

Supplementary Information of “Karasuyama et al. Understanding Colour Tuning Rules and Predicting Absorption Wavelengths of Microbial Rhodopsins by Data-Driven Machine-Learning Approach” submitted to Scientific Reports

Masayuki Karasuyama^{1,2,3}, Keiichi Inoue^{4,2,5,6}, Ryoko Nakamura⁷,
Hideki Kandori^{4,5}, and Ichiro Takeuchi^{1,3,7}

- ¹Department of Computer Science, Nagoya Institute of Technology, Gokiso, Showa-ku, Nagoya, Aichi, 466-8555, Japan
- ²PRESTO, Japan Science and Technological Agency (JST), 4-1-8 Honcho, Kawaguchi, Saitama, 332-0012, Japan
- ³Center for Materials research by Information Integration, National Institute for Materials Science (NIMS), Tsukuba 305-0047, Japan
- ⁴Department of Life Science and Applied Chemistry, Nagoya Institute of Technology, Gokiso, Showa-ku, Nagoya, Aichi, 466-8555, Japan
- ⁵OptoBioTechnology Research Center, Gokiso, Showa-ku, Nagoya, Aichi, 466-8555, Japan
- ⁶The Institute for Solid State Physics, The University of Tokyo, 5-1-5 Kashiwanoha, Kashiwa, Chiba 277-8581, Japan.
- ⁷RIKEN Center for Advanced Intelligence Project, Nihonbashi 1-chome, 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan

1 **Supplementary Information 1:** Abbreviations of the names of microbial
2 rhodopsins

