

A functional genomics study of extracellular protease production by Aspergillus niger

Braaksma, M.

Citation

Braaksma, M. (2010, December 15). *A functional genomics study of extracellular protease production by Aspergillus niger*. Retrieved from https://hdl.handle.net/1887/16246

Version:	Corrected Publisher's Version
License:	<u>Licence agreement concerning inclusion of doctoral</u> <u>thesis in the Institutional Repository of the University</u> <u>of Leiden</u>
Downloaded from:	https://hdl.handle.net/1887/16246

Note: To cite this publication please use the final published version (if applicable).

CHAPTER 4

METABOLOMICS AS A TOOL FOR TARGET IDENTIFICATION IN STRAIN IMPROVEMENT: THE INFLUENCE OF PHENOTYPE DEFINITION

Machtelt Braaksma, Sabina Bijlsma, Leon Coulier, Peter J. Punt and Mariët J. van der Werf

> This chapter has been accepted for publication in: Microbiology (2010), doi:10.1099/mic.0.041244-0

Supplementary data will be available with the online version of this paper at http://mic.sgmjournals.org/

ABSTRACT

For the optimization of microbial production processes, the choice of the quantitative phenotype to be optimized is crucial. For instance, for the optimization of product formation either product concentration or productivity can be pursued, potentially resulting in different targets for strain improvement. The choice of a quantitative phenotype is not only highly relevant for classical improvement approaches, but even more so for modern systems biology approaches.

In this study, the information content of a metabolomics data set was determined with respect to different quantitative phenotypes related to the formation of specific products. To this end, the production of two industrially relevant products by *Aspergillus niger* was evaluated; (i) the enzyme glucoamylase and (ii) the more complex product group of secreted proteases, consisting of multiple enzymes. For both products six quantitative phenotypes associated with activity and productivity were defined, taking also into account different time points of sampling during the fermentation. Both linear and non-linear relations between the metabolome data and the different quantitative phenotypes were considered.

The multivariate data analysis tool partial least squares (PLS) was used to evaluate the information content of the data sets for all the different quantitative phenotypes defined. Depending on the product studied, different quantitative phenotypes were found to have the highest information content in specific metabolomics data sets. A detailed analysis of the metabolites showing strong correlation with these quantitative phenotypes revealed that for glucoamylase activity various sugar-derivatives were found to be correlating. For the reduction of protease activity mainly as yet unidentified compounds were found to be correlating.

INTRODUCTION

The optimization of microbial production processes is an ongoing cycle of strain and/or process improvement. Traditionally, prior knowledge is the basis for identifying putative bottlenecks in the process. However, with the use of functional genomics technologies a more unbiased approach towards target selection for metabolic engineering or process optimization can be applied (van der Werf, 2005).

For optimization of the production process of a biological compound or enzymatic activity, a broad range of definitions of phenotypes can be selected for improvement. For instance, in studies reporting the production of glucoamylase by the filamentous fungus *Aspergillus niger* many different quantitative phenotypes for glucoamylase production were used. These included glucoamylase concentration (in g l⁻¹) (Withers *et al.*, 1998), activity (in U l⁻¹) (Wang *et al.*, 2008), yield (in mol product mol⁻¹ substrate) (Melzer *et al.*, 2007), specific concentration or activity (in g g⁻¹ DWT or U g⁻¹ DWT, respectively) (Swift *et al.*, 2000, Pedersen *et al.*, 2000; Schrickx *et al.*, 1993), and specific productivity (in mol, gram or units g⁻¹ DWT h⁻¹) (Melzer *et al.*, 2007; Withers *et al.*, 1998; Schrickx *et al.*, 1993).

The motivation for choosing a certain quantitative phenotype in bioprocess optimization is not always clear, and seems largely *ad libitum*. The choice of the quantitative phenotype to be pursued may have a major influence on the outcome of an optimization strategy. As stated by Kennedy & Krouse (1999) in their review on strategies for improving fermentation medium performance, some medium design studies flounder because the target variable to be improved is not clearly defined. Phenotype definition is not only important for classical optimization approaches, but perhaps even more so for modern, top-down systems biology approaches. In particular, as the enormous quantity of data that arise from these systems biology studies may easily result in a data overload (Braaksma *et al.*, 2010a). However, as far as we know, no systematic studies have been performed to study which quantitative phenotype is the most relevant in bioprocess optimization.

In bioprocess optimization a high quantity, e.g. concentration, of a product is not automatically the most desired result. In the case the substrate is an expensive part of the total fermentation costs, a high yield may be more relevant. However, improvement of the product yield is not always achieved by focussing on the yield itself during the strain improvement process. Focussing on the productivity may require fewer strain improvement steps during a particular bioprocess optimization process, thus resulting in an improved yield more quickly. Reduction of the fermentation time is another parameter to reduce production costs and can be realized by increasing the productivity. It is very likely that selection of either of these phenotypes for optimization will result in different targets to obtain the desired increase.

In this study, a metabolomics approach was used for target selection for process optimization and/or metabolic engineering of the host. Culture samples from *A. niger* fermentations were analyzed for the production of glucoamylase and protease. For both products different quantitative phenotypes associated with activity and productivity were defined. In a first step, we determined the information content of our metabolomics data set with respect to different quantitative phenotypes associated with the formation of either of the two different products. Subsequently, metabolites were identified showing the strongest correlation with the phenotype studied.

METHODS

Strain and cultivation conditions

Aspergillus niger N402, a cspA1 (conferring short conidiophores) derivative of ATCC 9029 (Bos et al., 1988), was used in this study.

Cultures were grown in batch fermentations in BioFlo 3000 (New Brunswick Scientific) bioreactors with a 5 litre working volume. Minimal medium (Bennett & Lasure, 1991) contained 7 mM KCl, 11 mM KH₂PO₄, 2 mM MgSO₄, 76 nM ZnSO₄, 178 nM H₃BO₃, 25 nM MnCl₂, 18 nM FeSO₄, 7.1 nM CoCl₂, 6.4 nM CuSO₄, 6.2 nM Na₂MoO₄ and 134 nM EDTA. This medium was supplemented with the appropriate carbon source or nitrogen source in concentrations as indicated below. To prevent foaming, 1 % (v/v) antifoam (Struktol J 673) was added to the medium and, when necessary, additional antifoam was added during the cultivation. The medium composition, cultivation conditions and operating procedure of the bioreactor have been described in detail previously (Braaksma *et al.*, 2009). Cultivations were performed according to a full factorial design (total 16 conditions, and 9 biological duplicates), varying the carbon source (277.5 mM glucose or 333.0 mM xylose), the nitrogen source (ammonium chloride or sodium nitrate), the nitrogen concentration (low (282.4 mM) or high (564.8 mM)), and the pH (4 or 5) (Braaksma *et al.*, 2009).

Enzyme assays

Protease activity. Extracellular proteolytic activities were measured at an assay pH of 4 as described previously (Braaksma *et al.*, 2009).

Glucoamylase activity. Glucoamylase activity was measured using PNPG (*p*-nitrophenyl α -D-glucopyranoside) (Sigma-Aldrich) as a substrate (Withers *et al.*, 1998). The procedure was fully automated using a COBAS MIRA Plus autoanalyser. 30 µl of cleared culture supernatant was incubated with 90 µl 0.1% (w/v) PNPG in 0.1 M sodium acetate buffer, pH 4.3, for 20 min. at 37 °C. The reaction was terminated by the addition of 135 µl 0.1 M borate buffer, pH 9.3, and the absorbance was read at 405 nm. One unit of glucoamylase activity was defined as the amount of enzyme that produces an absorbance at 405 nm equivalent to 1 µmol/l of *p*-nitrophenol in 1 minute under the given assay conditions.

