1 **Title: The puzzle of metabolite exchange and identification of putative octotrico peptide**

2 **repeat expression regulators in the nascent photosynthetic organelles of *Paulinella***

3 ***chromatophora***

4 **Running title:** Metabolic and genetic integration of the chromatophore

5 **Authors:** Linda Oberleitner[a], Gereon Poschmann[b], Luis Macorano[a], Stephan Schott-

6 Verdugo[c,d], Holger Gohlke[c,e], Kai Stühler[b,f], and Eva C. M. Nowack[a#]

7 [a]Department of Biology, Heinrich-Heine-Universität Düsseldorf, 40225 Düsseldorf, Germany

8 [b]Institute for Molecular Medicine, Proteome Research, Medical Faculty, Heinrich-Heine-

9 Universität Düsseldorf, 40225 Düsseldorf, Germany

10 [c]Institute for Pharmaceutical and Medicinal Chemistry, Heinrich Heine University Düsseldorf,

11 40225 Düsseldorf, Germany

12 [d]Centro de Bioinformática y Simulación Molecular (CBSM), Faculty of Engineering,

13 Universidad de Talca, 3460000 Talca, Chile

14 [e]John von Neumann Institute for Computing (NIC), Jülich Supercomputing Centre (JSC), and

15 Institute of Biological Information Processing (IBI-7: Structural Biochemistry),

16 Forschungszentrum Jülich GmbH, 52428 Jülich, Germany

17 [f]Molecular Proteomics Laboratory, BMFZ, Heinrich-Heine-Universität Düsseldorf, 40225

18 Düsseldorf, Germany

19

20 [#]**Correspondence to** Eva C. M. Nowack, e.nowack@uni-duesseldorf.de

21

22 **Word count Abstract:** 248 + 150

23 **Word count text (excluding references, table footnotes, and figure legends):** 5342

24

**Abstract**

The cercozoan amoeba *Paulinella chromatophora* contains photosynthetic organelles - termed chromatophores - that evolved from a cyanobacterium, independently from plastids in plants and algae. Despite the more recent origin of the chromatophore, it shows tight integration into the host cell. It imports hundreds of nucleus-encoded proteins, and diverse metabolites are exchanged across the two chromatophore envelope membranes. However, the limited set of chromatophore-encoded transporters appears insufficient for supporting metabolic connectivity or protein import. Furthermore, chromatophore-localized biosynthetic pathways as well as multiprotein complexes include proteins of dual genetic origin, suggesting coordination of gene expression levels between chromatophore and nucleus. These findings imply that similar to the situation in mitochondria and plastids, nuclear factors evolved that control metabolite exchange and gene expression in the chromatophore. Here we show by mass spectrometric analyses of enriched insoluble protein fractions that, unexpectedly, nucleus-encoded transporters are not inserted into the chromatophore inner envelope membrane. Thus, despite the apparent maintenance of its barrier function, canonical metabolite transporters are missing in this membrane. Instead we identified several expanded groups of short chromatophore-targeted orphan proteins. Members of one of these groups are characterized by a single transmembrane helix, and others contain amphipathic helices. We hypothesize that these proteins are involved in modulating membrane permeability. Furthermore, we identified an expanded family of chromatophore-targeted helical repeat proteins. These proteins show similar domain architectures as known organelle-targeted octotrico peptide repeat expression regulators in algae and plants suggesting their convergent evolution as nuclear regulators of gene expression levels in the chromatophore.

**Importance**

The endosymbiotic acquisition of mitochondria and plastids >1 billion years ago was central for the evolution of eukaryotic life. However, owing to their ancient origin, these organelles provide only limited insights into the initial stages of organellogenesis. The chromatophore in *Paulinella* evolved ~100 million years ago and thus, offers the possibility to gain valuable insights into early stages and common rules in organelle evolution. Critical to organellogenesis

55    appears to be the establishment of nuclear control over metabolite exchange and gene

56    expression in the endosymbiont. Here we show that the mechanism generating metabolic

57    connectivity of the chromatophore fundamentally differs from the one for mitochondria and

58    plastids, but likely rather resembles the poorly understood mechanism in various bacterial

59    endosymbionts in plants and insects. Furthermore, we describe a novel class of

60    chromatophore-targeted helical repeat proteins which evolved convergently to plastid-

61    targeted expression regulators and are likely involved in gene expression control in the

62    chromatophore.

**Introduction**

Endosymbiosis has been a major driver for the evolution of cellular complexity in eukaryotes. During organellogenesis, linkage of the previously independent biological networks of the former host and endosymbiont resulted in a homeostatic and synergistic association. Two critical factors during this dauntingly complex process appear to be the establishment of metabolic connectivity between the symbiotic partners, and nuclear control over protein levels within the organelle.

Besides mitochondria and primary plastids that evolved via endosymbiosis >1 billion years ago, recently, a third organelle of primary endosymbiotic origin has been identified (1, 2). The photosynthetically active 'chromatophore' of cercozoan amoebae of the genus *Paulinella* evolved around 100 million years ago from a cyanobacterium (3, 4). Hence, scrutiny of *Paulinella* can help to determine the degrees of freedom in the integration process of a eukaryotic organelle. Similar to the evolution of mitochondria and plastids, also in the chromatophore, reductive genome evolution resulted in the loss of many metabolic functions (5, 6), around 70 genes were transferred from the chromatophore to the nucleus (7-9), and functions lost from the chromatophore genome are compensated by import of nucleus-encoded proteins (10, 11). In a previous study, we identified by mass spectrometry (MS) around 200 nucleus-encoded, chromatophore-targeted proteins in *Paulinella chromatophora* that fall into two classes (10). Short import candidates [<90 amino acids (aa)] lack obvious targeting signals, whereas long import candidates (>250 aa) carry a conserved N-terminal sequence extension – likely a targeting signal – that is referred to as 'chromatophore transit peptide' (crTP). Bioinformatic identification of crTPs allowed to extend the catalogue of import candidates to >400 proteins (10).

Metabolic capacities of chromatophore and host cell are complementary resulting in the need for extensive exchange of metabolites such as sugars, amino acids, and cofactors across the two envelope membranes that surround the chromatophore (5, 10, 12). Furthermore, substrates for carbon, sulfur, and nitrogen assimilation (e.g. $HCO_3^-$, $SO_4^{2-}$, $NH_4^+$) and metal ions (e.g. $Mg^{2+}$, $Cu^{2+}$, $Mn^{2+}$, and $Co^{2+}$) that serve as cofactors of chromatophore-localized proteins have to be imported into the chromatophore. Whereas the chromatophore inner membrane (IM) clearly derives from the cyanobacterial plasma membrane, the outer membrane (OM) has been interpreted as being host-derived (13, 14). The nature of the

4

94 transporters underlying the deduced solute (and protein) transport processes across this
95 membrane system is unknown.

96      In plants and algae, transport across the plastid IM is mediated by a large set of multi-
97 spanning transmembrane (TM) proteins that are highly specific for their substrates. These
98 transporters contain usually four or more TM α-helices (TMHs) and are of the single subunit
99 secondary active or channel type (15). This set of transporters apparently evolved mainly via
100 the retargeting of existing host proteins to the plastid IM rather than the repurposing of
101 endosymbiont proteins (15-17). Transport across the plastid OM is enabled largely by (semi-
102 )selective pores formed by nucleus-encoded $\beta$-barrel proteins (18).

103      Another important issue during organellogenesis is the establishment of nuclear
104 control over organellar gene expression supporting (i) adjustment of the organelle to the
105 physiological state of the host cell, and (ii) assembly of organelle-localized protein complexes
106 composed of subunits encoded in either the organellar or nuclear genome in stoichiometric
107 amounts (19, 20). Also in *P. chromatophora*, the import of nucleus-encoded proteins resulted
108 in protein complexes of dual genetic origin [e.g. photosystem I (11)]. The difference in copy
109 numbers between chromatophore and nuclear genome (~100 vs one or two copies) (7) calls
110 for coordination of gene expression between nucleus and chromatophore.

111      To test the hypotheses that nuclear factors were recruited to establish (i) metabolic
112 connectivity between chromatophore and host cell and (ii) control over gene expression levels
113 within the chromatophore, here we analyzed the previously obtained proteomic dataset
114 derived from isolated chromatophores and a newly generated proteomic dataset derived
115 from enriched insoluble chromatophore proteins with a focus on chromatophore-targeted TM
116 proteins and putative expression regulators.

117

118 **Results**

119 **Paucity of chromatophore-encoded solute transporters**

120 Although diverse metabolites have to be exchanged constantly between the chromatophore
121 and cytoplasm, we identified genes for only 25 solute transporters on the chromatophore
122 genome (5, 10, 12). As judged from the localization of their cyanobacterial orthologs, only 19

5

123    of these transporters putatively localize to the chromatophore IM [**Fig. 1, Table S1**]. In

124    comparison, in *Synechococcus* sp. WH5701, a free-living relative of the chromatophore (21),

125    and the model cyanobacterium *Synechocystis* sp. PCC6803, genes for ~89 and >100 putative

126    envelope transporters were identified, respectively [**Fig. 1**, **Table S1**, (22)]. Substrates of most

127    of the chromatophore IM transporters are – according to annotation – restricted to inorganic

128    ions (e.g. $Na^+$, $K^+$, $Fe^{2+}$, $Mg^{2+}$, $PO_4^{2-}$, $HCO_3^-$). Notably, cyanobacterial uptake systems for

129    nitrogen and sulfur compounds such as nitrate (23), ammonium (24), urea (25), amino acids

130    (26) or sulfate (27) are missing. Only one transporter of the DME-family (10 TMS

131    Drug/Metabolite Exporter; PCC0734) could potentially be involved in metabolite export, and

132    one transporter of the DASS-family (Divalent Anion:$Na^+$ Symporter; PCC0664) could facilitate

133    import of either di-/tricarboxylates or sulfate via $Na^+$ symport. However, due to the multitude

134    of substrates transported by members of both families (28, 29), precise substrate specificities

135    cannot be predicted. Chromatophore-encoded β-barrel OM pores could not be identified.

136        In contrast, in plants and algae, a combination of bioinformatic and proteomic studies

137    identified 100-150 putative solute transporters in the plastid IM; 37 of these transporters have

138    been confidently assigned functions and many of them transport metabolites [(16, 30-32), **Fig.**

139    **1**, and **Table S1**]. Several porins are known to permit passage of solutes across the chloroplast

140    OM (18, 33-35). Almost all of these transport systems are encoded in the nucleus and

141    posttranslationally inserted into the plastid envelope membranes.

