

## Supplementary Data

A list of plastid assigned, manually curated gene models was prepared from the first release of the *Phaeodactylum tricornutum* genome v1.0<sup>1</sup> (Bowler et al., in preparation). Predicted signal peptides and signal peptide cleavage sites were displayed according to the results of SignalP<sup>2</sup> (Bendtsen et al., 2004), taking into account the different prediction methods Neuronal networks (NN) (Nielsen et al., 1997) or Hidden Markov models (HMM) (Nielsen and Krogh, 1998) (supplementary Fig. 6). The results were used to prepare sequence logos (Schneider and Stephens, 1990) using the WebLogo server<sup>3</sup> (Crooks et al., 2004) to illustrate the predictions of the different prediction algorithms with predictions combining computation and manual correction (Fig. 1).

## References

- Bendtsen JD, Nielsen H, von Heijne G, Brunak S (2004) Improved prediction of signal peptides: SignalP 3.0. *J Mol Biol* 340: 783–95.
- Crooks GE, Hon G, Chandonia JM, Brenner SE (2004) WebLogo: a sequence logo generator. *Genome Res* 14: 1188–90.
- Nielsen H, Engelbrecht J, Brunak S, von Heijne G (1997) A neural network method for identification of prokaryotic and eukaryotic signal peptides and prediction of their cleavage sites. *Int J Neural Syst* 8: 581–99.
- Nielsen H, Krogh A (1998) Prediction of signal peptides and signal anchors by a hidden Markov model. *Proc Int Conf Intell Syst Mol Biol* 6: 122–30.
- Schneider TD, Stephens RM (1990) Sequence logos: a new way to display consensus sequences. *Nucleic Acids Res* 18: 6097–100.