Collection of samples, extraction and sample clean-up

Samples for metabolome analysis (25-100 ml, depending on the dry weight concentration) were taken rapidly from the bioreactor by closing the gas outlet and opening the sampling port. Cells were immediately quenched at -45 °C in methanol and collected as described previously (Pieterse *et al.*, 2006). Cell pellets were stored at -45 °C until use. To allow correlation of the metabolite concentrations to cell dry weight, the internal standards phenylalanine-d₅, leucine-d₃ (Spectral Stable Isotopes, Columbia, USA) and labelled ¹³C₁₀,¹⁵N₅-GTP (Sigma-Aldrich, Zwijndrecht, the Netherlands) were added prior to extraction. The intracellular metabolites were extracted from the cell suspensions by chloroform extraction at -45 °C as described by Ruijter and Visser (Ruijter & Visser, 1996). The water/methanol phase was subsequently divided in two portions, one for GC- and one for LC-MS analysis. The LC-MS sample was deproteinized by filtration using a Microcon YM-10 (Millipore) filter centrifuged at 18000 g and -20 °C for 16 hours. Subsequently, all samples were lyophilized. To allow correction for the recovery of amino acids, the group of metabolites most susceptible to matrix effects (i.e. the effect that in complex samples the detection of some compounds is disturbed in the presence of other compounds), prior to lyophilizing the samples for GC-MS an internal standard mixture of ²D,¹⁵N-labeled amino acids (Spectra Stable Isotopes) was added.

Biomass determination

Cell culture samples. For the quantification of cell dry weight (DWT), a known volume of cell culture was filtered though a dried, pre-weighted filter paper, followed by washing with distilled water twice and then drying at 110 °C for 24 h.

Metabolome samples. The extracted mycelium was collected and dried at 110 °C for 24 h to determine the dry weight of the sample (Ruiter & Visser, 1996). The metabolite concentrations in the extracts were correlated to dry weight by the use of the above mentioned internal standards added prior to the extraction of the cell pellets.

Analytical procedures

IP-LC-MS method. Lyophilized metabolome samples were dissolved in 100 µl methanol/water (1:3 v/v) and analyzed as described by Coulier *et al.* (Coulier *et al.*, 2006). Samples (10 or 20 µl) were separated on a reversed phase column (Chrompack Inertsil 5 mm ODS-3 100 x 3 mm, Middelburg, The Netherlands) using a 40 min linear gradient from 100% 5 mM hexylamine (pH 6.3) to 100% of 90% methanol-10 mM ammonium acetate (pH 8.5) at a flow rate of 0.4 ml min⁻¹. Compounds were detected by electrospray ionization (negative ion mode) in the range m/z 150/1000 using a Thermo Finnigan LTQ linear ion-trap system (Thermo Electron Corp. San Jose, USA). During data acquisition, the mass spectrometer probe voltage was maintained at 3–4 kV, the heated capillary was kept at 250 °C.

RP-LC-MS method. After analysis with the IP-LC-MS method, the redissolved metabolome samples were used for analysis with the RP-LC-MS method. Samples (10 or 20 μ l) were separated on a reversed phase column (Waters Sunfire C18, 150 x 3 mm, 3.5 μ m) using a linear gradient from 100% water + 0.1% formic acid to 75% MeCN/water (80%/20%) + 0.1% formic acid in 18 minutes followed by a linear gradient to 100% MeCN/water (80%/20%) + 0.1% formic acid in 10 minutes at a flow rate of 0.3 ml min⁻¹. Compounds were detected by electrospray ionization (ESI; positive ion mode) in the range m/z 150-2000.

OS-GC-MS method. Lyophilized metabolome samples were derivatized using a solution of ethoxyamine hydrochloride in pyridine as the oximation reagent followed by silylation with *N*-methyl-*N*-(trimethylsilyl)trifluoroacetamide (MSTFA) as described by Koek *et al.* (Koek *et al.*, 2006). Before silylation, dicyclohexylphthalate (Sigma-Aldrich) was added as an internal standard for injection. GC-MS-analysis of the derivatized samples was performed using a temperature gradient from 70 °C to 320 °C at a rate of 10 °C min⁻¹ on an Agilent 6890 N GC and an Agilent 5973 mass selective detector (Agilent, Palo Alto, USA). 1 µl

aliquots of the derivatized samples were injected splitless on a HP5-MS capillary column (30 m x 0.25 mm, 0.25 μ m film thickness, Agilent). Detection was performed using MS detection in electron impact mode (70 eV).

Data preprocessing

The LC-MS data were converted to .cdf-files and imported in Matlab (version 7.7.0.471 (R2008b), The Mathworks, Inc., Natick, MA). The homemade software packages Impress V1.2, Winlin V2.4 and Equest V2.3XP (Vogels *et al.*, 1996; van der Greef *et al.*, 2004) were used to align and peak-pick the LC-MS data. Following preprocessing, all peaks in the obtained target tables (in the form of peak identifiers [mass.retention time] and peak areas) were normalized with respect to the amount of extracted biomass per sample.

Also the data from the GC-MS analyses were converted into target tables, i.e. spreadsheets containing relative peak areas for all significant metabolite peaks in all samples. Peak areas were obtained by automated peak integration, followed by manual inspection. To several of the peaks a (partial) chemical identity could be assigned by comparing retention time and mass spectrum with an in-house database, otherwise a unique peak identifier [AN codes] was assigned. All peak areas were corrected for the recovery of the internal standard for injection. Subsequently, the amino acids were corrected for the recovery of the labeled amino acids. Finally, peaks were normalized with respect to the amount of extracted biomass per sample.

Both preprocessed LC-MS and GC-MS data files were combined in one data matrix. As the presence of values equal to zero can disturb the statistical analysis, prior to this, a so-called 25%-rule was applied: only those variables were retained which were present in at least 25% of the samples (Rubingh *et al.*, 2009; Bijlsma *et al.*, 2006). Next, all remaining zero values in the separate GC-MS, IP-LC-MS and RP-LC-MS data sets were replaced by a threshold value of half the lowest value in the data set unequal to zero (Rubingh *et al.*, 2009). In total 489 individual peaks, i.e. 131 GC-MS, 176 IP-LC-MS and 182 RP-LC-MS peaks, were retained in the final data sets to be used as input for multivariate data analysis MVDA.

Multivariate data analysis

Before data analysis, the curves with glucoamylase and protease activity were corrected for noise and possible outliers using a smoothing algorithm as described previously (Braaksma et al., 2009). The phenotype data, e.g. protease or glucoamylase activity or productivity, were mean-centred $[(x - \overline{x})]$ prior to MVDA in order to remove the overall offset from the data (van den Berg et al., 2006). The metabolome data set was mean-centred and, in order to compare the metabolites relative to the biological response range, it was subsequently range scaled $[(x_i - \bar{x})/(x_{max} - x_{min})]$ prior to MVDA (van den Berg *et al.*, 2006). PLS analysis were performed in the Matlab environment using the PLS Toolbox (version 5.0.3, 2008; Eigenvector Research, Manson, WA). The PLS results were cross-validated by using a tenfold single cross validation procedure. In addition to PLS analysis on the original metabolome and phenotype data, PLS analysis was also performed after either natural logarithm transformation of the phenotype data in combination with the original metabolome data or after natural logarithm transformation of the metabolome data in combination with the original phenotype data. An automatic procedure was written in Matlab code in order to run the many PLS models in a short time. Every generated PLS model was inspected manually to judge if the number of latent variables (LV's) chosen by the algorithm seemed appropriate with respect to the Root Mean Square Error of Cross Validation (RMSECV) curve. In general, if more LV's are included in the PLS model, the given model will contain more noise. In the case too many LV's were chosen by the algorithm, a new PLS model was generated by choosing a smaller number of LV's.

Compound identification

The identity of relevant peaks was established by verifying peak retention time and mass spectrum against in-house and public databases. If a peak could not be identified in this way, in several cases it was subsequently reanalyzed using high resolution and/or tandem mass spectrometry (MS/MS) analytical instruments (van der Werf *et al.*, 2007).

RESULTS

Experimental setup

In order to evaluate whether the definition of the phenotype used influences the outcome of a metabolomics study, or for that matter any optimization approach, the production of two industrially relevant products, i.e. glucoamylase and proteases, by *A. niger* was studied. To this end, *A. niger* was grown at sixteen different environmental conditions, with nine randomly selected biological duplicates (see also Braaksma *et al.*, 2009). Samples for metabolome analyses were taken at three different time points of the growth curve based on cell dry weight concentrations. One sample was collected at the middle of the logarithmic growth phase (mid log), one at the end of the logarithmic growth phase (late log) and one during the stationary growth phase. Samples were immediately quenched in a methanol solution to prevent alterations in the metabolite composition of the samples. Subsequently, the metabolites present were analyzed using three analytical methods (see Methods section).

