

¹ **The UCYN-A transit peptide: a novel C-terminal
2 targeting system with distinctive biophysical substrate
3 signatures**

⁴ Author One^{1,*}, Author Two²

¹Institution One, City, Country

²Institution Two, City, Country

*Corresponding author: email@example.com

⁵ **Abstract**

⁶ UCYN-A is a nitrogen-fixing cyanobacterial endosymbiont that has evolved into an organelle-
⁷ like structure (nitroplast) in the marine alga *Braarudosphaera bigelowii*. Approximately 900
⁸ host-encoded proteins are imported into UCYN-A via a novel C-terminal transit peptide
⁹ (uTP), unlike the N-terminal signals used by mitochondria and chloroplasts.

10 **1 Introduction**

11 Mitochondria and chloroplasts import the vast majority of their proteins from the host cell
12 cytoplasm. This import is mediated by N-terminal transit peptides that are recognized by
13 dedicated translocon complexes (TOM/TIM in mitochondria, TOC/TIC in chloroplasts) and
14 typically cleaved upon import (?). These targeting systems are ancient, highly conserved, and
15 essential for organelle function. The evolution of efficient protein import was a key innovation that
16 enabled the extensive genome reduction characteristic of endosymbiont-to-organelle transitions.
17 UCYN-A (*Candidatus Atelocyanobacterium thalassa*) is a nitrogen-fixing cyanobacterial en-
18 dosymbiont that has recently been characterized as an organelle in the marine haptophyte
19 *Braarudosphaera bigelowii* (Coale et al., 2024). The symbiosis originated approximately 90–100
20 million years ago, and UCYN-A has undergone extreme genome reduction, retaining only ~22%
21 of genes compared to its free-living relative *Crocospaera watsonii* (Frail and Others, 2025). This
22 makes UCYN-A—termed the “nitroplast”—the first nitrogen-fixing organelle to be characterized
23 in eukaryotes.
24 Coale et al. identified approximately 368 host-encoded proteins that are imported into UCYN-A,
25 representing a substantial complement of the nitroplast proteome (Coale et al., 2024). These
26 imported proteins carry a novel C-terminal targeting sequence of ~100–150 amino acids, termed
27 the UCYN-A transit peptide (uTP). Using a hidden Markov model (HMM) derived from ex-
28 perimentally validated proteins, approximately 900 proteins in the *B. bigelowii* proteome are
29 predicted to carry uTP. The imported proteins fill critical metabolic gaps in UCYN-A, including
30 biosynthesis of threonine, serine, proline, pyrimidines, and tetrahydrofolate.
31 Despite the identification of uTP as a targeting signal, several fundamental questions remain
32 unanswered. What are the conserved sequence and structural features of uTP? Why is uTP
33 C-terminal, in contrast to the N-terminal signals used by mitochondria and chloroplasts? What
34 is the evolutionary origin of uTP—was it co-opted from an existing cellular function or did it
35 evolve de novo? And what determines which host proteins acquire uTP for nitroplast targeting?
36 Here, we present a systematic computational characterization of the uTP system. We identify
37 conserved motif architecture with invariant anchor sequences and a conserved U-bend structure.
38 We demonstrate that uTP has no detectable homologs in related haptophyte genomes, supporting
39 a de novo evolutionary origin. Most importantly, we show that uTP-containing proteins share
40 a distinctive biophysical signature—more disordered, acidic, and stable than the general pro-
41 teome—that enables prediction of uTP status with 90% accuracy from mature domain sequence
42 alone. These findings establish uTP as a novel solution to organellar protein targeting and provide
43 a foundation for understanding the mechanisms underlying nitroplast protein import.

44 **2 Results**

45 **2.1 uTP comprises conserved motif architecture with invariant anchor sequences**

46 To systematically characterize the uTP system, we applied the HMM profile from Coale et al. to
47 the full *B. bigelowii* proteome, identifying 933 proteins with putative uTP sequences. The uTP
48 region spans 58–1091 amino acids (median ~206 aa), located at the C-terminus of each protein.
49 De novo motif discovery using MEME identified 8 conserved motifs within the uTP region
50 (Table ??). Two motifs showed near-universal prevalence: motif 2 (93% of sequences) and motif 1
51 (91%). These “anchor motifs” appear in a stereotyped order, with motif 2 consistently preceding
52 motif 1 at the N-terminal end of the uTP region. The remaining motifs (3–9) appear in variable
53 combinations downstream, creating combinatorial diversity among uTP sequences.
54 Extended scanning with MAST across all 933 HMM-predicted proteins confirmed this architecture.
55 Of 745 proteins with detectable motif hits, 81.5% showed one of four canonical terminal motif

56 patterns (ending with motif 4, 5, 7, or 9), with terminal motif 7 being most common (61.7%).
57 The remaining proteins either lacked detectable motifs (20.2%) or showed non-canonical patterns.

58 **2.2 uTP adopts a conserved U-bend structure**

59 To assess structural conservation, we predicted three-dimensional structures for representative
60 uTP sequences using AlphaFold2. Despite sequence variation in the variable motif region, all
61 predicted structures converged on a conserved fold: two α -helices connected by a turn, forming a
62 characteristic “U-bend” configuration.

63 Structural alignment across predictions revealed low variance, with pairwise RMSD values
64 consistently below 4 Å. This structural conservation suggests that the U-bend fold is functionally
65 constrained, potentially reflecting requirements for recognition by the import machinery.