---

<sup>1</sup><http://genome.jgi-psf.org/Phatr1/Phatr1.home.html>

<sup>2</sup><http://www.cbs.dtu.dk/services/SignalP/>

<sup>3</sup><http://weblogo.berkeley.edu/>

Name	Protein Id	Sequence	NN Ymax	HMM Cmax
(a)				
Albino3	28289	<u>MARTRHGARLARCA</u> <u>FFVWLWV</u> <u>ASSTTA</u> <u>FTT</u> TS <sup>01</sup> SPRLAAHFRSASRTQRTTTT <sup>50</sup>	<b>0.758*</b>	0.720*
AtpC	21239	<u>MRSFCIAALLAVAS</u> <u>FTTQ</u> PT <sup>05</sup> SFTVKTANVGERASGVFPEQSSAHRTRKA	<b>0.701*</b>	0.595*
FbaC1	22933	<u>MKLSTAALFFIPAVVA</u> <u>FAPPQ</u> A <sup>10</sup> AFRSNPALFATE <sup>15</sup> TAAEKTTFSKMPASVK	<b>0.692*</b>	0.549*
FCP	48997	<u>MAKFSLAILAALVAT</u> <u>ASAFVA</u> SP <sup>20</sup> TTSSASTALRATQEGVWDPLGLMTLGT	<b>0.640*</b>	0.637*
FCP	50582	<u>MKSSAVLALAMAGST</u> <u>AA</u> FAP <sup>25</sup> TSSTQASTSSTSLQAAMPDRLWNTMVDKTE	<b>0.436*</b>	0.259*
FCP	51914	<u>MKITALCLTALV</u> <u>SASHA</u> FAP <sup>30</sup> STPSSASSSARALSTESTSDVPLLQIKEKV	<b>0.517*</b>	0.429*
FCP	49699	<u>MKYRTSILF</u> <u>SALATS</u> <u>ATA</u> FAP <sup>35</sup> TQSI <sup>40</sup> TRTLTQTNMADFNKDDFLSFKEKDK	<b>0.533*</b>	0.528*
FCP	62104	<u>MKLF</u> TIFLPLVLVGT <sup>45</sup> AAG <sup>50</sup> FASGPFSSKKASPSPEVSI <sup>55</sup> ESMPGIVAPTGFDD	<b>0.605*</b>	0.555*
FCP	47005	<u>MMKLALIASL</u> <u>VAGAAA</u> FAP <sup>50</sup> ASKQASSSALKAFENEAGVIQPTGFFD <sup>55</sup> PFGL	<b>0.664*</b>	0.632*
Hcf136	48926	<u>MKFS</u> TLVAVFVSG <sup>55</sup> SAA <sup>60</sup> FV <sup>65</sup> PN <sup>70</sup> SFATARTDALRMTSPGYEIESESSGARR	<b>0.717*</b>	0.690*
OEE3	48703	<u>MKLALVFS</u> <u>LFATAAA</u> F <sup>65</sup> SQ <sup>70</sup> EASRREALTKGAAAFGAAFLPVAANA <sup>75</sup> AVGES <sup>80</sup> P	<b>0.575*</b>	0.444*
petC	47867	<u>MKIIP</u> TVTSLALLAIS <sup>75</sup> IRA <sup>80</sup> FT <sup>85</sup> PLAPRHTHASSKTWAASLEEPAYEGTID	<b>0.641*</b>	0.529*
PsbM	50581	<u>MKS</u> FOLLTLFALIAA <sup>85</sup> SLA <sup>90</sup> FAP <sup>95</sup> NQAPQQVAKAAFKAAGALPVAALAAPAF	<b>0.618*</b>	0.586*
putative Tic62	35796	<u>MTV</u> STYFIFFFTLGRCT <sup>95</sup> AA <sup>100</sup> F <sup>105</sup> PS <sup>110</sup> LGASLTRSTAVPRFALSPSSNDDGIHS	<b>0.728*</b>	0.533*
recA	61939	<u>MMHRKIALVAGLW</u> <u>FCLL</u> <u>GSSCHG</u> FGL <sup>115</sup> LGGAAHSRTVRRWPQQWAVSSPIGP	<b>0.610*</b>	0.527*
Sep	62042	<u>MHHFAK</u> VILLV <sup>120</sup> CVAMLA <sup>125</sup> AVY <sup>130</sup> TE <sup>135</sup> FAT <sup>140</sup> PSRSVQSSAVSITNSMPFRTSALN	<b>0.688*</b>	0.682*