142

143    **Enrichment of insoluble protein fractions and proteomic analysis**

144    The scarcity of chromatophore-encoded solute transporters suggested that in *P.*

145    *chromatophora*, as in plastids, nucleus-encoded transport systems establish metabolic

146    connectivity of the chromatophore. However, among 432 previously identified import

147    candidates (10), only 3 proteins contained more than one predicted TMH (**Table 1**). One of

148    these proteins (identified by in silico prediction) contains two TMHs, only one of which is

149    predicted with high confidence. Of the other two proteins (identified by MS), one is short and

150    contains two predicted TMHs; the other contains eight predicted TMHs. However, this latter

151    protein was identified with one peptide only and shows no BlastP hits against the NCBI nr

152    database, whereas an alternative ORF (in the reverse complement) shows similarity to an

153   NAD-dependent epimerase/dehydratase. Therefore, this latter protein likely represents a

154   false positive (a false discovery rate of 1% was accepted in this analysis).

155

156   **Table 1: No nucleus-encoded solute transporters among previously identified import**

157   **candidates.** The table lists numbers of proteins previously identified to be imported into the

158   chromatophore [by in silico prediction (based on presence of a crTP), liquid chromatography

159   coupled to tandem MS (LC-MS/MS), and total] (10) sorted by the number of predicted TMHs

160   (outside of the crTP).

|                             | 0 TMH | 1 TMH | >1 TMH |
|-----------------------------|-------|-------|--------|
| **In silico predicted (crTP)** | 289   | 3     | 1      |
| **LC-MS/MS identified**     | 194   | 11    | 2      |
| **Total**                   | 416   | 13    | 3      |

161

162   The absence of multi-spanning TM proteins among import candidates could have two

163   reasons. (i) These proteins might lack a crTP, impairing their prediction as import candidates.

164   (ii) TM proteins are often underrepresented in LC-MS analyses owing to low abundance levels

165   as well as unfavorable retention and ionization properties. In fact, our previous MS analysis

166   identified 47% of the soluble but only 21% of TMH-containing chromatophore-encoded

167   proteins (**Fig. 2C**).

168   Thus, to enhance identification of TM proteins, we enriched TM proteins by collecting

169   the insoluble fractions from isolated chromatophores (CM samples) and intact *P.*

170   *chromatophora* cells (PM samples). Electron microscopic analysis of isolated chromatophores

171   suggested that the OM is lost during chromatophore isolation [**Fig. 2A**, compare (13, 14)].

172   Comparison of CM and PM samples to chromatophore lysates (CL samples) by SDS-PAGE

173   revealed distinct banding patterns between the three samples and high reproducibility

174   between three biological replicates (**Fig. 2B**). Further enrichment of membrane proteins or

175   separation of IM, OM, and thylakoids was not feasible owing the slow growth of *P.*

176   *chromatophora,* low yield of chromatophore isolations, and the loss of the OM. Two

7

177 consecutive, independent MS analyses of three replicates of each, CM, PM, and CL samples

178 led to the identification of 1,886 nucleus- and 555 chromatophore-encoded proteins over all

179 fractions (**Tables 2** and **S2**). Although most chromatophore-localized TM proteins were also

180 identified in our analyses in CL samples (**Table 2**), individual TM proteins were clearly enriched

181 in CM compared to CL samples (**Fig. S1**).

182

183 **Table 2: Proteins identified in this study by LC-MS/MS.** Numbers of chromatophore-encoded

184 (CE) and nucleus-encoded (NE) proteins identified in at least one out of two independent MS

185 experiments with ≥3 spectral counts (SpC) in chromatophore-derived samples (i.e. CM+CL) or

186 whole cell membranes (PM). The number of predicted TMHs (outside of the crTP) is indicated.

187 For proteins identified in CM samples, total number of proteins and number of proteins

188 enriched in CM as compared to PM samples (in brackets) is indicated separately.

| | All Proteins | | 1 TMH | | >1 TMH | |
|---|---|---|---|---|---|---|
| | CE | NE | CE | NE | CE | NE |
| CM | 533 (506) | 297 (236) | 28 (24) | 20 (13) | 70 (67) | 5 (2) |
| CL | 551 | 354 | 28 | 25 | 67 | 7 |
| **Chromatophore total** | 555 | 361 | 28 | 27 | 70 | 7 |
| PM | 385 | 1691 | 24 | 209 | 50 | 175 |
| **Total** | 555 | 1886 | 28 | 218 | 70 | 179 |

189

190     In CM samples, 46% (or 98/213) of the chromatophore-encoded TM proteins were

191 identified, representing a gain of 118% compared to our previous analysis (**Fig. 2C**); in

192 particular, of the 25 chromatophore-encoded solute transport systems, 72% (or 18) were

193 identified with at least one subunit, and 60% (or 15) were identified with their TM subunit

194 (**Fig. 2D**) while our previous study identified only three of these transporters. Highest

195 intensities (representing a rough estimation for protein abundances) were found in CM

8

196     samples for an ABC-transporter annotated as multidrug importer of the P-FAT family (level 4

197     to 5, placing the transporter among the 10% most abundant proteins in CM). Also the

198     bicarbonate transporter BicA, two multidrug efflux ABC-transporters, and a NhaS3

199     proton/sodium antiporter were found in the upper tiers of abundance levels (level 3 to 4,

200     placing them among the 30% most abundant proteins in CM). The remaining transporters

201     showed moderate to low abundance levels (**Fig. 2D**).

202

203     **No multi-spanning TM proteins appear to be imported into the chromatophore**

204     Determination of nucleus-encoded proteins enriched in CM compared to PM samples led to

205     the identification of 188 high confidence (HC) [and further 48 low confidence (LC); see

206     Methods and **Fig. S2**] import candidates (**Fig. 3A, Table S3**). Nucleus-encoded multi-spanning

207     TM proteins appeared invariably depleted in chromatophores (**Figs. 3B, C**). Only two of 236

208     import candidates were multi-spanning TM proteins (**Table 2**). However, one of these (with 7

209     predicted TMHs, scaffold1608-m.20717, arrowhead in **Fig. 3B**) was identified by only one

210     hepta-peptide and shows no similarity to other proteins in the NCBI nr database whereas an

211     overlapping ORF (in another reading frame) encodes a peroxidase that was MS-identified in

212     ref. (10) likely classifying the protein as a false positive. For the other import candidate

213     (scaffold18898-m.107131; with an enrichment level close to 0; arrowhead in **Fig. 3C**) a full-

214     length transcript sequence is missing precluding determination of the correct start codon.

215     Thus, this protein might represent in fact a short import candidate with a single TMH. Of the

216     three nucleus-encoded multi-spanning TM proteins that were present but appeared depleted

217     in CM compared to PM samples (**Table 2**), two were annotated as mitochondrial NAD(P)

218     transhydrogenase and mitochondrial ATP/ADP translocase, suggesting a mild contamination

219     of CM samples with mitochondrial membrane material.

220        In comparison, 70 chromatophore-encoded multi-spanning TM proteins were

221     identified in CM samples, and 67 of these appeared enriched in CM samples. In PM samples,

222     50 chromatophore- and 175 nucleus-encoded multi-spanning TM proteins were found (**Table**

223     **2**).

224

225 **Targeting of single-spanning TM proteins and antimicrobial peptide-like proteins to the**

226 **chromatophore**

227 In contrast to the striking lack of multi-spanning TM proteins, there were 13 (5 HC and 8 LC)

228 single-spanning TM proteins (containing one TMH outside of the crTP) among the identified

229 import candidates (**Table 2**). Three of these proteins contain a TMH close to their C-terminus

230 and likely represent tail-anchored proteins. One of these proteins is long and annotated as

231 low-density lipoprotein receptor-related protein 2-like, the other two (with N-terminal

232 sequence information missing) as polyubiquitin. However, most import candidates with one

233 TMH (10 proteins) represent short proteins. These short import candidates included two high

234 light-inducible proteins [i.e. thylakoid-localized cyanobacterial proteins involved in light

235 acclimation of the cell (9)]. The remaining eight proteins are orphan proteins lacking

236 detectable homologs in other species (BlastP against NCBI nr database, cutoff e$^{-03}$); all of these

237 contain a TMH with a large percentage of small aa (26-45% Gly, Ala, Ser) close to their

238 negatively charged N-terminus (**Fig. 4A**).

239 In our previous proteome analysis, short orphan proteins represented the largest

240 group of MS-identified import candidates (1/3 of total). However, most of these proteins did

241 not possess predicted TMHs. Based on the occurrence of specific Cys motifs (CxxC, CxxxxC)

242 and stretches of positively charged aas these short proteins were described as antimicrobial

243 peptide (AMP)-like proteins (10). Including the eight TMH-containing proteins (see above), the

244 recent study identified further 19 short orphan import candidates (or – only few proteins –

245 showing similarity to hypothetical proteins in other species). Scrutiny of all 88 short orphan

246 import candidates (resulting from both studies together) revealed that besides the TMH-

247 containing proteins (group 1, 10 proteins), these short import candidates form at least three

248 further distinct groups (**Fig. 4A**). Members of group 2 (12 proteins) contain a conserved motif

249 of unknown function that occurs also in bacterial proteins that often possess domains pointing

250 towards DNA processing-related functions (**Fig. 4A** and **B**). Members of group 3 (10 proteins)

251 contain another conserved motif of unknown function that encompasses two Cys-motifs

252 (CxxxxC and CxxC). Members of group 4 (30 proteins) show either one or two CxxC mini motifs

253 (one of these is often CPxCG) but no further sequence conservation. The remaining 26 short

254 orphan import candidates have no obvious common characteristics but several appear to have

255 a propensity to form amphipathic helices (**Fig. 4A**).

256    Screening a large nuclear *P. chromatophora* transcriptome dataset (7) revealed

257    additional putative members of groups 1 to 3 (**Fig. 4A** and **Fig. S3**): further 53 translated

258    transcripts represent short proteins with a predicted TMH in the N-terminal 2/3 of the

259    sequence that is rich (>20%) in small aa and have an N-terminus with a net charge ≤0. Notably,

260    the TMHs of >90% of all group 1 proteins comprise at least one (small)xxx(small) motif which

261    can promote oligomerization of single-spanning TM proteins (36). Furthermore, many of these

262    putative group 1 short import candidates are predicted to have antimicrobial activity and/or

263    pore-lining residues (**Table S4**). Further 192 and 28 translated transcripts contain the

264    conserved motifs of group 2 or 3, respectively. Importantly, all MS-identified members of

265    these extended protein groups were identified in chromatophore-derived samples in this and

266    our previous analysis.

267

268    **An expanded family of octotrico peptide repeat putative expression regulators is targeted**

269    **to the chromatophore**

270    Of the 235 import candidates (excluding the false positive, see above) identified in this study

271    (**Fig. 3A**), 159 were known import candidates (10) (**Fig. 5A**, **Table S3**), with 46 proteins

272    experimentally confirming import candidates previously only predicted in silico. 76 proteins

273    represent new import candidates, mostly lacking N-terminal sequence information (42

274    proteins) or representing short import candidates (22 proteins). A particularly large number

275    of newly MS-identified import candidates (24 proteins) fall into the category 'genetic

276    information processing' (**Fig. 5B**). Among these proteins an expanded group of 10 RNA-binding

277    or RAP domain-containing proteins [where RAP stands for **R**NA binding domain abundant in

278    **ap**icomplexans, (37)] stood out.