66 **2.3 No uTP homologs detected in related haptophyte genomes**

67 To investigate the evolutionary origin of uTP, we searched for homologous sequences in related
68 haptophyte species. We compiled proteomes from seven haptophyte species spanning estimated
69 divergence times of 330–920 million years from *B. bigelowii* (Table ??).

70 Profile HMM search with the uTP model yielded 84 total hits across all haptophyte proteomes.
71 However, detailed analysis of hit coordinates revealed that all 33 significant hits (E-value < 0.01)
72 aligned exclusively to the mature domain region of the HMM (positions 100–713), not the uTP-
73 specific region (positions 1–100). These hits represent orthologs of the mature protein domains,
74 not uTP homologs.

75 We also searched for individual uTP motifs using position weight matrices (MAST). While
76 530 individual motif hits were detected across haptophyte proteomes, the defining feature of
77 uTP—co-occurrence of the two anchor motifs (MEME-1 and MEME-2) at the C-terminus—was
78 completely absent. Only one protein in any haptophyte genome contained both anchor motifs,
79 and these were located internally (not C-terminal), in a protein 1269 amino acids long.

80 As a control, we performed the same motif search against the *Arabidopsis thaliana* proteome.
81 Similar hit rates confirmed that individual motif matches represent generic sequence patterns,
82 not uTP-specific features.

83 The complete absence of proteins with the uTP motif architecture in any haptophyte genome
84 strongly supports de novo evolution of uTP in the *B. bigelowii* lineage, likely after establishment
85 of the UCYN-A symbiosis.

86 **2.4 uTP-containing proteins have distinctive biophysical properties**

87 To identify features distinguishing uTP-containing proteins from the general proteome, we
88 constructed a carefully controlled comparison. From each uTP protein, we extracted the mature
89 domain by removing the uTP region identified by HMM. We then selected length-matched control
90 proteins from the *B. bigelowii* proteome, filtering by subcellular localization prediction (CELLO)
91 to exclude proteins with other targeting signals (signal peptides, mitochondrial/chloroplast transit
92 peptides). This yielded 605 uTP mature domains and 773 cytoplasmic/nuclear controls.

93 We trained binary classifiers using ProtT5 protein language model embeddings as features. Logistic
94 regression achieved 92.8% accuracy ($F_1 = 0.92$), with significance confirmed by permutation
95 testing ($p = 0.002$, 1000 permutations). The ROC AUC was 0.92, indicating strong discriminative
96 power.

97 To understand what drives this classification, we computed biophysical properties for all proteins.
98 Three properties showed large effect sizes distinguishing uTP from control proteins (Table ??):

- 99 • **Intrinsic disorder:** uTP proteins have significantly higher predicted coil/disordered
100 content (Cohen's $d = +1.05$, $p < 0.006$)
101 • **Isoelectric point:** uTP proteins are more acidic (mean pI = 6.5 vs 8.5; $d = -0.89$, $p <$
102 0.006)
103 • **Stability:** uTP proteins have lower instability index, indicating greater stability ($d = -0.81$,
104 $p < 0.006$)

105 All p-values were Bonferroni-corrected for 8 tests.

106 To validate these findings at scale, we applied the trained classifier to the entire *B. bigelowii*
107 proteome (933 uTP vs 43,430 non-uTP proteins). Despite the extreme class imbalance, the
108 classifier maintained strong performance (ROC AUC = 0.920, recall = 80.3

109 These biophysical signatures have plausible mechanistic interpretations. Increased disorder may
110 facilitate protein unfolding during translocation across the nitroplast membrane. The acidic
111 character could mediate electrostatic interactions with positively charged components of the
112 import machinery. The greater stability may reflect selection for function in the nitroplast
113 environment.

114 3 Discussion

115 4 Methods

116 4.1 Data sources

117 Proteomics data, genome annotations, and the uTP HMM profile were obtained from Coale et al.
118 (Coale et al., 2024). The *B. bigelowii* transcriptome contained approximately 44,000 predicted
119 proteins.

120 Haptophyte proteomes were downloaded from NCBI GenBank: *Emiliania huxleyi* CCMP1516
121 (GCA_000372725.1, ~38,500 proteins), *Chrysochromulina tobini* CCMP291 (GCA_001275005.1,
122 ~16,700 proteins), *Diacronema lutheri* (GCA_019448385.1, ~14,400 proteins), *Pavlovales* sp.
123 CCMP2436 (GCA_026770615.1, ~26,100 proteins), *Prymnesium parvum* 12B1 (GCA_041296205.1,
124 ~23,700 proteins), and *Prymnesium* sp. SGEUK-05 (GCA_046255225.1, ~15,000 proteins).

125 Acknowledgements

126 [To be added]

127 Author Contributions

128 [To be added]

129 Data Availability

130 [To be added]

131 References

- 132 T. H. Coale, V. Loconte, K. A. Turk-Kubo, J. P. Zehr, and D. Bhattacharya. Nitrogen-fixing
133 organelle in a marine alga. *Science*, 384(6692):217–222, 2024. doi: 10.1126/science.adk1075.
134 S. Frail and Others. Genomes of nitrogen-fixing eukaryotes reveal an alternate path for organello-
135 genesis. *Nature*, 2025. In press.