The production of glucoamylase and protease was monitored during the course of the fermentation by analyzing culture samples every six hours. The variation in maximum protease and glucoamylase activities under the different experimental conditions is shown in Fig. 1. For protease activity the variation is evenly distributed over the different experimental conditions (Braaksma *et al.*, 2009). For glucoamylase the experiments can be clearly separated in two groups. One group with very low activities of conditions where the fungus was grown under non-induced conditions (on xylose) and another group with high activities of growth under induced conditions (on glucose).



Fig. 1. (A) Maximum protease activity and (B) maximum glucoamylase activity in the different fermentations.

Quantitative phenotypes

Six different quantitative phenotype values for the three different products were determined. Glucoamylase and protease were expressed as activity (see A in Fig. 2), and for both products the rate of production, i.e. the productivity (see B in Fig. 2), was calculated. However, the amount of product formed also depends on the biomass concentration (DWT). Therefore, specific activity and specific productivity were also determined. These two specific phenotypes were calculated using the DWT at the time point of sampling (see A1 and B1, respectively, in Fig. 2). However, when a sample was collected during the stationary phase of the fermentation, the biomass concentration may already be declining due to autolysis of the fungal cells (White *et al.*, 2002), thus making specific activity and specific productivity dependent on the degree of lysis. Therefore, both phenotypes were also calculated in relation to the maximum biomass concentration (DWT_{max}) (see A2 and B2, respectively, in Fig. 2). By using DWT_{max}, the phenotypic value is not artificially increased when in certain fermentations severe cell lysis had occurred. In addition to the phenotypes described above, similar quantitative phenotypes values were also calculated using the *maximum* activity or productivity for these products (see also Braaksma et al., 2009). Thus, in this latter case, for all three metabolome time samples the phenotypic value was identical. For a detailed description of how each phenotype was defined and a complete overview of the phenotypic values corresponding to each metabolome sample, see Supplementary data file 1.

Analysis of the information content of the data set

The multivariate data analysis (MVDA) tool partial least squares (PLS) was used to determine the information content of the metabolome data sets for all the different quantitative phenotypes defined. PLS is a regression tool that results in a model that describes a quantifiable phenotype of interest, such as protease activity or productivity, based on the concentrations of each of the metabolites determined. In MVDA analysis of metabolomics data it is important to realize that due to the relatively large number of variables and few number of samples, chance correlations are a serious issue. Therefore, the cross-validated correlation coefficient, R^2_{CV} , obtained from a PLS model after cross validation, is a better measure for the information content of a PLS model than the initial correlation coefficient R^2_{fit} , because R^2_{CV} also reflects the robustness of the model. A high R^2_{CV} indicates a high information content of the metabolome data in relation to the quantitative phenotype. In this study, cross validated PLS models with a R^2_{CV} of 0.6 or higher were considered good

statistical models. For both products, cross validated PLS models were made for all different quantitative phenotypes (Table 1).

To investigate whether the information content of the metabolomics data set was growth phase specific, PLS models of these six quantitative phenotypes were calculated by including the metabolome data of different time samples in the PLS model. PLS models were determined using metabolome data of all three samples generated from the different fermentations as well as with the metabolome data of only the samples collected at one of the growth phases during the fermentation. In addition, also PLS models were generated evaluating non-linear relations between the quantitative phenotype and the metabolome data, in order to identify metabolites with a non-linear relation to the studied phenotype. An overview of the PLS models generated from the metabolome data of this study, including the R²_{CV} of each model, is shown in Table 1.



Fig. 2. A schematic representation of production in time to illustrate the various product-related phenotypes that can be defined. Solid line, product; dashed line, biomass concentration DWT. (A) activity at time point of sampling; (A1) specific activity – 1, based on the biomass at the time point of sampling; (A2) specific activity – 2, based on the maximal biomass concentration during the fermentation; (B) productivity at time point of sampling; (B1) specific productivity – 1, based on the biomass at the time point of sampling; (B2) specific productivity – 2, based on the maximal biomass concentration during the fermentation. (Adapted from Braaksma *et al.* (2009), Microbiology 155, 3430-3439.)

Table 1. Overview of the cross validation values (R^2_{CV}) of the PLS models made for glucoamylase (A) and protease (B).

Models with a R²_{CV} of 0.6 or higher are considered good statistical models and are indicated in bold.

Glucoamylase							
Table 1A	Phenotype *	Р	R ² cv	LN(P)	R ² _{cv}	LN(M)	R ² _{cv}
	Max.Act.	G1	0.59	G49	0.66	G97	0.75
	Max.Spec.Act1	G2	0.47	G50	0.64	G98	0.64
Maximum phenotype,	Max.Spec.Act2	G3	0.59	G51	0.64	G99	0.77
metabolome data of all samples	Max.Prod.	G4	0.59	G52	0.67	G100	0.73
	Max.Spec.Prod1	G5	0.60	G53	0.63	G101	0.78
	Max.Spec.Prod2	G6	0.59	G54	0.65	G102	0.74
	Max.Act.	G7	0.71	G55	0.76	G103	0.77
	Max.Spec.Act1	G8	0.47	G56	0.72	G104	0.67
Maximum phenotype,	Max.Spec.Act2	G9	0.62	G57	0.74	G105	0.71
metabolome data of mid log samples	Max.Prod.	G10	0.79	G58	0.75	G106	0.82
	Max.Spec.Prod1	G11	0.71	G59	0.74	G107	0.82
	Max.Spec.Prod2	G12	0.73	G60	0.74	G108	0.82
	Max.Act.	G13	0.43	G61	0.63	G109	0.50
	Max.Spec.Act1	G14	0.42	G62	0.62	G110	0.48
Maximum phenotype,	Max.Spec.Act2	G15	0.43	G63	0.61	G111	0.49
metabolome data of late log samples	Max.Prod.	G16	0.60	G64	0.72	G112	0.57
	Max.Spec.Prod1	G17	0.67	G65	0.71	G113	0.66
	Max.Spec.Prod2	G18	0.61	G66	0.70	G114	0.58
	Max.Act.	G19	0.00	G67	0.01	G115	0.41
	Max.Spec.Act1	G20	0.01	G68	0.03	G116	0.45
Maximum phenotype,	Max.Spec.Act2	G21	0.03	G69	0.04	G117	0.44
metabolome data of stationary samples	Max.Prod.	G22	0.02	G70	0.01	G118	0.40
	Max.Spec.Prod1	G23	0.00	G71	0.03	G119	0.44
	Max.Spec.Prod2	G24	0.00	G72	0.02	G120	0.39
	Act.	G25	0.40	G73	0.68	G121	0.51
	Spec.Act1	G26	0.38	G/4	0.67	G122	0.48
Phenotype at time point of sampling,	Spec.Act2	G27	0.41	G75	0.66	G123	0.53
metabolome data of all samples	Prod.	G28	0.55	G/6	0.59	G124	0.69
	Spec.Prod1	G29	0.56	G77	0.57	G125	0.66
	Spec.Prod2	G30	0.59	G78	0.57	G126	0.68
	Act.	G31 C22	0.67	G79 C80	0.69	G12/	0.67
Bhonotype at time point of compling	Spec.Act1	632	0.05	G80 C81	0.69	6128	0.67
metabolomo data of mid log complex	Brod	633	0.03	601	0.05	G125 G120	0.07
metabolome data or mid log samples	Spec Prod -1	635	0.78	683	0.05	G130 G131	0.78
	Spec Prod -2 †	636	0.77	G84	0.70	6132	0.81
	Act	637	0.22	685	0.70	6132	0.30
	Spec.Act1	G38	0.22	G86	0.48	G134	0.33
Phenotype at time point of sampling.	Spec.Act2 †	G39	0.23	G87	0.48	G135	0.33
metabolome data of late log samples	Prod.	G40	0.33	G88	0.28	G136	0.42
······································	Spec.Prod1	G41	0.29	G89	0.25	G137	0.38
	Spec.Prod2 †	G42	0.29	G90	0.25	G138	0.38
	Act.	G43	0.04	G91	0.00	G139	0.34
	Spec.Act1	G44	0.01	G92	0.01	G140	0.34
Phenotype at time point of sampling,	Spec.Act2	G45	0.02	G93	0.01	G141	0.37
metabolome data of stationary samples	Prod.	G46	0.05	G94	0.02	G142	0.40
	Spec.Prod1	G47	0.01	G95	0.02	G143	0.38
	Spec.Prod2	G48	0.01	G96	0.02	G144	0.40

Table 1. Continued.