279    These RNA-binding proteins encompass, in addition to the crTP, from N- to C-terminus

280    a variable region of 0-320 aa followed by a ~105 aa long conserved region (CR1), 2-13 repeats

281    of a degenerate 38 aa motif with the most conserved residues being

282    (xxxPxxxxLxxxxxxxxxxxxxFxxQxxxxxLNAxAKL), often followed by a 110 aa long conserved

283    sequence (CR2), and the 60 aa long RAP domain (**Fig. 6**). This domain organization resembles

284    the one of organelle-targeted octotrico peptide repeat (OPR; i.e. 38 aa peptide repeat) gene

285    expression regulators in green algae and plants (**Fig. 6B, D**) and repeat-containing T3SS

11

286  effector proteins described from symbiotic or pathogenic bacteria (**Fig. 6B, E, F**). All repeat

287  motifs share the prediction to fold into two antiparallel α-helices. 3D-structure prediction

288  suggests folding of the α-helical repeats in *Paulinella* OPR proteins into a super helix (or α-

289  solenoid) structure (**Fig. 6G**).

290      Screening the complete *P. chromatophora* transcriptome identified OPR proteins as

291  part of an expanded protein family containing at least 101 members with 1-13 individual OPR

292  motifs (**Table S5**). Besides the 12 chromatophore-localized OPR proteins identified by MS (**Fig.**

293  **6A**), of the further 12 OPR proteins identified only in the transcriptome for which full-length

294  N-terminal sequence information was available, seven proteins contained a crTP (**Fig. 6A**), the

295  remaining five a mitochondrial targeting signal.

296

297  **Discussion**

298  **Metabolite transport**

299  Despite the obvious need for extensive metabolite exchange between the chromatophore and

300  cytoplasm (12), the chromatophore likely lost on the order of 70 solute transporters following

301  symbiosis establishment (**Fig. 1**). The remaining transport systems do not appear apt to

302  establish metabolic connectivity (**Fig. 2D**). Solely two systems, a DME family and a DASS family

303  transporter, might be involved in metabolite transport. Furthermore, there are three ABC-

304  transporters for which substrate specificity is unknown. However, the high energy costs

305  associated with their ATP-consuming primary active mode of transport appears to be

306  incongruous with high-throughput metabolite shuttling. Some of these ABC-transporters

307  might have become specialized for protein import instead. In line with this idea, the ABC-half

308  transporter PCC0669 that showed highest ion intensities among all chromatophore-encoded

309  transporters (**Fig. 2D**), possesses 33% similarity to *Bradyrhizobium* BclA that functions as an

310  importer for nodule-specific cysteine-rich (NCR) peptides produced by the host plant (38).

311  However, since other transporters in the same family are involved in peroxisomal transport of

312  fatty acids or fatty acyl-CoA (39), similar substrates could also be transported by PCC0669.

313      In plastids, insertion of nucleus-encoded transporters into the IM is crucial for

314  metabolic connectivity (15-17). Also in more recently established endosymbiotic associations,

315    such as plant sap-feeding insects with nutritional bacterial endosymbionts, multiplication of

316    host transporters followed by their recruitment to the host/endosymbiont interface

317    apparently was involved in establishing metabolic connectivity (40, 41). However, these

318    transporters localize to the symbiosomal membrane, a host membrane that surrounds

319    bacterial endosymbionts. The mechanism enabling metabolite transport across the

320    symbionts' IM and OM, with symbiont-encoded transport systems being scarce, is a

321    longstanding, unanswered question (42).

322    Despite the import of hundreds of soluble proteins into the chromatophore, our work

323    provided no evidence for the insertion of nucleus-encoded transporters into the

324    chromatophore IM (or thylakoids). The possibility that such proteins escaped detection for

325    technical reasons appears improbable because: (i) 72% of the chromatophore-encoded

326    transporters were identified in CM samples. Assuming comparable abundances for nucleus-

327    encoded chromatophore-targeted transporters, a large percentage of these proteins should

328    have been detected, too. (ii) More than 100 nucleus-encoded transporters or transporter

329    components were detected in comparable amounts of PM samples showing that our method

330    is feasible to detect this group of proteins. (iii) IM transporters were repeatedly identified in

331    comparable analyses of cyanobacterial (43-47) or plastidial membrane fractions (48-50). Thus,

332    a mechanism to insert nucleus-encoded multi-spanning TM proteins into chromatophore IM

333    and thylakoids likely has not evolved (yet) in *P. chromatophora*.

334    The protein composition of the chromatophore OM is currently unclear. However, its

335    putative host origin and the notion that proteins traffic into the chromatophore likely via the

336    Golgi (11) suggest that nucleus-encoded transporters can be targeted to the OM by vesicle

337    fusion. Nonetheless, our findings spotlight the puzzling absence of suitable transporters that

338    would allow metabolite exchange across the chromatophore IM. The conservation of active

339    and secondary active IM transporters on the chromatophore genome (**Fig. 2D**) strongly implies

340    that the chromatophore IM kept its barrier function and there is an electrochemical gradient

341    across this lipid bilayer.

342    In contrast to the absence of multi-spanning TM proteins, we identified numerous

343    short single-spanning TM and AMP-like orphan proteins among chromatophore-targeted

344    proteins. These short import candidates fall into at least four expanded groups, suggesting

345    some degree of functional specialization. Interestingly, expanded arsenals of symbiont-

346 targeted polypeptides convergently evolved in many taxonomically unrelated symbiotic

347 associations and thus seems to represent a powerful strategy to establish host control over

348 bacterial endosymbionts (51). It has been suggested that these 'symbiotic AMPs' have the

349 ability to self-translocate across or self-insert into endosymbiont membranes and mediate

350 control over various biological processes in the symbionts including translation, septum

351 formation or modulation of membrane permeability and metabolite exchange (42, 51-56).

352      The discovery of TMH-containing group 1 proteins appears to be of particular interest

353 in the context of metabolite exchange. The frequent occurrence of (small)xxx(small) motifs

354 might indicate the potential of these proteins to oligomerize (36, 57). The predicted pore-

355 lining residues (**Table S4**) in many of these proteins further suggest that they could form

356 homo- or hetero-oligomeric channels. It has been previously reported that AMPs can arrange

357 in channel-like assemblies which facilitate diffusion along concentration gradients (58, 59),

358 though the lifetime and selectivity of such arrangements requires further investigation. Given

359 the size of the metabolites to be transported, they would be required to form multimer

360 arrangements in barrel-stave (**Fig. S4**), or shortly lived toroidal pores, while maintaining the

361 overall impermeability of the membrane. The formation of such pores still begs the question

362 of how they could maintain a selective metabolite transport. An interesting example in that

363 respect is the VDAC channel of the mitochondrial OM which has been described to follow a

364 stochastic gating mechanism, in which only bigger and, hence, slowly diffusing molecules

365 would be allowed to permeate (60).

366      An alternative mode of action involves soluble, short import candidates which could

367 interact with the chromatophore envelope membranes via stretches of positively charged aa

368 and amphipathic helices (**Fig. 4A**), and putatively modulate their permeability (42) in what is

369 known as carpet model (61). The mechanism by which such an interaction could cause a

370 transient permeabilization is still a matter of debate, although the asymmetric distribution on

371 the membrane bilayer has been pointed out as plausible reason (62). Further experimental

372 work with the identified proteins could shed light on the potential transport mechanism.

373      Other short import candidates might also attack intracellular targets. The group 2

374 sequence motif is found also in hypothetical bacterial proteins which include domains related

375 to DNA processing functions (**Fig. 4B**). Thus, group 2 proteins might provide the host with

376 control over aspects of genetic information processing in the chromatophore. The presence

14

377 of dozens to hundreds of similar proteins in the various groups, points to a functional

378 interdependence or reciprocal control of individual peptides. In insects, co-occurring AMPs

379 have been shown to synergize, e.g. some AMPs permeabilize membranes to enable entry of

380 other AMPs that have intracellular targets (63).

381

**Nuclear control over expression of chromatophore-encoded proteins**

383 Besides the establishment of metabolic connectivity, our analyses illuminated another

384 cornerstone in organellogenesis, the evolution of nuclear control over organellar gene

385 expression. Previously, we identified a large number of proteins annotated as transcription

386 factors among chromatophore-targeted proteins (10). Here we described a novel class of

387 chromatophore-targeted helical repeat proteins. Helical repeat proteins appear to represent

388 ubiquitous nuclear factors involved in regulation of organellar gene expression. These proteins

389 are generally characterized by the presence of degenerate 30-40 aa repeat motifs, each of

390 them containing two antiparallel α-helices. The succession of motifs underpins the formation

391 of a super helix that enables sequence specific binding to nucleic acids.

392 The *P. chromatophora* nuclear genome encodes at least 101 OPR helical repeat

393 proteins (**Fig. 6C**). OPR proteins have mostly been studied in the green alga *C. reinhardtii*,

394 where 44 OPR genes were identified in the nuclear genome. Almost all of these OPR proteins

395 are predicted to localize to organelles [(64); **Fig. 6D**] and five have been shown experimentally

396 to be involved in post-transcriptional steps of chloroplast gene expression. The only known

397 *Arabidopsis* OPR protein is AtRAP [(65); **Fig. 6B**] a factor promoting chloroplast rRNA

398 maturation.

399 Also the *Paulinella* OPR proteins seem to be mostly organelle-targeted. Many

400 *Paulinella* OPR proteins possess, in addition to the OPR stretches, a Fas-activated

401 serine/threonine (FAST) kinase-like domain (66) and a C-terminal RAP domain (**Fig. 6A**). This

402 domain combination is also present in some of the *Chlamydomonas* OPR proteins (e.g. CrRAP

403 in **Fig. 6B**), the *Arabidopsis* AtRAP protein (**Fig. 6B**), and the FASTK family of vertebrate

404 nucleus-encoded regulators of mitochondrial gene expression (67). Additionally, some

405 bacterial T3SS effector proteins (**Fig. 6B**) show similar domain architectures. However, the

406 exact molecular functions of FAST kinase-like and RAP domains as well as the two conserved

15

407 regions in *Paulinella* OPR proteins (CR1 and CR2) that share no similarity with known domains

408 remain unknown.

409       In conclusion, in parallel to the evolution of mitochondria and plastids, also during

410 chromatophore evolution an expanded family of chromatophore-targeted helical repeat

411 proteins evolved. Based on the similarity of their domain architecture to known organelle-

412 targeted expression regulators, the OPR proteins in *P. chromatophora* likely serve as nuclear

413 factors modulating chromatophore gene expression by direct binding to specific target RNAs.

414 Probably chromatophore-targeted OPR proteins evolved from pre-existing mitochondrial

415 expression regulators and were recruited to the chromatophore by crTP acquisition. However,

416 the RNA-binding ability of *Paulinella* OPR proteins, their specific target sequences as well as

417 their ability to modulate expression of chromatophore-encoded proteins remain to be tested

418 experimentally.