3 *AcetR1*: *acetabularia* rhodopsin of *Acetabularia acetabulum*

4 *AR1*: archaerhodopsin 1 of *Halorubrum sodomense*

5 *AR2*: archaerhodopsin 2 of *Halorubrum sodomense*

6 *AR3*: archaerhodopsin 3 of *Halorubrum sodomense*

7 *ASR*: *Anabaena* sensory rhodopsin of *Anabaena (Nostoc) sp.* PCC7120

8 *BPR*: blue proteorhodopsin of uncultured γ -proteobacterium

9 *BR*: bacteriorhodopsin of *Halobacterium salinarum*

10 *Chrimson*: channelrhodopsin-1 of *Chlamydomonas noctigama*

11 *Chronos*: channelrhodopsin of *Stigeoclonium helveticum*

12 *CnChR2*: channelrhodopsin-2 of *Chlamydomonas noctigama*

13 *CrChR1*: channelrhodopsin-1 of *Chlamydomonas reinhardtii*

14 *CrChR2*: channelrhodopsin-2 of *Chlamydomonas reinhardtii*

15 *DChR1*: channelrhodopsin-1 of *Dunaliella salina*

16 *PaR*: DTG rhodopsin of *Pantoea ananatis*

17 *PspR*: DTG rhodopsin of *Pseudomonas putida*

18 *PvR*: DTG rhodopsin of *Pantoea vagans*

19 *ESR*: proton pump rhodopsin of *Exiguobacterium sibiricum*

20 *FR*: chloride pump rhodopsin of *Fulvimarina pelagi*

21 *GINaR*: sodium pump rhodopsin of *Gillisia limnaea*

22 *GPR*: green proteorhodopsin of uncultured γ -proteobacterium

23 *GR*: *Gloeobacter* rhodopsin of *Gloeobacter violaceus* PCC 7421

24 *GtACR1*: anion channel rhodopsin of *Guillardia theta*

25 *GtACR2*: anion channel rhodopsin of *Guillardia theta*

26 *HmBRI*: bacteriorhodopsin of *Haloarcula marismortui*

27 *HmBRII*: bacteriorhodopsin of *Haloarcula marismortui*

28 *HmHR*: halorhodopsin of *Haloarcula marismortui*

29 *HmSRI*: sensory rhodopsin I of *Haloarcula marismortui*

30 *HmSRII*: sensory rhodopsin II of *Haloarcula marismortui*

31 *HmSRIII*: sensory rhodopsin III of *Haloarcula marismortui*

32 *HsHR*: halorhodopsin of *Halobacterium salinarum*
 33 *HsSRI*: sensory rhodopsin I of *Halobacterium salinarum*
 34 *HvSRI*: sensory rhodopsin I of *Haloarcula vallismortis*
 35 *HwBR*: bacteriorhodopsin of *Haloquadratum walsbyi*
 36 *IaNaR*: sodium pump rhodopsin of *Indibacter alkaliphilus*
 37 *IaR1*: proteorhodopsin of *Indibacter alkaliphilus*
 38 *KR1*: proteorhodopsin of *Krokinobacter eikastus*
 39 *KR2*: sodium pump rhodopsin of *Krokinobacter eikastus*
 40 *LaNaR*: sodium pump rhodopsin of *Lyngbya aestuarii*
 41 *LR*: *Leptosphaeria* rhodopsin of *Leptosphaeria maculans*
 42 *MR*: middle rhodopsin of *Haloquadratum walsbyi*
 43 *NdPR*: proteorhodopsin of *Nonlabens dokdonensis*
 44 *NdNaR*: sodium pump rhodopsin of *Nonlabens dokdonensis*
 45 *NmClR*: chloride pump rhodopsin of *Nonlabens marinus*
 46 *NmNaR*: sodium pump rhodopsin of *Nonlabens marinus*
 47 *NmPR*: proteorhodopsin of *Nonlabens marinus*
 48 *NpHR*: halorhodopsin of *Natronobacterium pharaonis*
 49 *NpSRII*: sensory rhodopsin II of *Natronobacterium pharaonis*
 50 *NR*: *Neurospora* rhodopsin of *Neurospora crassa*
 51 *PhaeoRD1*: putative proton pump rhodopsin of *Phaeosphaeria (Stagonospora)*
 52 *nodorum*
 53 *PhaeoRD2*: putative proton pump rhodopsin of *Phaeosphaeria (Stagonospora)*
 54 *nodorum*
 55 *PoXeR*: xenorhodopsin of *Parvularcula oceani*
 56 *PsChR2*: channelrhodopsin-2 of *Proteomonas sulcata*
 57 *SrSRI*: sensory rhodopsin I of *Salinibacter ruber*
 58 *TR*: thermophilic rhodopsin of *Thermus thermophilus*
 59 *TsChR*: channelrhodopsin of *Tetraselmis striata*
 60 *XR*: xanthorhodopsin of *Salinibacter ruber*

61 **Supplementary Information 2:** References for Supplementary Table 1

62 See MS Word file `References_Supplementary_Table_1.docx`

63 **Supplementary Information 3:** Software implementation

64 Our R code for the group-LASSO-based wavelength prediction is available
65 at <http://...> (The site will be public after acceptance). “R” is a standard
66 command-based environment for statistical computations that is freely available
67 and installable on Windows, Mac OS, and most major distributions of Linux.
68 In the R command-line, our package can be installed by the following command
69 at the directory in which the downloaded file exists:

```
70 > install.packages(c("grplasso", "seqinr", "foreach"))  
71 > install.packages("GrplassoSeq_1.0.tar.gz",  
72                   repos = NULL, type = "source")
```

73 A set of sequences and their wavelengths should be provided as two separated
74 input files. Sequence data should be in the FASTA format, which can be read by
75 the `read.fasta` function of the R `seqinr` package (<https://cran.r-project.org/web/packages/seqinr/index.html>). Note that, for both the training and
76 target proteins, the sequences should be identically arranged according to length
77 beforehand. Wavelengths should be stored in a text file in which each line con-
78 tains a wavelength for a sequence. The order of sequences in the FASTA file and
79 the wavelength file should be consistent. For wavelength-unknown sequences,
80 please specify “NA” in the wavelength file. Proteins specified as “NA” are
81 regarded as target proteins, and other proteins with wavelengths are used as
82 training proteins.
83

⁸⁴ **Supplementary Table 1:** Microbial rhodopsin database composed of the
⁸⁵ amino-acid sequences and absorption wavelengths λ_{\max} s of 519 proteins pre-
⁸⁶ viously reported in the literature and 277 proteins newly investigated by our
⁸⁷ group.