Table 18Phenotype*PR* rvLN(P)R* rvLN(M)R* rvMax MaxSpec.Act.1P10.70P490.75P370.78MaxSpec.Act.2P30.57P510.66P590.66metabolone data of all samplesMax.Poct.2P30.57P510.68P520.69P1000.63MaxSpec.Prod.2P50.58P530.50P1010.63Max.Spec.Prod.2P60.88P550.58P1020.65metabolone data of nid log samplesMax.Act.P170.46P550.72P1030.32MaxSpec.Prod.1P100.51P580.68P1060.38P1060.38MaxSpec.Prod.2P120.29P590.44P1070.160.88MaxSpec.Prod.2P120.23P600.45P1080.18MaxSpec.Prod.2P120.24P630.49P1110.47MaxSpec.Prod.2P120.42P630.49P1110.47MaxSpec.Prod.2P150.42P630.49P1110.47MaxSpec.Prod.2P150.42P630.49P1110.47MaxSpec.Prod.2P140.58P660.34P1110.47MaxSpec.Prod.2P150.42P630.49P1110.41MaxSpec.Prod.2P140.13P650.36P130.11MaxSpec.Prod.2P170.14	Protease							
Max.ht. P1 0.70 P49 0.75 P97 0.78 MaxSpec.At.1. P2 0.66 P99 0.66 P99 0.66 metabolome data of all samples Max.Frod. P4 0.71 P52 0.69 P100 0.80 MaxSpec.Prod.1 P5 0.58 P54 0.48 P102 0.65 MaxSpec.Prod.2 P6 0.58 P54 0.57 P103 0.47 MaxSpec.Prod.1 P8 0.38 P56 0.58 P104 0.32 metabolome data of mid log samples MaxSpec.Prod.1 P11 0.28 P57 0.46 P107 0.16 MaxSpec.Prod.2 P12 0.29 P60 0.45 P108 0.18 MaxSpec.Prod.1 P11 0.32 P61 0.65 P009 0.662 MaxSpec.Prod.1 P11 0.32 P60 0.48 P110 0.47 metabolome data of late log samples MaxSpec.Prod.2 P12 0.44 P10 </th <th>Table 1B</th> <th>Phenotype *</th> <th>Р</th> <th>R²_{cv}</th> <th>LN(P)</th> <th>R²cv</th> <th>LN(M)</th> <th>R²cv</th>	Table 1B	Phenotype *	Р	R ² _{cv}	LN(P)	R ² cv	LN(M)	R ² cv
Maxinum phenotype, metabolome data of all samplesMax Spec.Arc.1P20.60P500.660.72Max Spec.Arc.1P40.71P520.69P1000.80Max Spec.Prol.2P50.82P530.83P530.83P530.83Max Spec.Prol.2P50.84P550.72P1030.72Max Spec.Arc.1P70.46P550.72P1030.72Max Spec.Arc.1P10.83P560.58P1040.32Max Spec.Arc.2P90.28P570.55P1050.28Max Spec.Arc.1P110.72P610.55P1050.28Max Spec.Arc.1P110.72P610.55P1050.28Max Spec.Arc.1P110.72P100.44P1070.16Max Spec.Arc.1P110.72P100.44P1070.16Max Spec.Arc.1P140.72P100.44P1070.11Max Spec.Arc.1P140.72P100.44P1110.47Max Prod.P140.12P100.44P1120.11Max Spec.Arc.1P140.13P620.44P1140.11Max Spec.Arc.1P140.11P670.14P1140.11Max Spec.Arc.2P240.11P670.14P1140.11Max Spec.Arc.1P120.11P670.14P1140.11Max Spec.Arc.1		Max.Act.	P1	0.70	P49	0.75	P97	0.78
Maximum phenotype, metabolome data of all samplesMax.Spec.Arcl.2P30.57P510.60P990.60Max.Spec.Prod1P50.58P530.50P1010.63Max.Spec.Prod2P60.68P550.72P1030.47Max.Spec.Arcl.2P90.28P550.72P1050.28metabolome data of mid log samplesMax.Spec.Arcl.2P90.28P570.55P1050.28Max.Spec.Arcl.2P100.51P580.66P1070.16Max.Spec.Arcl.3P110.52P600.44P1070.16Max.Spec.Arcl.4P130.52P610.65P1090.62Max.Spec.Arcl.2P150.42P630.69P1110.47Max.Spec.Arcl.4P130.52P610.65P1090.62Max.Spec.Arcl.4P130.52P610.69P1110.47Max.Spec.Arcl.4P130.62P660.30P1130.47Max.Spec.Arcl.4P170.28P660.30P1130.47Max.Spec.Arcl.4P170.28P660.30P1130.61Max.Spec.Arcl.4P160.42P160.44P170.58Max.Spec.Arcl.4P170.28P660.34P1140.36Max.Spec.Arcl.4P130.11P680.64P1140.58Max.Spec.Arcl.4P210.11P680.64P11		Max.Spec.Act1	P2	0.66	P50	0.66	P98	0.72
metabolome data of all samplesMax.Sprc.Prod1P40.71P520.69P1000.63Max.Sprc.Prod2P60.58P540.48P1020.63Max.Sprc.Prod2P60.58P560.58P1040.32Max.Sprc.Art1P80.38P560.55P1050.38metabolome data of mid log samplesMax.Sprc.Art1P100.51P580.65P1060.43Max.Sprc.Prod1P110.28P590.44P1070.160.45P1080.46Max.Sprc.Prod1P110.28P590.44P1070.160.45P1080.46Max.Sprc.Prod1P110.28P510.42P1080.460.45P1080.46Max.Sprc.Art1P140.37P620.48P1110.470.44P1130.47Max.Sprc.Art2P150.42P630.48P1130.470.44P130.17Max.Sprc.Art1P140.28P640.59P1130.17P110.48P1140.40Max.Sprc.Art1P160.48P640.59P1150.48P1150.48P1150.48Max.Sprc.Art1P190.11P660.44P1150.46P140.46P140.46Max.Sprc.Art1P200.11P690.14P1150.58P1560.58P1560.58P1560.58P156	Maximum phenotype,	Max.Spec.Act2	P3	0.57	P51	0.60	P99	0.66
Max.Spcc.Prod1 P5 0.58 P53 0.50 P102 0.65 Max.Spcc.Prod2 P6 0.88 P54 0.84 P102 0.65 Max.Max.Ct. P7 0.46 P55 0.72 P103 0.47 Max.Spcc.At1 P8 0.38 P56 0.58 P104 0.22 metabolome data of mid log samples Max.Spcc.Prod1 P11 0.28 P59 0.45 P108 0.18 Max.Spcc.Prod1 P11 0.28 P50 0.45 P109 0.16 Max.Spcc.Prod1 P11 0.28 P61 0.45 P109 0.16 Max.Spcc.Prod1 P11 0.27 P62 0.48 P111 0.47 Max.Spcc.Prod1 P17 0.28 P65 0.30 P113 0.47 Max.Spcc.Prod1 P17 0.28 P65 0.30 P113 0.47 Max.Spcc.Prod1 P17 0.28 P16 0.48 P114 0.18	metabolome data of all samples	Max.Prod.	P4	0.71	P52	0.69	P100	0.80
Imach max Spec.Prod2P60.58P550.78P1020.67Max.LP70.46P550.72P1030.32MaxSpec.Act2P90.28P570.55P1050.38metabolome data of mid log samplesMax.Spec.Prod1P110.28P590.44P1070.16MaxSpec.Prod1P110.28P590.44P1070.160.18MaxSpec.Prod1P110.28P600.65P1090.62MaxSpec.Prod1P140.37P620.48P1010.47MaxSpec.Prod1P150.42P630.49P1110.47metabolome data of late log samplesMax.Spec.Prod1P160.48P640.59P1120.41MaxSpec.Prod1P170.28P660.34P1140.37P670.19P1130.47metabolome data of late log samplesMax.Spec.Prod2P180.29P660.34P1140.37MaxSpec.Prod1P170.28P660.34P1140.37P120.44MaxSpec.Prod1P180.29P660.34P1130.59metabolome data of stationary samplesMax.Spec.Prod1P220.17P700.14P1130.59MaxSpec.Prod1P230.14P710.59P1240.57P1240.57Phenotype,MaxSpec.Prod1P250.66P740.66P124		Max.Spec.Prod1	P5	0.58	P53	0.50	P101	0.63
Max.htt.P70.46P750.72P1030.43Max.prec.ht.1P80.38P560.52P1050.32metabolome data of mid log samplesMax.Spec.Att.2P90.28P570.55P1050.43Max.Prod.P100.51P580.69P1050.160.43Max.Spec.Prod.12P110.22P600.45P1080.18Max.Spec.Prod.2P120.23P600.45P1090.62Max.Spec.Att.2P150.42P650.30P1110.47Max.Spec.Att.1P170.28P650.30P1110.47Max.Spec.Att.2P150.42P650.30P1130.41Max.Spec.Prod.1P170.28P650.30P1130.41Max.Spec.Prod.1P170.28P650.30P1130.41Max.Spec.Prod.1P170.28P650.30P1130.41Max.Spec.Prod.1P120.11P660.34P1140.58Max.Spec.Prod.1P220.11P660.34P1130.44Max.Spec.Prod.1P230.14P770.19P1130.44Max.Spec.Prod.2P240.18P770.14P130.44Max.Spec.Prod.1P230.14P770.19P1130.44Max.Spec.Prod.1P240.36P770.57P1410.52Max.Spec.Prod.1P23 <t< th=""><th></th><th>Max.Spec.Prod2</th><th>P6</th><th>0.58</th><th>P54</th><th>0.48</th><th>P102</th><th>0.65</th></t<>		Max.Spec.Prod2	P6	0.58	P54	0.48	P102	0.65
