A rational protein redesign method for improved secretion yields in Aspergillus niger

B.A. van den Berg^{*1,3,4}, M.J.T Reinders^{1,3,4}, J.M. van der Laan², J.A. Roubos², and D. de Ridder^{1,3,4}

The Delft Bioinformatics Lab, Faculty of Electrical Engineering, Mathematics & Computer Science, Delft University of Technology, Delft, The Netherlands, ² DSM Biotechnology Center, Delft, The Netherlands, ³ Netherlands Bioinformatics Centre, Nijmegen, The Netherlands, ⁴ Kluyver Centre for Genomics of Industrial Fermentation, Delft, The Netherlands * b.a.vandenberg@tudelft.nl

Introduction	1) Data set	2) Predicting high secretion yield
In industrial biotechnology, en- zymes are produced by over- expressing genes in production	345 A. niger proteins	amino acid weights Amino acid composition most predictive
hosts such as <i>Aspergillus niger</i> . Ideally, the produced protein should be secreted, such that it can easily be recovered from a	167 proteins: 187 proteins: low secretion yield high secretion yield	0.12 0.01 0.06A C D E0.27 -0.13 0.51 0.040.27 -0.13 0.51 0.05To predict high-yield secretion, many features have been explored. The amino acid composition was found to be most predictive. High-yield secretion prediction is calculated by taking the sum of a protein's weighted amino acid composition. The

reactor. Previously, we exploited a dataset of proteins to predict successful high-yield secretion. Here, we develop a method to rationally redesign the sequence of a low-yield protein to make it more similar to that of highyield proteins, with the aim to improve secretion yields.



All proteins are predicted to be secreted, i.e. contain a signal peptide and do not contain an ERretention signal or transmembrane regions. A high secretion yield is assigned in case of a high extracellular concentration after overexpression (visible band on gel), low secretion yield otherwise.



0.04

0.63

-0.28

-0.06

-0.25

-0.50

the sum of a protein's weighted amino acid composition. The amino acid weights are retrieved from a linear support vector machine that was trained using the proteins in the data set.

prediction outcomes data set 50



To reduce the probability that a mutation affects enzymatic activity, both buried (white) and ligandbinding residues (orange) are fixed.

=t∥tac_≠cze@agq

TNRA7N

≤ac c

🖉 aae 🎢 🗖 aat oac

cac

Constraints

Fixed residues

Design method Genetic algorithm

genetic algorithm is used to

Can secretion yields be improved by making the amino acid composition of a negative protein more similar to that of positive proteins?

composition To test this, we developed a protein design method that aims to alter the amino acid composition of a protein without affecting protein structure and enzymatic activity. The

Objective 1

Prediction outcome

The first objective is to increase the prediction outcome, i.e. promote mutations from amino acids with a negative weight to amino acids with a positive weight. As illustrated in the histogram above, the aim is to increase the prediction from 0.30 to around 1.10.

Objective 2

Often observed mutations

optimize the objectives under the given constrainst.

3) Protein redesign

designs will be tested in the lab.

The second objective promotes mutations that are often observed in homologous proteins. A multiple sequence alignment for a set of homologous proteins is used to create a position frequency matrix.

Multiple sequence alignment 55 415 *wt* ... EFGPTCIGVGEEISPLVGEDCLFINVFTPSHATTLSRLPVWVHIQ... *hom 1* ... EFGPICVGTGQSATSMRAEDCLFINVFTPSDATKHSKLPVWVFIQ... *hom 2* ... KRKPVCIGTSSDPIGTEDEDCLFLNIWAPTHASSKSKLPVYFYIQ... hom x ... HHGPVCYGVQELAV-PLSEDCLFIDVYAPSNATDSGRLPVMLWLQ... Position frequency matrix 1.0 of occu G 0.5 frequency 90 92 93 87 88 89 91 94 95 sequence position

4) In silico validation

Increasing sequence similarity compared to a protein from the postive set

As an *in silico* test, a negative protein (red) is redesigned with the number of mutations set to 30. The redesigned version (green) is afterwards compared to a positive protein (blue) with a similar structure. Mutated positions are shown below, in blue the residues at the mutated positions in the structurally aligned positive protein. The dots indicate identical residues, clearly showing an increasing number for the redesigned sequence, thereby showing an increasing sequence similarity with the positive protein.

11 10 D V A O D R A R A H A A P R R O A P P S T S R A E K D A D T T Q S Y T K YYGTTE Original negative (low secretion yield) Structarally similar positive(high secretion yield) **Redesigned sequence**

Objective 3

Amino acid composition

The third objective aims to optimize the similarity between the amino acid composition of the redesigned protein to the average amino acid composition of the proteins in the positive set.



tryptophan frequency

Mutations from rarely observed amino acids to frequently observed amino acids are prefered. For example, a mutation at position 89 from leucine (L) to tyrosine (Y) would be promoted, because Y occurs more often in homologous proteins at this position.

netherlands

centre





Kluyver Centre for Genomics of Industrial Fermentation