Trx-f	36974	<u>MMQQQQQQHAA</u> <u>ORRAGV</u> <u>GTILL</u> <u>SLWFV</u> <u>PAATA</u> FAP <sup>145</sup> STPLAFRTQT <sup>150</sup> PVVT	<b>0.739*</b>	0.738*
Trx-m	62165	<u>MYRAQQYCRS</u> <u>RTLFI</u> <u>HYAIVL</u> <u>LVTRYCSA</u> <u>FCS</u> LE <sup>155</sup> PVLRPSRWISNRKSL	<b>0.706*</b>	0.472*
GLNA2	14218	<u>MKLNIAAIALFA</u> <u>ASASA</u> FAP <sup>160</sup> RFASPRSHATVLSAVLEERTGQSQ <sup>165</sup> LDPAVI	<b>0.515*</b>	<b>0.515*</b>
ADK	45920	<u>MLRSLAFTAS</u> <u>VALLFS</u> <u>LDPLLA</u> FAP <sup>170</sup> IRTTTSVAVSPSGVSI <sup>175</sup> RAQTGDQLFA	0.520*	<b>0.767*</b>
CPE	46734	<u>MKLPWL</u> GPSAAALLSSQ <sup>180</sup> TMA <sup>185</sup> FL <sup>190</sup> PS <sup>195</sup> SLPSQSARNAGVTLQEKPSADSS <sup>200</sup> FF	0.624*	<b>0.793*</b>
Elip	18985	<u>MAPLR</u> TTFALLLSLV <sup>200</sup> SASA <sup>205</sup> FAP <sup>210</sup> VQNVARKQTSVSAFKIDPQLYDDAVSDW	0.599*	<b>0.725*</b>
Enolase	44169	<u>MLFKP</u> STLLALFAVAG <sup>215</sup> TTLA <sup>220</sup> FAP <sup>225</sup> RSTT <sup>230</sup> PLTSTTRGSASSSVTTLAMS <sup>235</sup> GI	0.697*	<b>0.704*</b>
FbaC2	14856	<u>MKI</u> AVVAVFVIAQCG <sup>240</sup> AFAP <sup>245</sup> PAHY <sup>250</sup> SRTVTSSTLLGAKEKGGTSKELDLPCAD	0.488*	<b>0.631*</b>
FBPC1	45363	<u>MEK</u> WGF <sup>255</sup> SRGQVPLLLSVIALSFVLLP <sup>260</sup> TNS <sup>265</sup> FQT <sup>270</sup> STGQSQPQLPASSASL	0.762*	<b>0.884*</b>
FBPC2	41361	<u>MKI</u> ALLPVVFS <sup>275</sup> SAISVRA <sup>280</sup> FLP <sup>285</sup> TRPSPATQYFRYGPRLATASLSQAAGAAVS	0.491*	<b>0.897*</b>
FCP	22723	<u>MKT</u> SAIVAILAVSG <sup>285</sup> ASAF <sup>290</sup> TP <sup>295</sup> NTNAPQQLTKVGATAELDNMLGV <sup>300</sup> DIETGKK	0.724*	<b>0.961*</b>
FCP	46231	<u>MK</u> CIAAIALLATTASA <sup>305</sup> FN <sup>310</sup> AFGAAKKAAPKPKVFSI <sup>315</sup> ETIPGALAPVGFIDP	0.654*	<b>0.894*</b>
Fcp related	62176	<u>MK</u> KILKRSTICCLGLY <sup>320</sup> TGSPCSYA <sup>325</sup> FR <sup>330</sup> PP <sup>335</sup> ISEEIGCSPTLAI <sup>340</sup> FLKEDCH	0.328*	<b>0.923*</b>
Flavodoxin	23541	<u>MNTQ</u> FVSALLLASAAIT <sup>345</sup> NGFA <sup>350</sup> FV <sup>355</sup> NTHRYTASTTAL <sup>360</sup> EAGVKIYYSSTGNT	0.608*	<b>0.623*</b>
FtrB	43703	<u>MKV</u> FLSFVALLVLSL <sup>365</sup> TQA <sup>370</sup> FMP <sup>375</sup> VSKPSFGRVGGTVYMAKEMTPEEBE <sup>380</sup> IAVE	0.799*	<b>0.985*</b>
FtsZ	31286	<u>MR</u> ISSKLM <sup>385</sup> AVLTTTVG <sup>390</sup> IASA <sup>395</sup> FQ <sup>400</sup> PAT <sup>405</sup> PSRIRQSPRTT <sup>410</sup> VVLAMSDDKPPA	0.729*	<b>0.798*</b>