Maximup henotype, metabolome data of mid log samplesMax.Spec.Act1 Max.Spec.Prod1P100.28 P570.55P1050.28 0.28 P57Max.Spec.Prod1P110.28 P58P590.44P1070.16 0.18Max.Spec.Prod1P110.28 P59P560.44P1090.16Max.Spec.Prod1P110.28 P56P560.48P1090.62Max.Spec.Act1P140.37P660.48P1100.47matabolome data of late log samplesMax.Spec.Act2P150.42P660.30P1120.44Max.Spec.Prod1P170.28P660.30P1130.11		Max.Act.	P7	0.46	P55	0.72	P103	0.47
Max.bype, metabolome data of mid log samplesMax.Spec.Arcd. Max.Spec.Prod1P100.28P730.55P1050.28metabolome data of mid log samplesMax.Spec.Prod2P110.28P590.44P1070.16Max.Spec.Prod2P120.29P610.65P1080.13Max.Max.P130.32P610.65P1080.47Max.Spec.Arct2P150.42P630.48P110.47Max.Spec.Arct2P150.42P630.49P1110.47Max.Spec.Prod1P170.28P630.30P1130.17Max.Spec.Prod1P170.28P630.30P1130.17Max.Spec.Prod1P190.11P670.19P1150.68Max.Spec.Prod1P200.11P630.25P1160.58Max.Spec.Arct.2P200.11P680.25P1160.58Max.Spec.Prod1P220.17P700.14P1180.60Max.Spec.Prod1P230.14P710.190.44Max.Spec.Prod1P230.14P710.14P1190.44Max.Spec.Prod1P240.18P700.14P120.44Max.Spec.Prod1P250.70P730.44P120.57Metabolome data of all samplesSpec.Arc1.1P260.66P120.140.12Prenotype at time point of sampling, meta		Max.Spec.Act1	P8	0.38	P56	0.58	P104	0.32
metabolome data of mid log samplesMax.Prod.P100.51P830.69P1060.43Max.Spec.Prod1P110.28P590.44P1080.18Max.Spec.Prod2P120.29P600.45P1080.18Max.Spec.At1P130.52P610.65P1090.62Max.Spec.At2P150.42P630.49P1110.47metabolome data of late log samplesMax.Spec.Prod1P170.48P660.30P1130.17Max.Spec.Prod1P170.11P670.19P1160.58Max.Spec.Prod2P180.29P660.34P1160.59Max.Spec.At1P200.11P670.19P1160.58Max.Spec.At1P210.11P670.14P110.59Max.Spec.At1P220.11P680.25P100.47Max.Spec.At1P220.11P690.14P1190.44Max.Spec.At1P230.14P110.590.77Max.Spec.At1P250.70P730.57P1210.50Max.Spec.Prod1P260.65P740.64P1220.77Phenotype at time point of sampling, metabolome data of all samplesSpec.At1P320.31P340.51Spec.At1P320.36P780.44P120.510.51Phenotype at time point of sampling, metabolome data	Maximum phenotype,	Max.Spec.Act2	P9	0.28	P57	0.55	P105	0.28
Max.Spec.Prod1P110.28P590.44P1070.16Max.Spec.Prod2P120.29P600.45P1080.81Max.Max.P130.37P620.48P1090.62metabolome data of late log samplesMax.Spec.Atc2P150.42P630.49P1110.47Max.Prod.P150.42P630.30P1120.440.59P1120.44Max.Spec.Prod1P170.28P660.30P1130.170.160.160.160.160.160.160.170.170.160.170.160.580.30P1130.170.350.44P1150.680.30P1140.180.170.55P1160.580.36P1140.180.690.44P1170.590.14P1180.600.580.580.14P1170.590.14P1180.600.580.550.16P130.41P1180.600.44P1180.600.560.16P140.44P120.750.570.57P1210.570.570.57P1210.560.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1240.510.550.56P1	metabolome data of mid log samples	Max.Prod.	P10	0.51	P58	0.69	P106	0.43
Max.Spec.Prod2 P12 0.29 P60 0.45 P108 0.18 Max.Act. P13 0.52 P61 0.65 P109 0.62 Max.Spec.Act1 P14 0.37 P62 0.48 P110 0.47 metabolome data of late log samples Max.Spec.Arct2 P15 0.42 P63 0.49 P112 0.44 Max.Spec.Prod1 P17 0.28 P65 0.30 P113 0.17 Max.Spec.Prod2 P18 0.29 P66 0.34 P144 0.18 Max.Spec.Arct1 P20 0.11 P66 0.34 P114 0.68 Max.Spec.Arct1 P20 0.11 P68 0.25 P116 0.58 Max.Spec.Arct2 P21 0.11 P69 0.14 P117 0.59 metabolome data of stationary samples Max.Spec.Prod2 P24 0.18 P22 0.37 Max.Spec.Prod2 P24 0.18 P710 0.19 P119		Max.Spec.Prod1	P11	0.28	P59	0.44	P107	0.16
Max.Act.Pi30.52P610.65P1090.62Max.Spec.Act1P140.37P620.48P110.47metabolome data of late log samplesMax.Spec.Act2P150.42P630.49P1110.47Max.Spec.Prod1P170.28P650.30P1120.44Max.Spec.Prod2P180.29P660.34P1140.17Max.Spec.Prod2P180.29P660.34P1150.58Max.Spec.Act2P210.11P670.19P1150.58Max.Spec.Act1P200.11P680.25P1160.58Max.Spec.Prod1P220.17P700.14P1180.60Max.Spec.Prod1P220.17P700.14P1180.60Max.Spec.Prod1P240.18P710.18P1200.47Max.Spec.Prod1P240.66P730.57P120.57Max.Spec.Prod1P260.66P730.64P120.57Max.Spec.Prod1P260.66P740.65P1240.65Spec.Prod.1P270.32P760.55P1240.51Spec.Prod.1P280.35P760.55P1240.51Spec.Prod.1P320.30P800.01P1250.55Prod.P330.31P310.01P1260.51Spec.Prod.1P350.16P8		Max.Spec.Prod2	P12	0.29	P60	0.45	P108	0.18
Max, Spec, Act1Pi40.37P620.48P100.47Max, Spec, Act2P150.420.430.430.430.47metabolome data of late log samplesMax, Prod.P170.28P630.30P1130.17Max, Spec, Prod1P170.28P650.30P1140.180.17Max, Spec, Prod2P180.29P660.30P1140.180.17Max, Spec, Act1P200.11P670.19P1150.58Max, Spec, Act1P200.11P680.25P1160.58Max, Spec, Act2P210.11P690.14P1170.59metabolome data of stationary samplesMax, Spec, Prod1P230.14P710.19P1180.60Max, Spec, Prod1P230.14P710.19P1190.44P1180.60Max, Spec, Prod1P230.14P710.19P1190.44Max, Spec, Prod1P230.14P710.19P1190.44Max, Spec, Prod1P260.66P740.14P1220.57Phenotype at time point of sampling, metabolome data of all samplesAct.P230.32P770.49P1230.31Phenotype at time point of sampling, metabolome data of all samplesAct.P330.03P810.01P1270.17Phenotype at time point of sampling, metabolome data of all log samplesAct		Max.Act.	P13	0.52	P61	0.65	P109	0.62