419

420 **Materials and Methods**

421 **Cultivation of *P. chromatophora* and chromatophore isolation**

422 *P. chromatophora* CCAC0185 [axenic version (7)] was grown (11) and chromatophores isolated

423 as described previously (10). In brief, *P. chromatophora* cells were washed three times with

424 isolation buffer (IB: 50 mM HEPES pH 7.5, 2 mM EGTA, 2 mM $MgCl_2$, 250 mM sucrose, 125

425 mM NaCl) and depleted of dead cells on a discontinuous 20-80% Percoll gradient. The resulting

426 pellet of intact cells was resuspended in IB, cells were broken in a cell disruptor (Constant

427 Systems) at 0.5 kbar, and intact chromatophores isolated on another discontinuous 20-80%

428 Percoll gradient. To increase purity, isolated chromatophores were re-isolated from a fresh

429 Percoll gradient. Recovered chromatophores were washed three times in IB, supplemented

430 with protease inhibitor cocktail (Roche cOmplete), frozen in liquid nitrogen, and stored at -

431 80°C until further use.

432

433 **Transmission electron microscopy (TEM)**

434 Isolated chromatophores were fixed in IB containing 1.25% glutaraldehyde for 45 min on ice

435 followed by 30 min post-fixation in 1% $OsO_4$ in IB at room temperature. Fixed chromatophores

436 were washed, mixed with 14.5% (w/v) BSA, pelleted, and the pellet fixed with 2.5%

437 glutaraldehyde for 20 min at room temperature. The fixed pellet was dehydrated in rising

438 concentrations of ethanol (from 60% to 100% at -20°C) and then infiltrated with Epon using

439 propylene oxide as a transition solvent. Epon was polymerized at 60°C for 24 h. 70 nm ultrathin

440 sections were prepared and contrasted with uranyl acetate and lead citrate according to (68).

441 A Hitachi H7100 TEM (Hitachi, Tokyo, Japan) with Morada camera (EMSIS GmbH, Münster,

442 Germany) operated at 100 kV was used for TEM analysis. Essentially the same protocol was

443 used for intact *P. chromatophora* cells, however, IB was replaced by growth medium [WARIS-

444 H (69) supplemented with 1.5 mM $Na_2SiO_3$].

445

446 **Protein fractionation**

447 *CM and PM samples:* Isolated chromatophores or *P. chromatophora* cells were washed with

448 Buffer I (50 mM HEPES pH 7.5, 125 mM NaCl, 0.5 mM EDTA) at 20,000 x g or 200 x g,

449 respectively. Pellets were resuspended in Buffer I and broken by two passages in a cell

450 disrupter at 2.4 kbar. Lysates were supplemented with 500 mM NaCl (final concentration) and

451 passed five times through a 0.6 mm cannula. Cell debris was removed by two successive

452 centrifugations at 15,500 x g. The supernatant was subjected to ultracentrifugation for 1 h at

453 150,000 x g (Beckmann L-80XL optima ultracentrifuge, Rotor 70.1 Ti at 50,000 rpm). Pellets

454 were resuspended in 100 mM $Na_2CO_3$ pH >11 and incubated for 1 h intermitted by 15 passes

455 through a 0.6 mm cannula. Then, insoluble proteins were collected by ultracentrifugation, and

456 subsequently washed with Buffer II (10 mM Tris-HCl pH 7.5, 150 mM NaCl, 0.5 mM EDTA) by

457 passage through a cannula until no particles were visible. Finally, the insoluble fraction was

458 pelleted by ultracentrifugation and solubilized at 36°C in Buffer II supplemented with 1%

459 TritonX-100, 1% Na-deoxycholate, and 0.1% SDS.

460 *CL samples:* Protein was extracted from intact isolated chromatophores by

461 precipitation with 10% trichloracetic acid for 30 min on ice and pelleted at 21,000 x g for 20

462 min. Pellets were washed twice with ice cold acetone for 10 min and finally resuspended in

463 Buffer II plus detergents.

464    Protein concentration was determined in a Neuhoff assay (70). Aliquots were

465    supplemented with SDS sample buffer (final conc. 35 mM Tris-HCl pH 7.0, 7.5% Glycerol, 3%

466    SDS, 150 mM DTT, Bromophenol blue), frozen in liquid nitrogen, and stored at -80°C until MS-

467    analysis. All steps were performed at 4°C, protease inhibitor cocktail (Roche cOmplete) was

468    added to all buffers used.

469

470    **MS analysis and protein identification**

471    Sample preparation and subsequent MS/MS analysis of 3 independent preparations of

472    CM, PM, and CL samples was essentially carried out as described (10). Briefly, proteins were

473    in-gel digested in (per sample) 0.1 μg trypsin in 10 mM ammonium hydrogen carbonate

474    overnight at 37°C and resulting peptides resuspended in 0.1% trifluoroacetic acid. Two

475    independent MS analyses were performed. In MS experiment 1, 500 ng protein per sample,

476    and in MS experiment 2, 500 ng protein per lysate and 1.5 μg protein per membrane sample

477    was analyzed. Peptides were separated on C18 material by LC, injected into a QExactive plus

478    mass spectrometer, and the mass spectrometer was operated as described (10). Raw files

479    were further processed with MaxQuant (MPI for Biochemistry, Planegg, Germany) for protein

480    identification and quantification using standard parameters. MaxQuant 1.6.2.10 was used for

481    the MS experiment 1 analysis and MaxQuant 1.6.3.4 for MS experiment 2. Searches were

482    carried out using 60,108 sequences translated from a *P. chromatophora* transcriptome and

483    the 867 translated genes predicted on the chromatophore genome (10). Peptides and proteins

484    were accepted at a false discovery rate of 1%. Proteomic data have been deposited to the

485    ProteomeXchange Consortium via the PRIDE (71) partner repository with the dataset

486    identifier PXD021087.

487

488    **Protein enrichment analysis**

489    Intensities of individual proteins were normalized by division of individual intensities in each

490    replicate by the sum of intensities of all proteins identified with ≥2 peptides in the same

491    replicate. Each protein was assigned an intensity level representing its log10 transformed

492    mean    normalized    intensity    from    three    replicates    in    either    fraction    added    7

493     $(log10\ (\overline{normInt}) + 7)$, enabling a simple ranking of intensities in a logarithmic range from 0

494     to 6.

495          The enrichment factor for each protein in CM as compared to PM or CL samples ($E_{CM/PM}$

496     or $E_{CM/CL}$, respectively) was calculated as $E_{CM/PM} = \overline{normInt_{CM}}/\overline{normInt_{PM}}$ or $E_{CM/CL}$

497     $=\overline{normInt_{CM}}/\overline{normInt_{CL}}$ [**Table S2;** missing values (intensity = 0) were excluded from the

498     calculation of means]. Proteins with ≥3 SpC in the chromatophore (i.e. CM+CL fractions) and

499     either $E_{CM/PM}$>1.5 in at least one of two MS experiments or 0.5<$E_{CM/PM}$<1.5 in both MS

500     experiments were considered as enriched in chromatophores (see **Fig. S2**). Correspondingly,

501     $E_{CM/CL}$>1 indicate protein enrichment, $E_{CM/CL}$<1 depletion in CM samples.

502          Furthermore, a statistic approach was applied to visualize differences between

503     proteins enriched or exclusively found in a certain fraction. In pairwise comparisons, only

504     proteins were considered showing valid *normInt* values in all three replicates of at least one

505     of the samples being compared. *NormInt* values were log2 transformed and missing values

506     imputed by values from a down shifted normal distribution (width 0.3 SD, down shift 1.8 SD)

507     followed by a pairwise sample comparison based on Student's t-tests and the significance

508     analysis of microarrays algorithm ($S_0$ = 0.8, FDR 5%) (72). Differences between individual

509     proteins in CM vs. PM or CM vs. CL samples were calculated as $\overline{log2(normInt_{CM})} -$

510     $\overline{log2(normInt_{PM})}$ or $\overline{log2(normInt_{CM})} - \overline{log2(normInt_{CL})}$, respectively.

511

**512 Sequence and structural bioinformatics analyses**

513     TMHs were predicted with TMHMM 2.0 (73), pore-lining residues were predicted with

514     MEMSAT-SVM-pore (74), and AMP peptides were predicted with AmpGram (75). Sequence

515     motifs were discovered using MEME 5.0.5 algorithm (76) in classic mode and visualized with

516     WebLogo (77), number and position of motifs in protein sequences were determined with

517     MAST 5.0.5 using default settings (78). The *P. chromatophora* transcriptome was screened for

518     (i) conserved motifs shown in **Fig. 4A** group 2 and 3 and (ii) the degenerate 38 aa motif shown

519     in **Fig. 6C** using FIMO 5.0.5 with default settings (79). Proteins that contain at least 5 repeats

520     of the 38 aa motif with a p-value < $e^{-10}$ and / or at least 1 repeat with a p-value < $e^{-20}$ were

521     considered candidate OPR-proteins. α-helices in repetitive elements or AMP-like proteins

19

522  were predicted with Jpred4 (80) and NetSurfP-2.0 (81), respectively. Helical wheel projections

523  were created with HeliQuest (82). Functional protein domains were found with DELTA-BLAST

524  (83). Targeting signals were predicted with PredAlgo (84) for CrRAP and CrTab1, and TargetP

525  2.0 (85), WoLFPSORT (86), and Predotar (87) for *P. chromatophora* proteins. Tertiary structure

526  predictions were obtained using Phyre2 (88) in normal mode. Area-proportional Venn

527  diagrams were calculated with eulerAPE (89).

528      Transporters were classified according to the Transporter Classification Database (90).

529  Complete lists of the transporters depicted in Figures 1 and 2D and methods for their

530  identification and classification are provided in Table S1. No OM porins could be identified in

531  the chromatophore genome based on sequence similarity or topology predictions using

532  MCMBB (91).

533

534  **Acknowledgments**

538

539  **References**

540  1.  **Nowack ECM.** 2014. *Paulinella chromatophora* - rethinking the transition from endosymbiont
541      to organelle. Acta Soc Bot Pol **83:**387-397.
542  2.  **Gabr A, Grossman AR, Bhattacharya D.** 2020. *Paulinella*, a model for understanding plastid
543      primary endosymbiosis. J Phycol **56:**837-843.
544  3.  **Delaye L, Valadez Cano C, Pérez Zamorano B.** 2016. How really ancient is *Paulinella
545      chromatophora*?. PLOS Currents Tree of Life Ed 1.
546  4.  **Marin B, Nowack ECM, Melkonian M.** 2005. A plastid in the making: Evidence for a second
547      primary endosymbiosis. Protist **156:**425-432.
548  5.  **Nowack ECM, Melkonian M, Glöckner G.** 2008. Chromatophore genome sequence of
549      *Paulinella* sheds light on acquisition of photosynthesis by eukaryotes. Curr Biol **18:**410-418.
550  6.  **Reyes-Prieto A, Yoon HS, Moustafa A, Yang EC, Andersen RA, Boo SM, Nakayama T, Ishida K,
551      Bhattacharya D.** 2010. Differential gene retention in plastids of common recent origin. Mol
552      Biol Evol **27:**1530-1537.
553  7.  **Nowack ECM, Price DC, Bhattacharya D, Singer A, Melkonian M, Grossman AR.** 2016. Gene
554      transfers from diverse bacteria compensate for reductive genome evolution in the
555      chromatophore of *Paulinella chromatophora*. Proc Nati Acad Sci USA **113:**12214-12219.