GLRXC2	43651	<u>MT</u> TNRVLF <sup>415</sup> LAILAL <sup>420</sup> TCL <sup>425</sup> TRAF <sup>430</sup> LPV <sup>435</sup> TTFVSR <sup>440</sup> TVAGSHGVHSLSMADMADD	0.715*	<b>0.927*</b>
GLRXC3	41758	<u>MAIL</u> SRTL <sup>445</sup> VGVL <sup>450</sup> TAT <sup>455</sup> FACIPHE <sup>460</sup> TNA <sup>465</sup> FLK <sup>470</sup> PHLGTVSSLSVIGEPS <sup>475</sup> STRLHI	0.448*	<b>0.767*</b>
GLRXC4	41850	<u>MKS</u> FIACAFILASAAA <sup>480</sup> FAP <sup>485</sup> QPGSPVLRST <sup>490</sup> SALSASPS <sup>495</sup> DDRS <sup>500</sup> PNPIK <sup>505</sup> VM	0.488*	<b>0.648*</b>
Hmox	45923	<u>MRNI</u> TKRL <sup>510</sup> LAVLAVAAS <sup>515</sup> PRTSQA <sup>520</sup> FAP <sup>525</sup> SVSRKQTNCHK <sup>530</sup> TPLTTSVSIKN	0.688*	<b>0.820*</b>
HSP70D	44604	<u>MRSQ</u> ILLCFRILFL <sup>535</sup> FLP <sup>540</sup> TTTIA <sup>545</sup> FST <sup>550</sup> HPKFRARRKQPCASSQ <sup>555</sup> LAVTDKHE	0.917*	<b>0.923*</b>
NiR	48100	<u>MVQ</u> TLPVFRSLSVMAIGS <sup>560</sup> FLLALG <sup>565</sup> SPSADAF <sup>570</sup> AF <sup>575</sup> STRLP <sup>580</sup> TQSLSLRSTST	0.363*	<b>0.442*</b>
NTT1	50541	<u>MI</u> PTTSVASHRSV <sup>585</sup> VLVGLACT <sup>590</sup> TLLLAVLSP <sup>595</sup> SSTEA <sup>600</sup> FAP <sup>605</sup> SAHRYSANQAVP	0.453*	<b>0.744*</b>
PAP fibrillin II	45766	<u>MPGR</u> VVPLSLLAVVCAS <sup>610</sup> DLFR <sup>615</sup> STLAF <sup>620</sup> HE <sup>625</sup> KVSTSSRTIPPSTTQLHAFGFL	0.632*	<b>0.894*</b>
PetJ	43773	<u>MKL</u> AVIATLLATASAF <sup>630</sup> FSIQAE <sup>635</sup> FSKVAKGAAAVGVGAVIAAAPALAGDVGA	0.789*	<b>0.839*</b>
PPdK	49184	<u>MK</u> FSSAAATVGLLLSGHAP <sup>640</sup> MIFS <sup>645</sup> FVT <sup>650</sup> PPSRFASGHQASGERSI <sup>655</sup> SHST	0.433*	<b>0.823*</b>
Prk	43777	<u>MK</u> FAV <sup>660</sup> FASLTATAAA <sup>665</sup> FAP <sup>670</sup> TAFVPSNLRGVAPSASSLN <sup>675</sup> MALKEGQTPIIIG	0.502*	<b>0.539*</b>
PtCA1	45670	<u>MK</u> FLSASIALLACATS <sup>680</sup> VEAF <sup>685</sup> FN <sup>690</sup> ANKAFRF <sup>695</sup> GAKAMPEV <sup>700</sup> SE <sup>705</sup> SATSALSAGGA	0.792*	<b>0.951*</b>
Rpe	24020	<u>MK</u> FTIVSLAAVVASA <sup>710</sup> SAFAP <sup>715</sup> ATKSVRSV <sup>720</sup> SALNVWGD <sup>725</sup> KDYLIAPSILSADF	0.484*	<b>0.602*</b>
RPI	47776	<u>MRL</u> TAGASILLASSA <sup>730</sup> HAF <sup>735</sup> TN <sup>740</sup> PAFFPRTATFASTSSYTSALILKMGVDQDE	0.550*	<b>0.890*</b>