Maxium phenotype, metabolome data of late log samplesMax.Spec.Arc.1-2P150.48P630.49P1110.47Max.Spec.Prod1P160.28P650.30P1130.17Max.Spec.Prod2P180.29P660.34P1140.18Max.Max.LP190.11P670.19P1150.58Maximum phenotype, metabolome data of stationary samplesMax.Spec.Arc.1-P200.11P680.25P1160.58Max.Spec.Arc.1P220.11P690.14P1170.590.44P1170.59Max.Spec.Prod1P230.14P120.14P1190.440.47Max.Spec.Prod1P240.18P720.18P1200.47Max.Spec.Prod1P240.18P720.18P1200.47Phenotype at time point of sampling, metabolome data of all samplesSpec.Arc.1P250.66P750.44P1230.77Pred.P280.45P760.65P1240.6591240.650.6591240.65Phenotype at time point of sampling, metabolome data of all samplesAct.P310.03P780.01P1250.45Phenotype at time point of sampling, metabolome data of ind log samplesAct.P310.02P850.92P130.11Phenotype at time point of sampling, metabolome data of ind log samplesAct.P340.24P1320.120.12 <th></th> <th>Max.Spec.Act1</th> <th>P14</th> <th>0.37</th> <th>P62</th> <th>0.48</th> <th>P110</th> <th>0.47</th>		Max.Spec.Act1	P14	0.37	P62	0.48	P110	0.47
metabolome data of late log samples Max.Prod. P16 0.48 P64 0.59 P112 0.14 Max.Spec.Prod1 P17 0.28 P65 0.30 P113 0.17 Max.Spec.Prod2 P18 0.29 P66 0.34 P114 0.18 Max.Spec.Act1 P10 0.11 P67 0.19 P115 0.68 Max.Spec.Act1 P20 0.11 P68 0.25 P116 0.59 metabolome data of stationary samples Max.Spec.Act1 P20 0.11 P69 0.14 P11 0.19 P118 0.60 Max.Spec.Prod1 P23 0.14 P71 0.18 P122 0.75 Max.Spec.Prod1 P26 0.66 P74 0.66 P124 0.61 metabolome data of all samples Spec.Act1 P26 0.66 P74 0.64 P122 0.75 Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Prod.1	Maximum phenotype,	Max.Spec.Act2	P15	0.42	P63	0.49	P111	0.47
Max.Spec.Prod1 P17 0.28 P65 0.30 P113 0.17 Max.Spec.Prod2 P18 0.29 P66 0.34 P114 0.18 Max.Max. P19 0.11 P67 0.19 P115 0.68 Max.Spec.Act2 P21 0.11 P69 0.14 P113 0.59 metabolome data of stationary samples Max.Spec.Act2 P21 0.17 P70 0.14 P113 0.59 Max.Spec.Prod1 P23 0.14 P113 0.60 0.47 P114 0.60 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Prod2 P26 0.66 P74 0.46 P122 0.75 Phenotype at time point of sampling, Spec.Act2 P27 0.67 P75 0.44 P125 0.45 Spec.Prod1 P29 0.32 P77 0.49 P125 0.45 Phenotype at time point of sampling, Spec.Act1 <	metabolome data of late log samples	Max.Prod.	P16	0.48	P64	0.59	P112	0.44
Max.Spec.Prod2 P18 0.29 P66 0.34 P114 0.18 Max.Act. P19 0.11 P67 0.19 P115 0.68 Max.Spec.Act1 P20 0.11 P68 0.25 P116 0.58 metabolome data of stationary samples Max.Spec.Act2 P21 0.11 P69 0.14 P117 0.59 metabolome data of stationary samples Max.Spec.Act2 P21 0.18 P72 0.17 P70 0.14 P118 0.60 Max.Spec.Prod1 P23 0.14 P71 0.19 P119 0.44 Max.Spec.Prod2 P24 0.18 P72 0.67 P73 0.57 P121 0.80 metabolome data of all samples Spec.Act2 P27 0.67 P75 0.44 P123 0.75 prod. P28 0.45 P76 0.65 P124 0.61 Spec.Prod1 P31 0.09 P79 0.44 P125 0.45		Max.Spec.Prod1	P17	0.28	P65	0.30	P113	0.17
Max.Act.P190.11P670.19P1150.68Max.Spec.Act2P200.11P680.25P1160.58metabolome data of stationary samplesMax.Spec.Act2P210.11P700.14P1180.60Max.Spec.Prod1P230.14P710.19P1190.44Max.Spec.Prod2P240.18P720.16P1220.77Max.Spec.Prod1P250.70P730.57P1210.80Spec.Act2P270.66P740.46P1220.77metabolome data of all samplesSpec.Act2P270.67P750.44P1230.77metabolome data of all samplesSpec.Act2P270.67P750.44P1230.77metabolome data of all samplesSpec.Act2P270.67P760.44P1230.45Spec.Prod1P260.65P740.61P120.610.61Spec.Prod1P330.35P780.44P120.17Max.Spec.Prod1P330.03P810.01P120.17Max.Spec.Prod1P330.03P810.01P120.17Max.Spec.Prod1P340.21P850.24P1300.18Spec.Art.2*P330.03P810.01P120.17Max.Spec.Art.2*P330.03P810.01P120.12Max.Spec.Art.4P39 <td< th=""><th></th><th>Max.Spec.Prod2</th><th>P18</th><th>0.29</th><th>P66</th><th>0.34</th><th>P114</th><th>0.18</th></td<>		Max.Spec.Prod2	P18	0.29	P66	0.34	P114	0.18
Max.Spec.Act1 P20 0.11 P68 0.25 P116 0.58 Max.Spec.Act2 P21 0.11 P69 0.14 P117 0.59 metabolome data of stationary samples Max.Spec.Prod1 P23 0.14 P71 0.19 P118 0.600 Max.Spec.Prod1 P23 0.14 P71 0.19 P119 0.44 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Prod. P25 0.70 P73 0.57 P124 0.75 Phenotype at time point of sampling, metabolome data of all samples Spec.Act2 P27 0.67 P124 0.45 Spec.Act.1 P30 0.45 P76 0.44 P125 0.45 Spec.Act.2 P30 0.36 P78 0.48 P125 0.45 Spec.Act.1 P31 0.09 P79 0.01 P127 <th></th> <th>Max.Act.</th> <th>P19</th> <th>0.11</th> <th>P67</th> <th>0.19</th> <th>P115</th> <th>0.68</th>		Max.Act.	P19	0.11	P67	0.19	P115	0.68
Maximum phenotype, metabolome data of stationary samples Max.Spec.Act2 P21 0.11 P69 0.14 P11 0.59 Max.Spec.Prod1 P22 0.17 P70 0.14 P118 0.60 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.44 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Art1 P25 0.70 P73 0.57 P121 0.80 Spec.Art1 P26 0.66 P74 0.46 P122 0.75 metabolome data of all samples Spec.Art2 P27 0.67 P75 0.44 P123 0.77 metabolome data of all samples Spec.Art1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod1 P32 0.36 P78 0.01 P129 0.05 metabolome data of mid log samples		Max.Spec.Act1	P20	0.11	P68	0.25	P116	0.58
metabolome data of stationary samples Max.Prod. P22 0.17 P70 0.14 P118 0.60 Max.Spec.Prod1 P23 0.14 P71 0.19 P119 0.44 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Prod2 P24 0.18 P72 0.16 P121 0.80 Phenotype at time point of sampling, Spec.Act2 P27 0.67 P75 0.44 P122 0.77 metabolome data of all samples Spec.Prod1 P28 0.32 P77 0.49 P123 0.77 metabolome data of all samples Spec.Prod1 P29 0.36 P76 0.65 P124 0.61 Spec.Prod1 P23 0.36 P78 0.48 P125 0.45 Spec.Prod1 P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P33 0.33 P81 0.01 P129 0.55 metabolome data of mi	Maximum phenotype,	Max.Spec.Act2	P21	0.11	P69	0.14	P117	0.59