8.  **Nowack ECM, Vogel H, Groth M, Grossman AR, Melkonian M, Glöckner G.** 2011. Endosymbiotic gene transfer and transcriptional regulation of transferred genes in *Paulinella chromatophora*. Mol Biol Evol **28:**407-422.

9.  **Zhang R, Nowack ECM, Price DC, Bhattacharya D, Grossman AR.** 2017. Impact of light intensity and quality on chromatophore and nuclear gene expression in *Paulinella chromatophora*, an amoeba with nascent photosynthetic organelles. Plant J **90:**221–234.

10. **Singer A, Poschmann G, Mühlich C, Valadez-Cano C, Hänsch S, Rensing SA, Stühler K, Nowack ECM.** 2017. Massive protein import into the early evolutionary stage photosynthetic organelle of the amoeba *Paulinella chromatophora*. Curr Biol **27:**2763-2773.

11. **Nowack ECM, Grossman AR.** 2012. Trafficking of protein into the recently established photosynthetic organelles of *Paulinella chromatophora*. Proc Nati Acad Sci USA **109:**5340-5345.

12. **Valadez-Cano C, Olivares-Hernandez R, Resendis-Antonio O, DeLuna A, Delaye L.** 2017. Natural selection drove metabolic specialization of the chromatophore in *Paulinella chromatophora*. BMC Evol Biol **17:**99.

13. **Sato N, Yoshitomi T, Mori-Moriyama N.** 2020. Characterization and biosynthesis of lipids in *Paulinella micropora* MYN1: Evidence for efficient integration of chromatophores into cellular lipid metabolism. Plant Cell Phys **61:**869-881.

14. **Kies L.** 1974. Electron microscopical investigations on *Paulinella chromatophora* Lauterborn, a thecamoeba containing blue-green endosymbionts (cyanelles). Protoplasma **80:**69-89.

15. **Facchinelli F, Weber APM.** 2011. The metabolite transporters of the plastid envelope: An update. Front Plant Sci **2:**50.

16. **Karkar S, Facchinelli F, Price DC, Weber APM, Bhattacharya D.** 2015. Metabolic connectivity as a driver of host and endosymbiont integration. Proc Nati Acad Sci USA **112:**10208-10215.

17. **Fischer K.** 2011. The import and export business in plastids: Transport processes across the inner envelope membrane. Plant Physiology **155:**1511-1519.

18. **Breuers FKH, Bräutigam A, Weber APM.** 2011. The plastid outer envelope - a highly dynamic interface between plastid and cytoplasm. Front Plant Sci **2:**97.

19. **Hammani K, Bonnard G, Bouchoucha A, Gobert A, Pinker F, Salinas T, Giegé P.** 2014. Helical repeats modular proteins are major players for organelle gene expression. Biochimie **100:**141-150.

20. **Woodson JD, Chory J.** 2008. Coordination of gene expression between organellar and nuclear genomes. Nature Rev Genetics **9:**383-395.

21. **Marin B, Nowack ECM, Glöckner G, Melkonian M.** 2007. The ancestor of the *Paulinella* chromatophore obtained a carboxysomal operon by horizontal gene transfer from a *Nitrococcus*-like gamma-proteobacterium. BMC Evol Biol **7:**85.

22. **Paulsen IT, Nguyen L, Sliwinski MK, Rabus R, Saier Jr MH.** 2000. Microbial genome analyses: Comparative transport capabilities in eighteen prokaryotes. J Mol Biol **301:**75-100.

23. **Omata T, Andriesse X, Hirano A.** 1993. Identification and characterization of a gene-cluster involved in nitrate transport in the cyanobacterium *Synechococcus* sp. PCC7942. Mol General Genetics **236:**193-202.

24. **Montesinos ML, Muro-Pastor AM, Herrero A, Flores E.** 1998. Ammonium/Methylammonium permeases of a cyanobacterium - Identification and analysis of three nitrogen-regulated amt genes in *Synechocystis* sp. PCC 6803. J Biol Chem **273:**31463-31470.

25. **Valladares A, Montesinos ML, Herrero A, Flores E.** 2002. An ABC-type, high-affinity urea permease identified in cyanobacteria. Mol Microbiol **43:**703-715.

26. **Quintero MJ, Montesinos ML, Herrero A, Flores E.** 2001. Identification of genes encoding amino acid permeases by inactivation of selected ORFs from the *Synechocystis* genomic sequence. Genome Research **11:**2034-2040.

27. **Laudenbach DE, Grossman AR.** 1991. Characterization and mutagenesis of sulfur-regulated genes in a cyanobacterium - evidence for function in sulfate transport. J Bacteriol **173:**2739-2750.
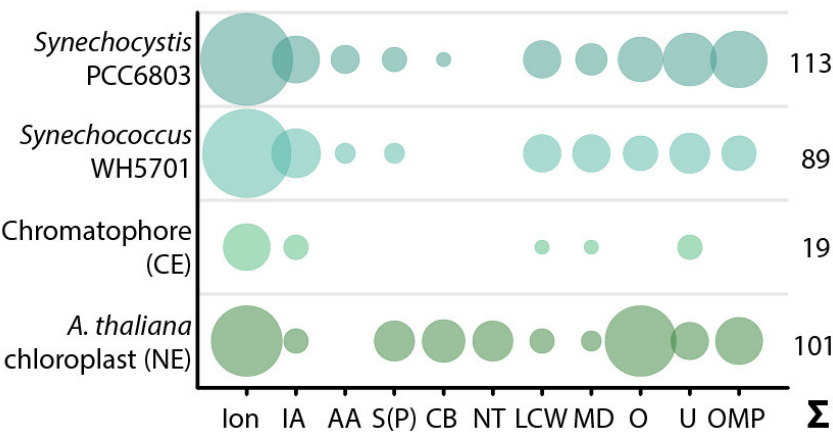
608    28.    **Jack DL, Yang NM, Saier MH.** 2001. The drug/metabolite transporter superfamily. European J
609           Biochem **268:**3620-3639.
610    29.    **Markovich D.** 2012. Sodium-Sulfate/Carboxylate Cotransporters (SLC13), p 239-256. *In*
611           Bevensee MO (ed), Co-Transport Systems, vol 70.
612    30.    **Marchand J, Heydarizadeh P, Schoefs B, Spetea C.** 2018. Ion and metabolite transport in the
613           chloroplast of algae: lessons from land plants. Cell Mol Life Sci **75:**2153-2176.
614    31.    **Mehrshahi P, Stefano G, Andaloro JM, Brandizzi F, Froehlich JE, DellaPenna D.** 2013.
615           Transorganellar complementation redefines the biochemical continuity of endoplasmic
616           reticulum and chloroplasts. Proc Nati Acad Sci USA **110:**12126-12131.
617    32.    **Weber APM, Schwacke R, Flügge UI.** 2005. Solute transporters of the plastid envelope
618           membrane. Annu Rev Plant Biol vol 56, p 133-164.
619    33.    **Goetze TA, Patil M, Jeshen I, Bölter B, Grahl S, Soll J.** 2015. Oep23 forms an ion channel in the
620           chloroplast outer envelope. BMC Plant Biol **15:**47.
621    34.    **Harsman A, Schock A, Hemmis B, Wahl V, Jeshen I, Bartsch P, Schlereth A, Pertl-Obermeyer
622           H, Goetze TA, Soll J, Philippar K, Wagner R.** 2016. OEP40, a regulated glucose-permeable β-
623           barrel solute channel in the chloroplast outer envelope membrane. J Biol Chem **291:**17848-
624           17860.
625    35.    **Wang Z, Anderson NS, Benning C.** 2013. The phosphatidic acid binding site of the *Arabidopsis*
626           trigalactosyldiacylglycerol 4 (TGD4) protein required for lipid import into chloroplasts. J Biol
627           Chem **288:**4763-4771.
628    36.    **Teese MG, Langosch D.** 2015. Role of GxxxG motifs in transmembrane domain interactions.
629           Biochemistry **54:**5125-5135.
630    37.    **Lee I, Hong W.** 2004. RAP - A putative RNA-binding domain. Trends Biochem Sci **29:**567-570.
631    38.    **Guefrachi I, Pierre O, Timchenko T, Alunni B, Barrière Q, Czernic P, Villaecija-Aguilar JA, Verly
632           C, Bourge M, Fardoux J, Mars M, Kondorosi E, Giraud E, Mergaert P.** 2015. *Bradyrhizobium*
633           BclA is a peptide transporter required for bacterial differentiation in symbiosis with
634           *Aeschynomene* legumes. Mol Plant-Microbe Interactions **28:**1155-1166.
635    39.    **Linka N, Esser C.** 2012. Transport proteins regulate the flux of metabolites and cofactors across
636           the membrane of plant peroxisomes. Front Plant Sci **3:**3.
637    40.    **Duncan RP, Husnik F, Van Leuven JT, Gilbert DG, Dávalos LM, McCutcheon JP, Wilson ACC.**
638           2014. Dynamic recruitment of amino acid transporters to the insect/symbiont interface. Mol
639           Ecol **23:**1608-1623.
640    41.    **Price DRG, Duncan RP, Shigenobu S, Wilson ACC.** 2011. Genome expansion and differential
641           expression of amino acid transporters at the aphid/*Buchnera* symbiotic interface. Mol Biol Evol
642           **28:**3113-3126.
643    42.    **Mergaert P, Kikuchi Y, Shigenobu S, Nowack ECM.** 2017. Metabolic integration of bacterial
644           endosymbionts through antimicrobial peptides. Trends Microbiol **25:**703-712.
645    43.    **Baers LL, Breckels LM, Mills LA, Gatto L, Deery MJ, Stevens TJ, Howe CJ, Lilley KS, Lea-Smith
646           DJ.** 2019. Proteome mapping of a cyanobacterium reveals distinct compartment organization
647           and cell-dispersed metabolism. Plant Phys **181:**1721-1738.
648    44.    **Choi J-S, Park YH, Oh JH, Kim S, Kwon J, Choi Y-E.** 2020. Efficient profiling of detergent-assisted
649           membrane proteome in cyanobacteria. J Applied Phycol. **32:**1177-1184.
650    45.    **Liberton M, Saha R, Jacobs JM, Nguyen AY, Gritsenko MA, Smith RD, Koppenaal DW, Pakrasi
651           HB.** 2016. Global proteomic analysis reveals an exclusive role of thylakoid membranes in
652           bioenergetics of a model cyanobacterium. Mol Cell Proteomics **15:**2021-2032.
653    46.    **Pisareva T, Kwon J, Oh J, Kim S, Ge C, Wieslander A, Choi J-S, Norling B.** 2011. Model for
654           membrane organization and protein sorting in the cyanobacterium *Synechocystis* sp. PCC 6803
655           inferred from proteomics and multivariate sequence analyses. J Proteome Research **10:**3617-
656           3631.
657    47.    **Plohnke N, Seidel T, Kahmann U, Rögner M, Schneider D, Rexroth S.** 2015. The proteome and
658           lipidome of *Synechocystis* sp. PCC 6803 cells grown under light-activated heterotrophic
659           conditions. Mol Cell Proteomics **14:**572-584.

