# PREDALGO, a new multi-subcellular localization prediction tool dedicated to Algae

Tardif, M., Atteia, A., Vallon, O., Specht, M., Cogne, G., Rolland, N., Brugière, S., Hippler, M., Ferro, M., Bruley, C., Peltier, G. and Cournac, L.*
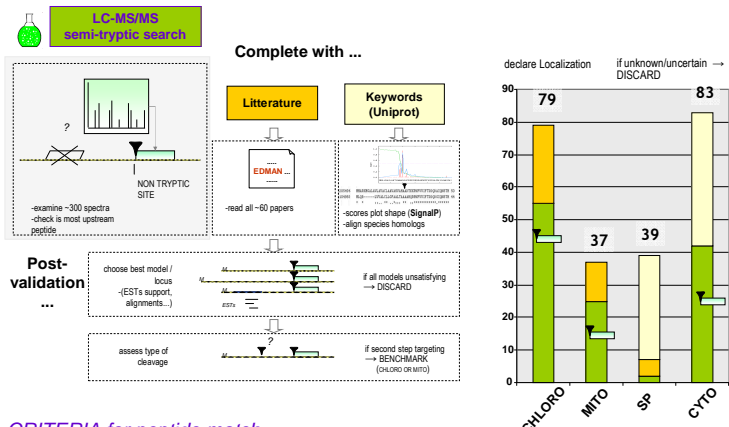
EDyP / CEA Grenoble

## Introduction

Despite of organella-specific proteomics efforts, assigning a subcellular localization to proteins still need to be assisted by **localization-predicting softwares**. However, the existing Plant-dedicated tools tend to mispredict algal chloroplast-localized proteins to mitochondria. We thus developed a tool **adapted to Algae**, using the feature of a cleavable N-terminal peptide present in proteins targeted to either of the mitochondria, chloroplast or endoplasmic reticulum compartments. PredAlgo was implemented using training sets of *Chlamydomonas reinhardtii* proteins, which principally consisted in **new N-terminal peptides of mature proteins identified from screening MS/MS data** of the recently published mitochondria and chloroplast proteomics surveys (1,2).
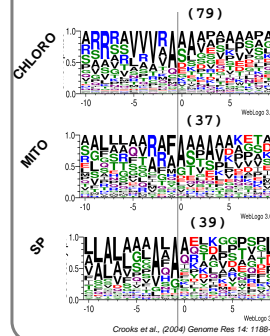
## 1  Identify maturation sites → Training sets

**LC-MS/MS semi-tryptic search**

Complete with ...

Litterature | Keywords (Uniprot)

EDMAN – read all ~60 papers

-examine ~300 spectra
-check is most upstream peptide

NON TRYPTIC SITE

-scores plot shape (SignalP)
-align species homologs

**Post-validation ...**

choose best model / locus -(ESTs support, alignments...)  →  if all models unsatisfying → DISCARD

assess type of cleavage  →  if second step targeting → BENCHMARK (CHLORO OR MITO)

declare Localization ; if unknown/uncertain → DISCARD



### CRITERIA for peptide match

- **non tryptic cleavage at Nter**
- **most upstream valid match on the protein**
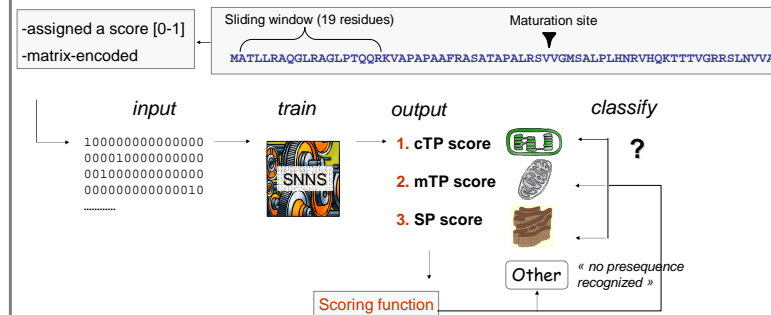- **within the first 150 amino acids**

The training Sets were finally constituted of candidates with MS/MS (green), Edman (orange) or SignalP prediction (yellow) **evidence of cleavage site**. The MS/MS fraction of the "CYTO" set was constituted by proteins for which the N-terminal of the sequence in the database was matched by a proteomic peptide, supporting the **absence of cleavage**.

## 2  Sequences around maturation sites



CHLORO (79), MITO (37), SP (39)

WebLogo 3.0

Crooks et al., (2004) Genome Res 14: 1188-1190

The sequences surrounding the cleavage site were converted into WebLogos. In the "chloro" and "mito" logos, negative charges predominated upstream of cleavage while practically no positively charged (acidic) residues were present. Too weak differential patterns between chloro and mito, make the application of **Artificial Neural Networks** a solution of choice.

