-type amoeba that contaminated the original culture (SAG46.84) for the MMETSP1089 transcriptome. **Association of KNOX-Red2 class sequences to Amorphea PBC sequences is attributed to a shared WFGN motif determining DNA-binding specificity of the homeodomain via convergent evolution.
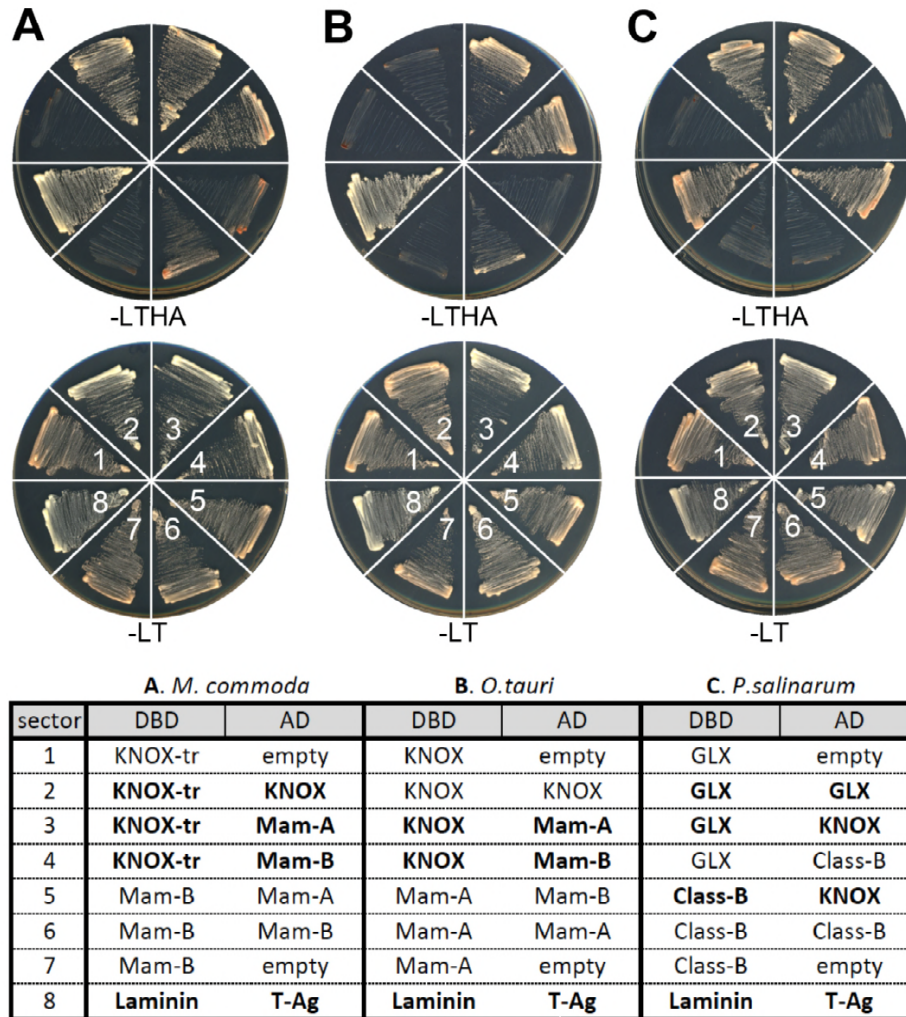
**Fig 3.**

**Fig 3. Archaeplastida Non-KNOX group TALEs possess a PBL domain sharing homology with metazoan PBC class TALEs.** (A) PBL-Glauco domain alignment. Two Glaucophyta non-KNOX sequences possess a PBC-homology domain spanning MEINOX and C-terminal half of the PBC-B domains which is shared among the three outgroup TALE sequences analyzed in this study. (B**)** PBL-Red domain alignment. All Rhodophyta non-KNOX sequences possess a PBC-homology domain that can be aligned to the MEINOX/PBC-C domains. (C**)** PBL-Chloro domain. Four non-KNOX sequences show >10% amino acid identity to one of the other PBC-homology blocks presented in (A) and (B). Picocystis_salinarum_02499 is a founding member of GLX class with a PBL-Chloro domain. **D**. Comparison among PBC-homology domains. The top row shows the consensus made from the alignment of (A), (B), and (C) combined and the lower consensus sequences are collected from the individual alignments presented in (A), (B), and GLX alignment (S3 Fig). Amino acid letters in black with Gray shades, in white with light shades, and in white with black shades show more than 60%, 80%, or 100% similarity in each column.

70 **Fig 4**.



| sector | A. *M. commoda* DBD | AD | B. *O.tauri* DBD | AD | C. *P.salinarum* DBD | AD |
|--------|------|------|------|------|------|------|
| 1 | KNOX-tr | empty | KNOX | empty | GLX | empty |
| 2 | **KNOX-tr** | **KNOX** | KNOX | KNOX | **GLX** | **GLX** |
| 3 | **KNOX-tr** | **Mam-A** | KNOX | Mam-A | GLX | KNOX |
| 4 | **KNOX-tr** | **Mam-B** | KNOX | Mam-B | GLX | Class-B |
| 5 | Mam-B | Mam-A | Mam-A | Mam-B | **Class-B** | **KNOX** |
| 6 | Mam-B | Mam-B | Mam-A | Mam-A | Class-B | Class-B |
| 7 | Mam-B | empty | Mam-A | empty | Class-B | empty |
| 8 | **Laminin** | **T-Ag** | **Laminin** | **T-Ag** | **Laminin** | **T-Ag** |

71

72

73 **Fig 4**. **All Chlorophyta TALE TFs engage in heterodimerization networks.** The bait constructs
74 conjugated to the GAL4 DNA-binding domain and the prey constructs conjugated to the GAL4
75 transcriptional activation domain are listed in the table. Construct combinations, numbered 1-8, are
76 arranged in wedges clock-wise, starting at 9 o'clock as labeled in (A). Interacting pairs confer yeast
77 growth in Leu-/Trp-/His-/Ade- (-LTHA) medium. Confirmed interacting pairs are shown in bold faces in the
78 table. The laminin and T-Antigen (T-Ag) pair, known to be interacting partners, was plated in the 8th
79 sector as a positive control. (A) Assays using *M. commoda* TALEs. (B**)** Assays using *O. tauri* TALEs. **C**.
80 Assays using *P. salinarum TALEs*. Class-A refers to the GLX-Chloro homolog. Details of the construct
81 information are found in S5 Spreadsheet.

82

83

24