⁸⁸ See MS Excel file: `Supplementary_Table_1.xlsx`

89 **Supplementary Table 2:** List of all fitted coefficient parameters

90 See MS Excel file: [Supplementary_Table_2.xlsx](#)

91 **Supplementary Figure 1**

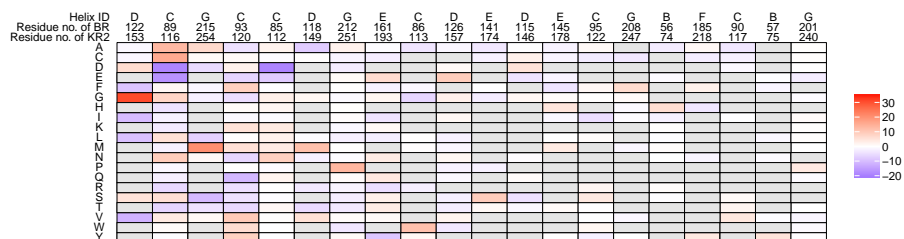


Figure S1: **Coefficient parameters of the statistical model fitted with all 796 rhodopsin proteins as the training set.** The coefficient parameters at 20 active residues in decreasing order of s_j are given. This figure has the same format as Figure 5, except that all 796 rhodopsin proteins, including those in the KR2 group, were used as the training data. Although the resulting coefficient parameters depend on the given training data, most of the coefficients parameters in Supplementary Figure S3 are similar to those in Figure 5 in the main text. The correlation coefficient between these two results was 0.9650.

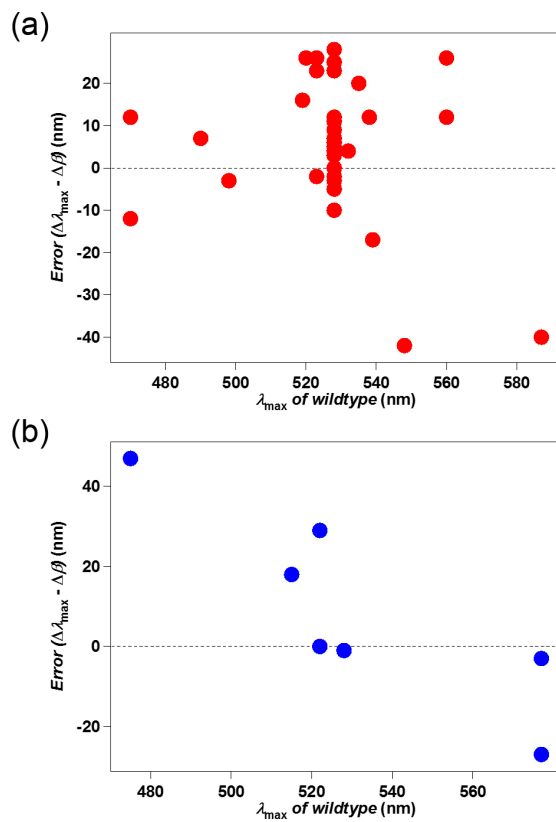


Figure S2: The correlation of the difference between the experimentally observed absorption shift (λ_{\max}) and predicted shift by ML ($\Delta\beta$) against the λ_{\max} of wildtype protein for the removal (a) or introduction (b) of counterion (D or E) at homologous positions to BR 85 or 89.

93 **Supplementary Figure 3**

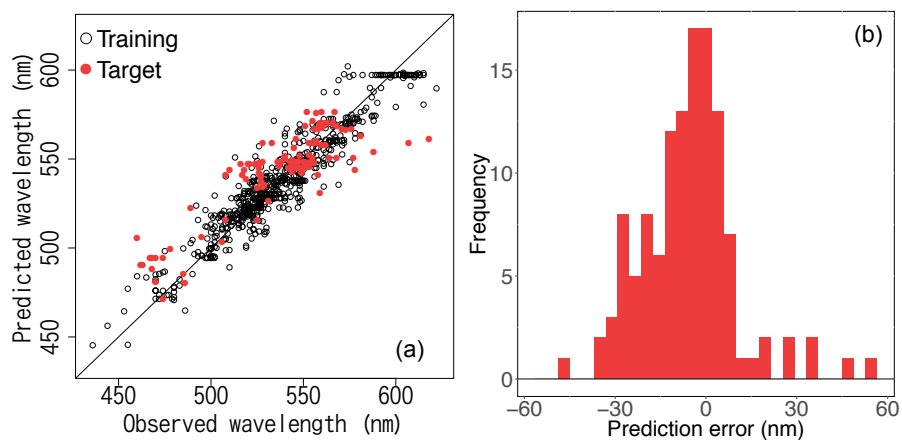


Figure S3: **Absorption wavelength prediction results for GR wildtype and its 121 variants.** The plots (a) and (b) have the same format as Figure 6 (e) and (f). GR wildtype and its 121 variants are the target proteins and the other 674 proteins are contained in the training data. The mean absolute error for the target proteins was 12.4 nm, and the correlation coefficient ρ was $\rho = 0.8596$ (p -value for the null hypothesis $\rho = 0$ was 8.7154×10^{-37}).