## 3  Implementation of PredAlgo (Neural Network)

-assigned a score [0-1]
-matrix-encoded

Sliding window (19 residues)  Maturation site

MATLLRAQGLRAGLPTQQRKVAPAPAAFRASATAPALRSVVGMSALPLHNRVHQKTTTVGRRSLNVVA

input → train → output → classify

100000000000000
00010000000000
00100000000000
000000000000010
...........

SNNS

1. cTP score
2. mTP score
3. SP score

Scoring function → Other « no presequence recognized »

## 4  Performance evaluation in *C. reinhardtii*



$$Accuracy = \frac{TP+TN}{(TP+FN)+(TN+FP)}$$

$$Matthews\ Corr.\ Coeff. = \frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$$

PredAlgo was assessed on independant **Benchmark sets** of *C. reinhardtii* sequences, and compared to existing softwares. PredAlgo circumvented the weakness of most current tools as it generated outputs with **highly improved discrimination between the chloroplast and mitochondria**, resulting in 85% sensitivity for the chloroplast, 72% precision for the mitochondria. The precision calculated for the mitochondrion is not to be taken as an absolute criteria due to the small size of the "mito" set relatively to other sets. The discriminating power between compartments is better reflected by the MCC values which were fairly improved again for "Chloro" and "Mito" outputs (0.77 and 0.69).
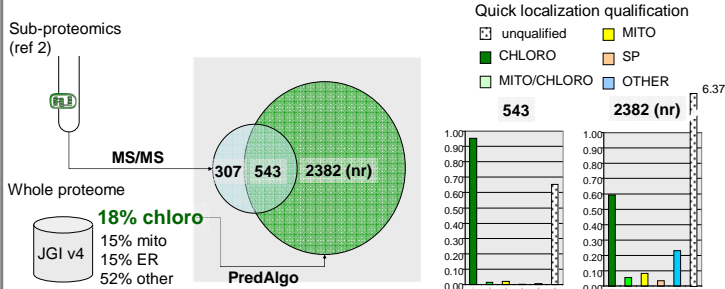
## 5  First application : the *C. reinhardtii* chloroplast proteome

Sub-proteomics (ref 2)

MS/MS

Whole proteome

JGI v4

**18% chloro**
15% mito
15% ER
52% other

307 | 543 | 2382 (nr)

PredAlgo

Quick localization qualification: unqualified, CHLORO, MITO/CHLORO, MITO, SP, OTHER

543 | 2382 (nr) | 6.37

PredAlgo predicted a chloroplast localization for 64% (543 / 850 ) of proteins previously identifed as chloroplastic by MS/MS (2) (while TargetP recovered only 19%, not shown). The core of a *C. reinhardtii* chloroplast could be delimited by the common set of 543 proteins. In addition, PredAlgo predicted 2382 (non redundant) proteins the vast majority has no obvious localization but which represent **potentially "new" chloroplastic proteins**.

## 6  Performances in other Green Algae



Volvox carteri, Ostreo. tauri, Chlorella, Ostreo. lucimarinus, Coccomyxa, Micro. pusilla

Trebouxiophyceae / Prasinophyceae

To assess whether PredAlgo is a **suitable tool for green algae in general**, orthologs pairs from whole protomes were computed as Blast Best Reciprocal Hits between *C. reinhardtii* and each of six Chlorophyta algal species. PredAlgo was then run on the algal proteomes and the predicted localization was compared to that of the *Chlamydomonas* ortholog. The counts shown here were restricted to sets of pairs for which the localization of the Chlamydomonas protein is known with certainty, i.e. to the Training and Benchmark sets.
i) Predictions in *Volvox carteri* correlated very well
ii) 'chloroplast' recall of at least 80% in all six algae
iii) Except for *Volvox*, consistency in 'mito' predictions was moderate (*Chlorella*, *Coccomyxa*) or not evaluable (*Ostreoccocus tauri*, *O. lucimarinus*, *Micromonas pusilla*)

## Conclusion

The identification of N-terminal peptides of mainly mitochondria- and chloroplast-targeted proteins via **MS/MS "semi-trypsic" strategies** brought up sufficient experimental data to implement a localization predictor specific to Algae. **PredAlgo** appeared as a **first position software to be used in *C. reinhardtii*** and its closest homolog **V. carteri** with a great discrimination capacity between chloroplast and mitochondria. It is also well-suited in **more distant algal species** if dedicated at chloroplast prediction recovery.
Noticeably, in organisms tested, a fair number of predictions errors were due to erroneous gene models. We assume that PredAlgo would perform even better as predicted models will become more accurate.

## Perspectives

- ☐ **Localization of *C. reinhardtii* Central metabolic enzymes**  (G. Cogne & co)
- ☐ **Mitochondria proteome refining in *C. reinhardtii***  (A. Atteia & co)
- ☐ ...

## References

1) Atteia A, et al. N (2009) Mol Biol Evol 26: 1533-1548
2) Terashima M, et al. (2010) Mol Cell Proteomics 9: 1514-1532

Abbreviations: cTP, chloroplast transit peptide; mTP, mitochondria transit peptide; SP, signal peptide; JGI v4, Joint Genome Institute, protein models version 4