48. **Bouchnak I, Brugière S, Moyet L, Le Gall S, Salvi D, Kuntz M, Tardif M, Rolland N.** 2019. Unraveling hidden components of the chloroplast envelope proteome: Opportunities and limits of better MS sensitivity. Mol Cell Proteomics **18:**1285-1306.

49. **Bräutigam A, Hoffmann-Benning S, Weber APM.** 2008. Comparative proteomics of chloroplast envelopes from $C_3$ and $C_4$ plants reveals specific adaptations of the plastid envelope to $C_4$ photosynthesis and candidate proteins required for maintaining $C_4$ metabolite fluxes. Plant Phys **148:**568-579.

50. **Simm S, Papasotiriou DG, Ibrahim M, Leisegang MS, Mueller B, Schorge T, Karas M, Mirus O, Sommer MS, Schleiff E.** 2013. Defining the core proteome of the chloroplast envelope membranes. Front Plant Sci **4:**11.

51. **Mergaert P.** 2018. Role of antimicrobial peptides in controlling symbiotic bacterial populations. Natural Product Rep **35:**336-356.

52. **Carro L, Pujic P, Alloisio N, Fournier P, Boubakri H, Hay AE, Poly F, François P, Hocher V, Mergaert P, Balmand S, Rey M, Heddi A, Normand P.** 2015. *Alnus* peptides modify membrane porosity and induce the release of nitrogen-rich metabolites from nitrogen-fixing *Frankia*. ISME J **9:**1723-1733.

53. **Farkas A, Maróti G, Dürgo H, Györgypál Z, Lima RM, Medzihradszky KF, Kereszt A, Mergaert P, Kondorosi E.** 2014. *Medicago truncatula* symbiotic peptide NCR247 contributes to bacteroid differentiation through multiple mechanisms. Proc Nati Acad Sci USA **111:**5183-5188.

54. **Login FH, Balmand S, Vallier A, Vincent-Monégat C, Vigneron A, Weiss-Gayet M, Rochat D, Heddi A.** 2011. Antimicrobial peptides keep insect endosymbionts under control. Science **334:**362-365.

55. **Mergaert P, Uchiumi T, Alunni B, Evanno G, Cheron A, Catrice O, Mausset AE, Barloy-Hubler F, Galibert F, Kondorosi A, Kondorosi E.** 2006. Eukaryotic control on bacterial cell cycle and differentiation in the Rhizobium-legume symbiosis. Proc Nati Acad Sci USA **103:**5230-5235.

56. **van de Velde W, Zehirov G, Szatmari A, Debreczeny M, Ishihara H, Kevei Z, Farkas A, Mikulass K, Nagy A, Tiricz H, Satiat-Jeunemaître B, Alunni B, Bourge M, Kucho KI, Abe M, Kereszt A, Maroti G, Uchiumi T, Kondorosi E, Mergaert P.** 2010. Plant peptides govern terminal differentiation of bacteria in symbiosis. Science **327:**1122-1126.

57. **Moore DT, Berger BW, DeGrado WF.** 2008. Protein-protein interactions in the membrane: Sequence, structural, and biological motifs. Structure **16:**991-1001.

58. **Rahaman A, Lazaridis T.** 2014. A thermodynamic approach to alamethicin pore formation. Biochim Biophys Acta - Biomembranes **1838:**98-105.

59. **Wang Y, Chen CH, Hu D, Ulmschneider MB, Ulmschneider JP.** 2016. Spontaneous formation of structurally diverse membrane channel architectures from a single antimicrobial peptide. Nature Comm **7:**13535.

60. **Berezhkovskii AM, Bezrukov SM.** 2018. Stochastic gating as a novel mechanism for channel selectivity. Biophys J **114:**1026-1029.

61. **Wimley WC.** 2010. Describing the mechanism of antimicrobial peptide action with the interfacial activity model. ACS Chem Biol **5:**905-917.

62. **Guha S, Ghimire J, Wu E, Wimley WC.** 2019. Mechanistic landscape of membrane-permeabilizing peptides. Chem Rev **119:**6040-6085.

63. **Rahnamaeian M, Cytrynska M, Zdybicka-Barabas A, Dobslaff K, Wiesner J, Twyman RM, Zuchner T, Sadd B, Regoes RR, Schmid-Hempel P, Vilcinskas A.** 2015. Insect antimicrobial peptides show potentiating functional interactions against Gram-negative bacteria. Proc Royal Soc B **282:**20150293.

64. **Eberhard S, Loiselay C, Drapier D, Bujaldon S, Girard-Bascou J, Kuras R, Choquet Y, Wollman FA.** 2011. Dual functions of the nucleus-encoded factor TDA1 in trapping and translation activation of *atpA* transcripts in *Chlamydomonas reinhardtii* chloroplasts. Plant J **67:**1055-1066.
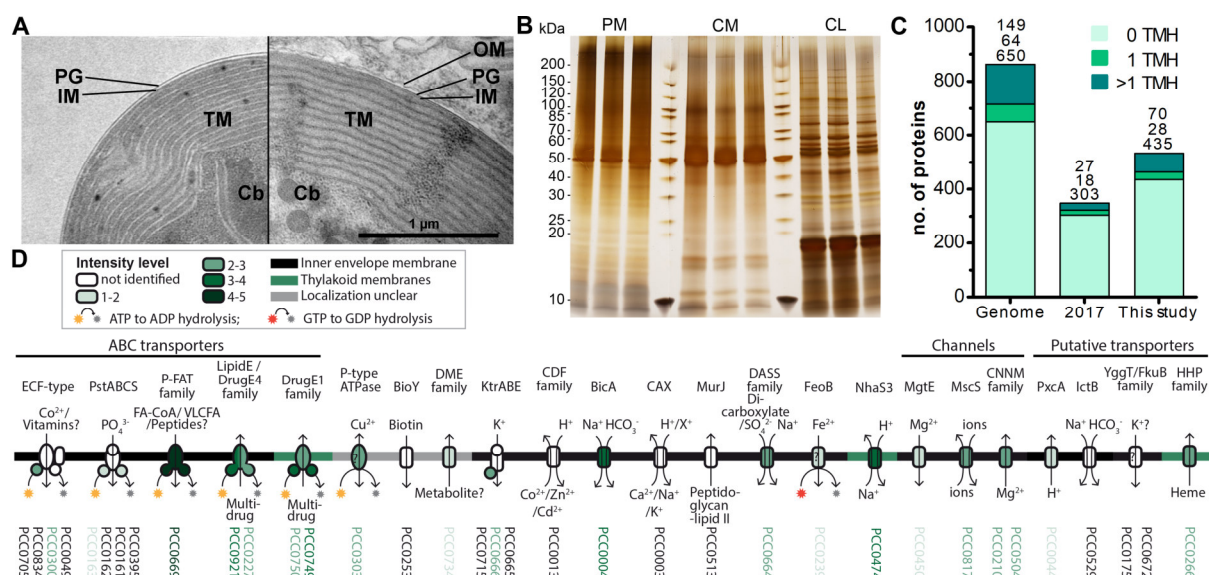
65. **Kleinknecht L, Wang F, Stübe R, Philippar K, Nickelsen J, Bohne AV.** 2014. RAP, the sole octotricopeptide repeat protein in Arabidopsis, is required for chloroplast 16S rRNA maturation. Plant Cell **26:**777-787.

66. **Tian QS, Taupin JL, Elledge S, Robertson M, Anderson P.** 1995. Fas-activated serine threonine kinase (FAST) phosphorylates TIA-1 during Fas-mediated apoptosis. J Exp Med **182:**865-874.

67. **Boehm E, Zornoza M, Jourdain AA, Magdalena AD, García-Consuegra I, Merino RT, Orduna A, Martín MA, Martinou JC, De la Fuente MA, Simarro M.** 2016. Role of FAST kinase domains 3 (FASTKD3) in post-transcriptional regulation of mitochondrial gene expression. J Biol Chem **291:**25877-25887.

68. **Reynolds ES.** 1963. The use of lead citrate at high pH as an electron-opaque stain in electron microscopy. J Cell Biol **17:**208-212.

69. **McFadden GI, Melkonian M.** 1986. Use of hepes buffer for microalgal culture media and fixation for electron-microscopy. Phycologia **25:**551-557.

70. **Neuhoff V, Philipp K, Zimmer HG, Mesecke S.** 1979. A simple, versatile, sensitive and volume-independent method for quantitative protein determination which is independent of other external influences. Hoppe-Seylers Zeitschrift Physiolo Chem **360:**1657-1670.

71. **Perez-Riverol Y, Csordas A, Bai J, Bernal-Llinares M, Hewapathirana S, Kundu DJ, Inuganti A, Griss J, Mayer G, Eisenacher M, Pérez E, Uszkoreit J, Pfeuffer J, Sachsenberg T, Yilmaz Ş, Tiwary S, Cox J, Audain E, Walzer M, Jarnuczak AF, Ternent T, Brazma A, Vizcaíno JA.** 2019. The PRIDE database and related tools and resources in 2019: Improving support for quantification data. Nuc Acids Res **47:**D442-D450.

72. **Tusher VG, Tibshirani R, Chu G.** 2001. Significance analysis of microarrays applied to the ionizing radiation response. Proc Nati Acad Sci USA **98:**5116-5121.

73. **Krogh A, Larsson B, von Heijne G, Sonnhammer ELL.** 2001. Predicting transmembrane protein topology with a hidden Markov model: Application to complete genomes. J Mol Biol **305:**567-580.

74. **Nugent T, Jones DT.** 2012. Detecting pore-lining regions in transmembrane protein sequences. BMC Bioinformatics **13:**169.

75. **Burdukiewicz M, Sidorczuk K, Rafacz D, Pietluch F, Chilimoniuk J, Rödiger S, Gagat P.** 2020. Proteomic screening for prediction and design of antimicrobial peptides with AmpGram. International J Mol Sci **21:**1-13.