(continuation on next page)

highest value  
bold

**BOLD:** prediction by NN, UNDERLINED: prediction by HMM, GREY: conserved motif at signal peptide cleavage site, BLACK: hydrophobic (ACFGILMPVWY), GREEN: hydrophilic (NQST), BLUE: basic (HKR), RED: acidic (DE)  
\* predicted signal peptide cleavage site coincides with conserved motif

supplementary Fig. 6, caption on page 4

Name	Protein Id	Sequence	NN Ymax	HMM Cmax
<b>(a, continuation)</b>				
secA	43680	<u>MRLTLTIQVALTLLLLPSSTVWA</u> <u>FRT</u> SPATQVFSRSRPMRSVASSSSSFS	0.803*	<b>0.978*</b>
SRP54	46102	<u>MRLQSGCVLTLAATFYPSTQA</u> <u>F</u> SIFSVGSPFASSSFASRESDVSRKSY	0.624*	<b>0.947*</b>
sufD	50465	<u>MKFTSVTLACFFV</u> <u>ESSG</u> <u>VSA</u> <u>F</u> FTSTPTQRQRSIPLPIVSSSRRTRSALH	0.690*	<b>0.802*</b>
TAL	51830	<u>MKSLVPLAFVLATSSG</u> <u>F</u> APREHHYRPNQPARTNALVLQAE <del>STAAV</del> LASA	0.385*	<b>0.795*</b>
Tic55	62146	<u>MALRRSISRLAMVYLVTLC</u> <u>LKTA</u> <u>FVSA</u> <u>F</u> SSSTNTPTPTSSVQQRSANQQA	0.674*	<b>0.702*</b>
Tkt	43792	<u>MKFS</u> <u>SIALATVFFAAK</u> <u>AG</u> <u>F</u> TCSSIKPRFGVQAHSVLKMSTATETDKMAA	0.595*	<b>0.782*</b>
Tpt3	18575	<u>MMKRALVVLTL</u> <u>SVGVSARASA</u> <u>F</u> APGAAVKNHAGATQSAIHKQTPFPTE	0.601*	<b>0.779*</b>
VDE-like2	62162	<u>MKLHRKGRYRLLVTA</u> <u>VLLG</u> <u>TVCS</u> <u>F</u> VPENLRSGSVRIPRKNANAGSVPGTH	0.748*	<b>0.858*</b>
Tic32	11808	<u>MKYSILGNILSLGLVLE</u> <u>TTKAW</u> <u>S</u> VPPPPQTSAPRESSGDSGTSTPVSP	0.741*	<b>0.538*</b>
FBPC4	27440	<u>MGRGVIIFCVKNFAVWLLI</u> <u>ITSAVSIQA</u> <u>W</u> IPLPLSATVKARIDSTTLFFS	0.523*	<b>0.556*</b>
FSA	48116	<u>MAISRSRRRSNGLGIVLVW</u> <u>TIFISAVW</u> <u>G</u> WTPPLRGSSRFLDSADPNVWN	0.825*	<b>0.886*</b>
GLXC1	44755	<u>MKLVSR</u> <u>SFC</u> <u>T</u> LAWTCGAAQAW <del>SV</del> APRTVARNRPAWVAARHPHTTACAFS	0.537*	<b>0.659*</b>
Hlip2	50908	<u>MRWTC</u> <u>AF</u> LWCVVVPTLHAW <del>VP</del> STTNPASRIGTRRWEALGDRLEEEPRMNP	0.847*	<b>0.921*</b>
PAP fibrillin I	49299	<u>MMRE</u> <u>QR</u> MLAILWAGLWFGGSGVHAW <del>Q</del> EPNLF <del>T</del> VPIQPSQKFSQGSTAKS	0.601*	<b>0.951*</b>
GAPDH	29119	<u>MKFS</u> <u>AA</u> TFAALVGSAAAYSSSFTGSALKSSASNDASMSMATGMGVNGFG	0.439*	<b>0.390*</b>
VDE-like1	62158	<u>MRFAWVAAGVVL</u> <u>TTT</u> <u>TQALV</u> <u>L</u> LDCTGMGETR <del>T</del> SGIRPIRGL <del>E</del> SNMARYA	0.760*	<b>0.858*</b>
<b>(b)</b>				
FbaC5	44147	<u>MRF</u> <u>S</u> LOSSLAVLLVLQASHAAAF <del>S</del> APVSSSNKNGIRSFAPLSMSLDKYA	<b>0.709*</b>	0.347
FCP	43697	<u>MKFAATILALIGSAAA</u> <u>F</u> APAQTSRASTSLQYAKEDLVGAIPPVGFDDPLG	<b>0.561*</b>	0.312
FCP	24872	<u>MKTAVIASLIAGAAA</u> <u>F</u> APAKNAARTSVATNMAFEDLGAQPPLGFFDPLG	<b>0.448*</b>	0.330
PGM	43607	<u>MKGYLFATWACLTISSNA</u> <u>TEA</u> <u>F</u> HRGPRAPCGLHASKLKTMDSGKLV <del>D</del> V	<b>0.380*</b>	0.375