Max.Spec.Prod1 P23 0.14 P71 0.19 P119 0.44 Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Max.Spec.Prod2 P25 0.70 P73 0.57 P121 0.80 Phenotype at time point of sampling, Spec.Act1 P26 0.66 P74 0.46 P122 0.77 metabolome data of all samples Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Prod1 P29 0.32 P77 0.48 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P125 0.45 Phenotype at time point of sampling, Spec.Act2 † P33 0.03 P80 0.01 P128 0.05 Phenotype at time point of sampling, Spec.Act2 † P33 0.03 P81 0.01 P129 0.05 Prod. P34 0.21 P83 0.24 P130 0.16 Spec.Prod1	metabolome data of stationary samples	Max.Prod.	P22	0.17	P70	0.14	P118	0.60
Max.Spec.Prod2 P24 0.18 P72 0.18 P120 0.47 Act. P25 0.70 P73 0.57 P121 0.80 Spec.Act1 P26 0.66 P74 0.46 P122 0.75 metabolome data of all samples Spec.Act2 P27 0.67 P75 0.44 P123 0.77 metabolome data of all samples Spec.Act2 P27 0.67 P75 0.44 P123 0.75 Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Act1 P29 0.32 P77 0.49 P125 0.45 Spec.Act1 P32 0.03 P78 0.48 P126 0.45 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 Spec.Act1 P33 0.03 P81 0.01 P129 0.05 metabolome data of mid log samples Spec.Act1 P33 0.16 P84 0.24		Max.Spec.Prod1	P23	0.14	P71	0.19	P119	0.44
Act. P25 0.70 P73 0.57 P121 0.80 Spec.Act1 P26 0.66 P74 0.46 P122 0.75 metabolome data of all samples Spec.Act2 P27 0.67 P75 0.44 P123 0.77 metabolome data of all samples Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Act1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P126 0.45 Phenotype at time point of sampling, Spec.Act1 P32 0.03 P80 0.01 P127 0.17 metabolome data of mid log samples Spec.Act2 † P33 0.03 P81 0.01 P129 0.05 Spec.Prod1 P34 0.21 P82 0.42 P130 0.18 Spec.Prod1 P35 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling,		Max.Spec.Prod2	P24	0.18	P72	0.18	P120	0.47
Spec.Act1 P26 0.66 P74 0.46 P122 0.75 Phenotype at time point of sampling, metabolome data of all samples Spec.Act2 P27 0.67 P75 0.44 P123 0.77 metabolome data of all samples Prod. P28 0.45 P76 0.44 P124 0.61 Spec.Prod1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P126 0.45 Phenotype at time point of sampling, metabolome data of mid log samples Act. P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 metabolome data of mid log samples Spec.Act2 † P33 0.21 P82 0.24 P131 0.12 Spec.Prod1 P35 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling, Spec.Act1 P38 0.25 P86 0.09 <th></th> <th>Act.</th> <th>P25</th> <th>0.70</th> <th>P73</th> <th>0.57</th> <th>P121</th> <th>0.80</th>		Act.	P25	0.70	P73	0.57	P121	0.80
Phenotype at time point of sampling, metabolome data of all samples Spec.Act2 P27 0.67 P76 0.48 P123 0.77 metabolome data of all samples Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Prod1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P126 0.45 Phenotype at time point of sampling, metabolome data of mid log samples Act. P31 0.09 P79 0.01 P128 0.05 Phenotype at time point of sampling, Spec.Act1 P32 0.03 P81 0.01 P128 0.05 Phenotype at time point of sampling, Spec.Act2 † P33 0.03 P81 0.01 P129 0.05 Spec.Prod1 P35 0.16 P83 0.24 P131 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P132 0.12 Spec.Prod1 P37 0.23 P85 0.		Spec.Act1	P26	0.66	P74	0.46	P122	0.75
metabolome data of all samples Prod. P28 0.45 P76 0.65 P124 0.61 Spec.Prod1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P126 0.45 Act. P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 Phenotype at time point of sampling, Spec.Act2 † P33 0.03 P81 0.01 P129 0.05 metabolome data of mid log samples Prod. P34 0.21 P82 0.42 P130 0.18 Spec.Prod2 † P36 0.16 P83 0.24 P131 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P133 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P133 0.12 Spec.Art1 P37 0.23 P85 0.	Phenotype at time point of sampling,	Spec.Act2	P27	0.67	P75	0.44	P123	0.77
Spec.Prod1 P29 0.32 P77 0.49 P125 0.45 Spec.Prod2 P30 0.36 P78 0.48 P126 0.45 Act. P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 metabolome data of mid log samples Spec.Act2 † P33 0.03 P81 0.01 P129 0.03 Spec.Prod1 P34 0.21 P82 0.42 P130 0.18 Spec.Prod1 P35 0.16 P83 0.24 P132 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling, Spec.Act1 P38 0.25 P86 0.09 P134 0.29 Phenotype at time point of sampling, Spec.Prod2 † P40 0.49 P88 0.51 P136 0.51 Spec.Prod2 † P42 0.32 P90	metabolome data of all samples	Prod.	P28	0.45	P76	0.65	P124	0.61
Spec.Prod2 P30 0.36 P78 0.48 P126 0.43 Act. P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 metabolome data of mid log samples Spec.Act2 † P33 0.03 P81 0.01 P129 0.05 Spec.Prod1 P34 0.21 P82 0.42 P130 0.18 Spec.Prod1 P35 0.16 P83 0.24 P132 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P132 0.12 Act. P37 0.23 P85 0.29 P133 0.41 Spec.Prod2 † P38 0.25 P86 0.09 P135 0.29 metabolome data of late log samples Spec.Act2 † P39 0.25 P87 0.09 P135 0.29 metabolome data of late log samples Spec.Prod1 P41 0.32 P89		Spec.Prod1	P29	0.32	P77	0.49	P125	0.45
Act. P31 0.09 P79 0.01 P127 0.17 Spec.Act1 P32 0.03 P80 0.01 P128 0.05 metabolome data of mid log samples Spec.Act2 + P33 0.03 P81 0.01 P129 0.05 metabolome data of mid log samples Prod. P34 0.21 P82 0.42 P130 0.18 Spec.Prod1 P35 0.16 P84 0.24 P132 0.12 Spec.Prod2 + P36 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling, Spec.Act1 P37 0.23 P85 0.29 P134 0.29 Phenotype at time point of sampling, Spec.Act2 + P39 0.25 P87 0.09 P135 0.29 metabolome data of late log samples Prod. P40 0.49 P88 0.23 P137 0.26 Spec.Prod2 + P42 0.32 P90 0.23 P138 0.26 <tr< th=""><th></th><th>Spec.Prod2</th><th>P30</th><th>0.36</th><th>P78</th><th>0.48</th><th>P126</th><th>0.45</th></tr<>		Spec.Prod2	P30	0.36	P78	0.48	P126	0.45