76. **Bailey TL, Elkan C.** 1994. Fitting a mixture model by expectation maximization to discover motifs in biopolymers. Proceedings / International Conference on Intelligent Systems for Molecular Biology ; ISMB International Conference on Intelligent Systems for Molecular Biology **2:**28-36.

77. **Crooks GE, Hon G, Chandonia JM, Brenner SE.** 2004. WebLogo: A sequence logo generator. Genome Res **14:**1188-1190.

78. **Bailey TL, Gribskov M.** 1998. Combining evidence using p-values: application to sequence homology searches. Bioinformatics **14:**48-54.

79. **Grant CE, Bailey TL, Noble WS.** 2011. FIMO: scanning for occurrences of a given motif. Bioinformatics **27:**1017-1018.

80. **Drozdetskiy A, Cole C, Procter J, Barton GJ.** 2015. JPred4: a protein secondary structure prediction server. Nuc Acids Res **43:**W389-W394.

81. **Schantz Klausen M, Jespersen MC, Nielsen H, Jensen KK, Jurtz VI, Sonderby CK, Sommer MOA, Winther O, Nielsen M, Petersen B, Marcatili P.** 2019. NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning. Proteins-Structure Function and Bioinformatics **87:**520-527.

82. **Gautier R, Douguet D, Antonny B, Drin G.** 2008. HELIQUEST: A web server to screen sequences with specific α-helical properties. Bioinformatics **24:**2101-2102.

83. **Boratyn GM, Schäffer AA, Agarwala R, Altschul SF, Lipman DJ, Madden TL.** 2012. Domain enhanced lookup time accelerated BLAST. Biol Direct **7:**12.

84. **Tardif M, Atteia A, Specht M, Cogne G, Rolland N, Brugiere S, Hippler M, Ferro M, Bruley C, Peltier G, Vallon O, Cournac L.** 2012. PredAlgo: A New Subcellular Localization Prediction Tool Dedicated to Green Algae. Mol Biol Evol **29:**3625-3639.

85. **Almagro Armenteros JJ, Salvatore M, Emanuelsson O, Winther O, von Heijne G, Elofsson A, Nielsen H.** 2019. Detecting sequence signals in targeting peptides using deep learning. Life Science Alliance **2:**e201900429.

86. **Horton P, Park KJ, Obayashi T, Fujita N, Harada H, Adams-Collier CJ, Nakai K.** 2007. WoLF PSORT: protein localization predictor. Nuc Acids Res **35:**W585-W587.

87. **Small I, Peeters N, Legeai F, Lurin C.** 2004. Predotar: A tool for rapidly screening proteomes for N-terminal targeting sequences. Proteomics **4:**1581-1590.

88. **Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE.** 2015. The Phyre2 web portal for protein modeling, prediction and analysis. Nature Protocols **10:**845-858.

89. **Micallef L, Rodgers P.** 2014. eulerAPE: drawing area-proportional 3-Venn diagrams using ellipses. Plos One **9:** e101717.

90. **Saier MH, Reddy VS, Tsu BV, Ahmed MS, Li C, Moreno-Hagelsieb G.** 2016. The Transporter Classification Database (TCDB): recent advances. Nuc Acids Res **44:**D372-D379.

91. **Bagos PG, Liakopoulos TD, Hamodrakas SJ.** 2004. Finding beta-barrel outer membrane proteins with a markov chain model. WSEAS Transact Biol Biomed **2:**186-189.

92. **Rahire M, Laroche F, Cerutti L, Rochaix JD.** 2012. Identification of an OPR protein involved in the translation initiation of the PsaB subunit of photosystem I. Plant J **72:**652-661.

93. **Mukaihara T, Tamura N.** 2009. Identification of novel *Ralstonia solanacearum* type III effector proteins through translocation analysis of hrpB-regulated gene products. Microbiol-Sgm **155:**2235-2244.

94. **Cline SG, Laughbaum IA, Hamel PP.** 2017. CCS2, an octatricopeptide-repeat protein, is required for plastid cytochrome c assembly in the green Alga *Chlamydomonas reinhardtii*. Front Plant Sci **8:**1306.

95. **Okazaki S, Okabe S, Higashi M, Shimoda Y, Sato S, Tabata S, Hashiguchi M, Akashi R, Gottfert M, Saeki K.** 2010. Identification and functional analysis of type III effector proteins in *Mesorhizobium loti*. Mol Plant-Microbe Interactions **23:**223-234.

96. **Teper D, Burstein D, Salomon D, Gershovitz M, Pupko T, Sessa G.** 2016. Identification of novel *Xanthomonas euvesicatoria* type III effector proteins by a machine-learning approach. Mol Plant Pathol **17:**398-411.
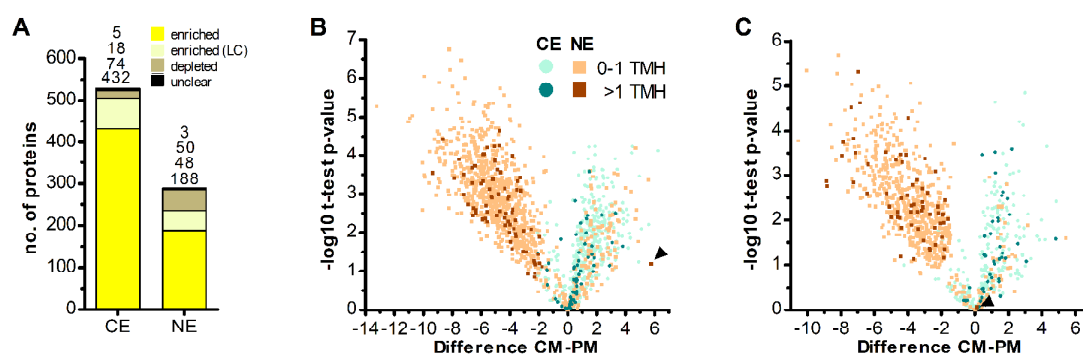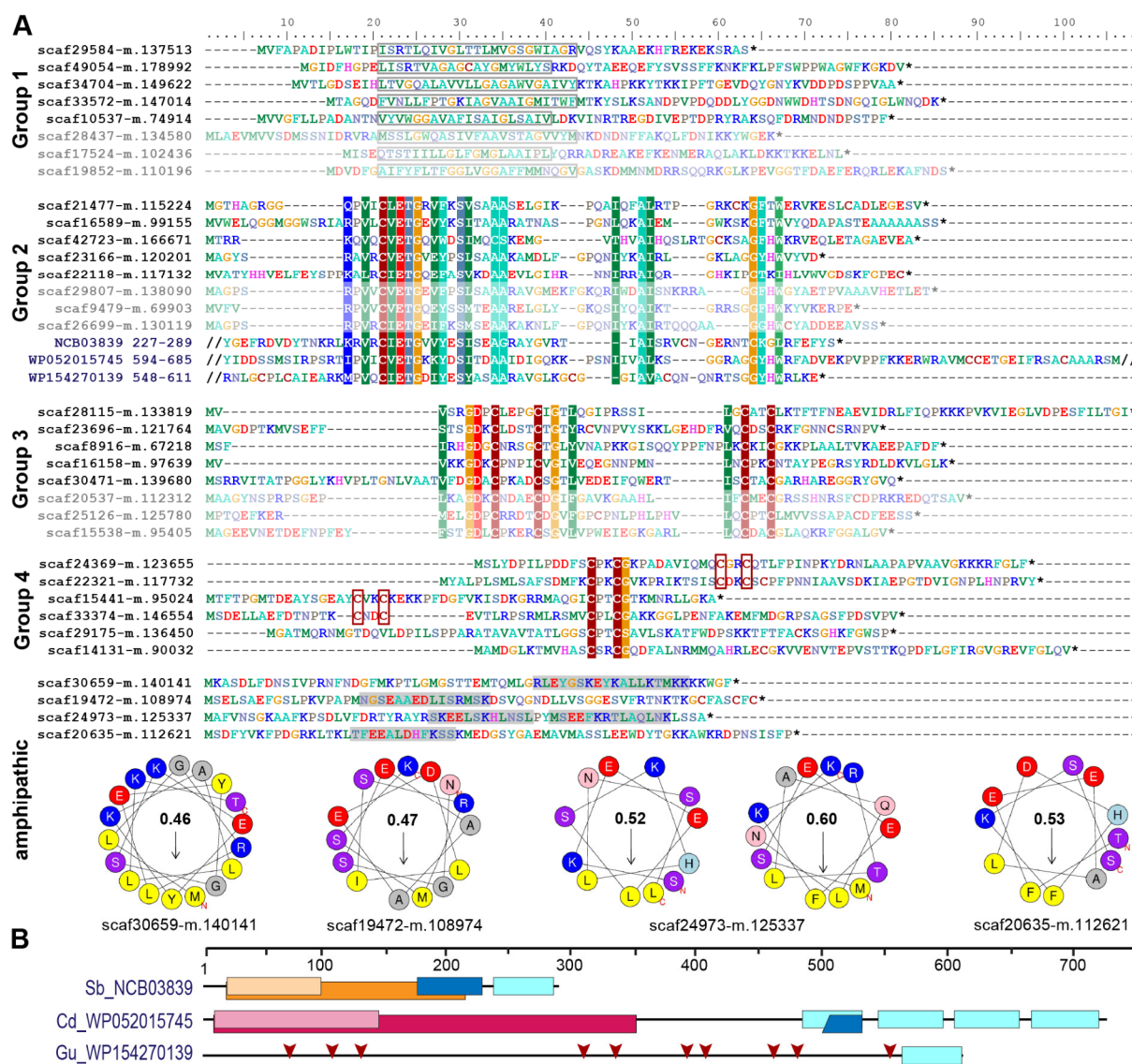
**Figures**



**Figure 1: Predicted solute transport capacities of the chromatophore, *Synechococcus* sp. WH5701, *Synechocystis* sp. PCC6803, and the *Arabidopsis thaliana* chloroplast.** Only transport systems for which experimental evidence suggests localization to the plasma membrane or the organellar envelope are shown. CE, chromatophore-encoded; NE, nucleus-encoded; Ion, ions/metals; IA, inorganic anions (phosphate, sulfate, nitrate, bicarbonate); AA, amino acid; S(P), sugars (hexoses, oligosaccharides) or sugar-phosphates; CB, mono-/di-/tricarboxylates; NT, nucleotides; LCW, lipid and lipopolysaccharide; MD, multidrug; O, other; U, unknown; OMP, outer membrane pores; Σ, total predicted transporters.