DDE	62131	<u>MKFLGV</u> <u>T</u> SLCLWVSVNRENVSEAFAPRHQSLSRPSSRTTSAFSRAPILS	0.614	<b>0.737*</b>
PsbO	28739	<u>MKF</u> <u>T</u> AACSI <del>ALAASA</del> SAFAFIP <del>S</del> VSRTTDLMSIQKDLANVGKVAAGAL	0.322	<b>0.535*</b>
FCP	22033	<u>MKSIIFASLLTSAAA</u> <u>F</u> APASSSTTRTATPTALNEEF <del>CR</del> GYVGGESVEPMF	0.458*	<b>0.506</b>
FCP	18759	<u>MKFAILASMLSAAAA</u> <u>F</u> APASQAGKASVALNAEKSPAMPFLPYPENLKG <del>Y</del>	0.400*	<b>0.563</b>
FCP	18195	<u>MMRSTILAALASAAA</u> <u>F</u> APASMQSQRAGSVSLNAEEMSKSIPFLVKPDKL	0.459*	<b>0.629</b>
Hlip1	62173	<u>MLTLILMTRLSLSE</u> <u>S</u> FGVTT <del>PR</del> IFRPAPCRH <del>TR</del> PLK <del>TT</del> LRHSTLPSET	0.542*	<b>0.611</b>
<b>(c)</b>				
FCP	47183	<u>MKYAVFASLLASAAAFAPA</u> <u>AKPA</u> ASTSALNAEMSKSMPFLTAPKNTGGYV	<b>0.467</b>	0.430
FCP	44994	<u>MKLSLAILALCASTNA</u> <u>A</u> FAPSVSQRTPRDLAGVVAPTGFDPAGFAARAD	<b>0.676</b>	0.672
FCP	24752	<u>MKFVAVFALLASAAAFAPA</u> <u>Q</u> SARTSVATNMAFENEIGAQQPLGYWDPLG	<b>0.585</b>	0.330
Tpt1	25658	<u>MKVATTLTLAFICCA</u> <u>S</u> AFGLNGQTTSMKKGVDAGSKPMQAI <del>D</del> VQGNR	<b>0.673</b>	0.511
FBPC3	16659	<u>MFILKSPALWLLLYPVVA</u> <u>FTAARA</u> <u>N</u> SIRPAAALSVDLSSVEAVPSR <del>KT</del> K	0.587	<b>0.832</b>
HY2	34372	<u>MKLVP</u> <u>AW</u> TMTRNAIF <del>S</del> ARNPSHFLEVTALFCILIASGRGARTNTAFVNP	0.492	<b>0.758</b>
Tic110	62144	<u>MKFTGVALAVSLAQQQALS</u> <u>N</u> TGGGIVGAYTVSSPSFFTPKSFSSFVRGP	0.424	<b>0.535</b>
FCP	44082	<u>MKLVGVSVVIVVAVAV</u> <u>SL</u> SDSVAFAFAMHGGSHKLSNTALRVTFEDEL	<b>0.500</b>	0.231
chID	43587	<u>MVSSSKSTWMMAGACLILMA</u> <u>FQ</u> VQSFTFVPA <del>TR</del> ATSTVKRVAPAFMSAVA	<b>0.635</b>	0.486

highest value  
bold

**BOLD:** prediction by NN, UNDERLINED: prediction by HMM, GREY: conserved motif at signal peptide cleavage site, BLACK: hydrophobic (ACFGILMPVWY), GREEN: hydrophilic (NQST), BLUE: basic (HKR), RED: acidic (DE)  
\* predicted signal peptide cleavage site coincides with conserved motif

supplementary Fig. 6, caption on page 4

supplementary **Fig. 6** (pages 2 and 3): List of 81 plastid assigned, manually curated gene models from the *Phaeodactylum tricorutum* genome used to construct the sequence logos (Fig. 1) comparing SignalP's Neuronal networks (NN) and Hidden Markov models (HMM). **(a)** sequences with identical predictions between NN and HMM and coincidence with an "ASAFAP"-motif (grey) in both prediction models (62 sequences). **(b)** sequences with differing predictions from NN and HMM but coincidence with an "ASAFAP"-motif in one of the prediction models (10 sequences). **(c)** sequences with no coincidence between the predicted cleavage sites of both models and an "ASAFAP"-motif (9 sequences). Protein IDs refer to the first release of the *P. tricorutum* genome v1.0.