**Figure 2: Increased recovery of TM proteins by MS analysis of enriched insoluble chromatophore proteins. (A)** TEM micrographs of isolated chromatophore (left) and chromatophore in the context of a *P. chromatophora* cell (right). The outer envelope membrane (OM) observed in intact cells was lost during the isolation process. IM, inner envelope membrane; PG, peptidoglycan; TM, thylakoid membranes; Cb, carboxysomes. **(B)** 1 µg of protein from three replicates of each, chromatophore lysates (CL) as well as high salt and carbonate-washed *P. chromatophora* (PM) and chromatophore membranes (CM) was resolved on a 4-20 % polyacrylamide gel and silver stained. **(C)** Numbers of proteins encoded on the chromatophore genome (Genome) and chromatophore-encoded proteins identified with ≥3 SpC in chromatophore-derived samples in our previous (2017) and current (This study) proteome analysis. The number of predicted TMHs is indicated by a color code. **(D)** Detection of chromatophore-encoded transport systems. Annotation or TCDB-family, predicted mode of transport, substrates, and probable subcellular localization are provided. For each protein, the mean normalized intensity in CM (over both MS experiments) is indicated by a color code (see also **Table S1**).
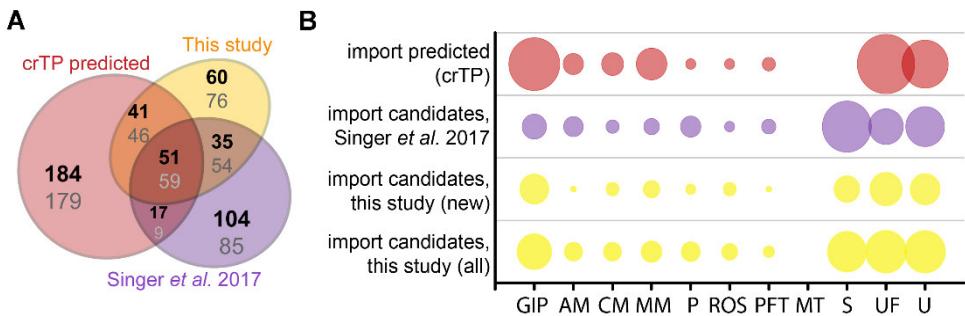
**Figure 3: No evidence for import of host-encoded multi-spanning TM proteins into chromatophores. (A)** Chromatophore-encoded (CE) and nucleus-encoded (NE) proteins enriched in CM compared to PM samples. Yellow, proteins enriched with high confidence; light yellow, proteins enriched with low confidence (LC); brown, proteins depleted in CM; black, proteins classified as "unclear" (see Methods and Fig. S2). Only proteins identified with ≥3 SpC in the chromatophore samples in at least one out of two independent MS experiments were considered. **(B, C)** The difference of intensities of individual proteins between CM and PM samples $(\overline{log2(normInt_{CM})} - \overline{log2(normInt_{PM})}$; Difference) is plotted against significance (-log10 p-values in Student´s t-test) for proteins detected with ≥3 SpC in the chromatophore samples (for proteins detected in CM only or CM and PM) or in whole cell samples (proteins detected in PM only). Values for proteins detected only in one sample have been imputed and are only shown when their difference is significant. The number of predicted TMHs (outside of the crTP) is indicated by a color code. Data from MS experiment 1 **(B)** and 2 **(C)** are shown separately. Scaffold1608-m.20717 and scaffold18898-m.107131 (see text) are marked by arrowheads in B and C, respectively. In both analyses, among the proteins enriched in CM (Difference CM-PM > 0), the proportion of identified multi-spanning TM proteins encoded in the chromatophore (49 of 409 in B; 39 of 134 in C) as compared to the nucleus (0 of 132 excluding the false positive in B; 1 of 54 in C) is significantly higher (both: p-value = 0.002, Fishers's Exact Test).
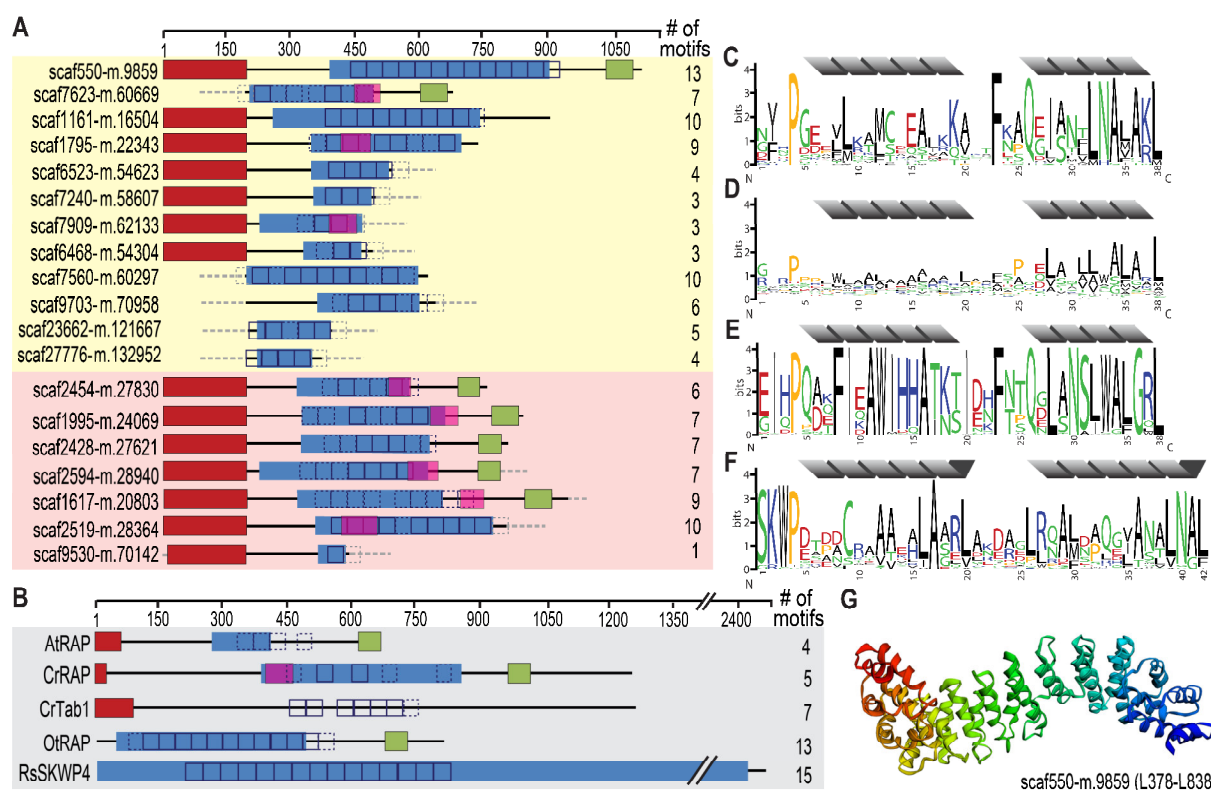
**Figure 4: Short orphan import candidates form distinct groups. (A)** For each group, representative MS-identified proteins (bright colors) and, if applicable, similar proteins identified among translated nuclear transcripts (pale colors) are displayed. Group 1: boxes indicate position of the predicted TMH. Group 2-4: colored background indicates ≥70% amino acid identity over alignments containing all MS-identified proteins of the respective group. The conserved sequence motif in group 2 was identified in diverse bacterial proteins (three examples with their NCBI accession number and aa positions are provided). Group 4: CxxC motifs are highlighted. Amphipathic: some short import candidates that do not belong to group 1 to 4 feature amphipathic helices. Areas highlighted in grey contain predicted alpha-helices. Corresponding wheel diagrams and hydrophobic moments are provided below. **(B)** Domain structure of bacterial proteins shown in A. Light blue boxes, conserved group 2 sequence motif; orange, group I intron endonuclease domain; light orange, GIY-YIG excision

858 nuclease domain; pink, Superfamily II DNA or RNA helicase domain (SSL2); light pink, DEXH-

859 box helicase domain of DEAD-like helicase restriction enzyme family proteins; blue, DNA-

860 binding motif found in homing endonucleases and related proteins (NUMOD); red arrows,

861 individual CxxC motifs. Sb, *Spirochaetia bacterium*; Cd, *Clostridioides difficile*; Gu,

862 *Gordonibacter urolithinfaciens*.

863



864

865 **Figure 5: Newly identified import candidates. (A)** Numbers of newly identified import

866 candidates in this study [see Fig. 3A, yellow], previously MS-identified import candidates

867 (Singer *et al.* 2017, purple), and in silico predicted import candidates [(10), red]. Numbers in

868 bold indicate distribution of proteins considering only HC import candidates, numbers in grey

869 considering all import candidates. **(B)** Functional categories of import candidates in (A). GIP,

870 genetic information processing; AM, amino acid metabolism; CM, carbohydrate metabolism;

871 MM, miscellaneous metabolism; P, photosynthesis and light protection; ROS, response to

872 oxidative stress; PFT, protein folding and transport; MT, metabolite transport; S, short

873 proteins (<90 aa) without functional annotation/homologs; UF, unspecific function; U,

874 unknown function. "New" import candidates were MS-identified in this study, but not in

875 Singer *et al.* 2017.

876

**Figure 6: Identification of an expanded family of putative OPR expression regulators targeted to the chromatophore (and mitochondrion) in *P. chromatophora*. (A)** Domain structure of 12 OPR-containing import candidates identified by MS (yellow background) and further 7 predicted import candidates with a similar domain structure (red background). The number of motif repeats identified in individual proteins is indicated. **(B)** Domain structure and motif repeats in (putative) expression regulators from other organisms. AtRAP, *A. thaliana* RAP domain-containing protein, NP_850176.1, (65); CrTab1, *Chlamydomonas reinhardtii* PsaB expression regulator, ADY68544.1, (92). OtRAP, *Orientia tsutsugamushi* uncharacterized RAP domain-containing protein, KJV97331.1, and RsSKWP4, *Ralstonia soleraceum* RipS4-family effector, AXW63421.1, (93) appear as the highest scoring BlastP/DELTA Blast hits (in the NCBI nr database) for *P. chromatophora* OPR proteins. **(C)** 38-aa repetitive motif found in *P. chromatophora* import candidates. **(D)** OPR motif found in *C. reinhardtii* expression regulators (designed according to (94)). **(E)** Motif derived from *O. tsutsugamushi* OPR proteins. **(F)** 42 aa SKWP motif derived from RipS-family effectors in *R. soleraceum, Xanthomonas euvesicatoria,* and *Mesorhizobium loti* (93, 95, 96). Individual repeats are predicted to fold into two α-helices (grey). Red, targeting signal (crTP for *P. chromatophora* proteins, cTP for AtRAP and CrTab1, mTP for CrRAP); blue, PRK09169-multidomain (Pssm-ID 236394); pink, FAST-kinase like domain (Pssm-ID 310980); green, RAP domain (Pssm-ID 312021); boxes, individual repeats of

31

896    the motifs shown in C-F ($p<e^{-20}$; $p<e^{-10}$ for CrRAP and CrTab1); dashed boxes, weak motif

897    repeats ($p<e^{-10}$; $p<e^{-7}$ for CrRAP and CrTab1); grey dashed boxes/lines, sequence information

898    incomplete. **(G)** Predicted 3D-structure of the OPR-containing region in scaffold550-m.9859.

899