Phenotype at time point of sampling, metabolome data of mid log samples Spec.Act2 + Prod. P33 0.03 P81 0.01 P128 0.03 Phenotype at time point of samples Prod. P34 0.21 P82 0.03 P81 0.01 P129 0.03 Spec.Act2 + P34 0.21 P82 0.16 P83 0.24 P130 0.12 Spec.Prod2 + P36 0.16 P84 0.24 P132 0.12 Spec.Prod2 + P36 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling, Spec.Act1 P37 0.23 P85 0.29 P134 0.29 Phenotype at time point of sampling, Spec.Act1 P38 0.25 P86 0.09 P134 0.29 metabolome data of late log samples Prod. P41 0.32 P89 0.23 P137 0.26 Spec.Prod1 P41 0.32 P89 0.23 P138 0.26 Spec.Prod2 <t< th=""><th></th><th>Act.</th><th>P31</th><th>0.09</th><th>P79</th><th>0.01</th><th>P127</th><th>0.17</th></t<>		Act.	P31	0.09	P79	0.01	P127	0.17
Phenotype at time point of sampling, metabolome data of mid log samples Spec.Act2 P33 0.03 P81 0.01 P129 0.03 Pred. P34 0.21 P82 0.42 P130 0.18 Spec.Prod1 P35 0.16 P84 0.24 P131 0.12 Spec.Prod2 † P36 0.16 P84 0.24 P132 0.12 Phenotype at time point of sampling, metabolome data of late log samples Act. P37 0.23 P85 0.29 P133 0.41 Spec.Act1 P38 0.25 P86 0.09 P134 0.29 Prod. P40 0.49 P88 0.51 P136 0.51 Spec.Prod1 P41 0.32 P89 0.23 P137 0.26 Spec.Prod2 † P42 0.32 P90 0.23 P138 0.26 Spec.Prod2 t P44 0.20 P92 0.14 P140 0.57 Phenotype at time point of sampling, Spec.Act1	Dhenetrine at time point of compling	Spec.Act1	P32	0.03	P60	0.01	P120	0.05
Interabolitine data of milding samples Prod.	metaboleme date of mid los comples	Spec.Act2 /	P33	0.05	P01 D02	0.01	P129 D120	0.05
Act. P33 0.10 P83 0.24 P131 0.12 Phenotype at time point of sampling, Act. P37 0.23 P85 0.29 P133 0.41 Spec.Act1 P38 0.25 P86 0.09 P134 0.29 Phenotype at time point of sampling, Spec.Act2 therapy P39 0.25 P87 0.09 P135 0.29 Prod. P40 0.49 P88 0.51 P136 0.51 Spec.Prod2 therapy P40 0.49 P88 0.51 P135 0.29 Prod. P40 0.49 P88 0.51 P136 0.51 Spec.Prod2 therapy P40 0.49 P88 0.51 P136 0.51 Spec.Prod2 therapy P40 0.49 P88 0.51 P137 0.26 Spec.Prod2 P42 0.32 P90 0.23 P137 0.26 Spec.Act1 P44 0.20 P92 0.14 P140	metabolome data of mid log samples	FIGU.	P 34	0.21	P02	0.42	P130	0.10
Act. P39 0.23 P85 0.29 P133 0.41 Phenotype at time point of sampling, Spec.Act1 P38 0.25 P86 0.09 P134 0.29 metabolome data of late log samples Spec.Act2 † P39 0.25 P87 0.09 P135 0.29 Spec.Prod1 P40 0.49 P88 0.51 P136 0.51 Spec.Prod2 † P40 0.49 P88 0.51 P135 0.29 Metabolome data of late log samples Prod. P40 0.49 P88 0.51 P135 0.21 Spec.Prod2 P41 0.32 P89 0.23 P137 0.26 Spec.Prod2 P42 0.32 P90 0.23 P138 0.26 Phenotype at time point of sampling, Spec.Act1 P44 0.20 P92 0.14 P140 0.57 Phenotype at time point of sampling, Spec.Act2 P45 0.19 P93 0.15 P141 0.59		Spec.Prod. 2 t	P35	0.10	P03	0.24	P131	0.12
Att. P37 0.23 P83 0.29 P133 0.41 Spec.Att1 P38 0.25 P86 0.09 P134 0.29 metabolome data of late log samples Spec.Att2 † P39 0.25 P87 0.09 P134 0.29 metabolome data of late log samples Prod. P40 0.49 P88 0.51 P136 0.29 Spec.Art2 † P40 0.49 P88 0.51 P136 0.51 Spec.Prod1 P41 0.32 P89 0.23 P137 0.26 Spec.Prod2 † P42 0.32 P90 0.23 P138 0.26 Phenotype at time point of sampling, Act. P43 0.18 P91 0.18 P139 0.69 Spec.Act1 P44 0.20 P92 0.14 P140 0.57 Phenotype at time point of sampling, Spec.Act2 P45 0.19 P93 0.15 P141 0.59 metabolome data of stationary samples <th></th> <th>Apt.</th> <th>F30</th> <th>0.10</th> <th>P04</th> <th>0.24</th> <th>P132</th> <th>0.12</th>		Apt.	F30	0.10	P04	0.24	P132	0.12
Phenotype at time point of sampling, metabolome data of late log samples Spec.Act2 † P39 0.25 P87 0.09 P135 0.29 Prod. P40 0.49 P88 0.51 P136 0.51 Spec.Prod1 P41 0.32 P89 0.23 P137 0.26 Spec.Prod2 † P42 0.32 P90 0.23 P138 0.26 Phenotype at time point of sampling, Act. P43 0.18 P91 0.18 P139 0.69 Phenotype at time point of sampling, Spec.Act1 P44 0.20 P92 0.14 P140 0.57 Phenotype at time point of sampling, Spec.Act2 P45 0.19 P93 0.15 P141 0.59 metabolome data of stationary samples Prod. P46 0.18 P94 0.39 P142 0.55 Spec.Prod2 P48 0.03 P944 0.34 0.38		Spec Act -1	P37	0.25	P85	0.29	P133 P134	0.41
metabolome data of late log sampling, metabolome data of late log samples Prod. P40 0.49 P88 0.51 P136 0.52 Prod. P40 0.49 P88 0.51 P136 0.51 Spec.Prod1 P41 0.32 P89 0.23 P137 0.26 Spec.Prod2 † P42 0.32 P90 0.23 P138 0.26 Phenotype at time point of sampling, metabolome data of stationary samples Act. P43 0.18 P91 0.18 P140 0.57 Phenotype at time point of sampling, metabolome data of stationary samples Spec.Act1 P44 0.20 P92 0.14 P140 0.57 Spec.Prod2 P45 0.19 P93 0.15 P141 0.59 Metabolome data of stationary samples Frod. P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42	Phenotype at time point of campling	Spec.Act1	P30	0.25	P 80	0.05	P134	0.20
Inclusion and of factory and period Instruction Instrution Instruction <th< th=""><th>metabolome data of late log samples</th><th>Prod</th><th>P40</th><th>0.25</th><th>P88</th><th>0.51</th><th>P136</th><th>0.51</th></th<>	metabolome data of late log samples	Prod	P40	0.25	P88	0.51	P136	0.51
Act. P43 0.18 P90 0.23 P138 0.26 Phenotype at time point of sampling, Act. P43 0.18 P91 0.18 P139 0.69 Spec.Act1 P44 0.20 P92 0.14 P140 0.57 Phenotype at time point of sampling, Spec.Act2 P45 0.19 P93 0.15 P141 0.59 Prod. P46 0.18 P94 0.39 P142 0.55 Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P94 0.45 P144 0.42	metabolome data of late log samples	Spec Prod -1	P41	0.45	P89	0.23	P137	0.26
Phenotype at time point of sampling, metabolome data of stationary samples Prod. P43 0.18 P91 0.18 P139 0.69 Phenotype at time point of sampling, metabolome data of stationary samples Spec.Act2 P45 0.19 P93 0.15 P140 0.57 Phenotype at time point of sampling, metabolome data of stationary samples Spec.Act2 P45 0.19 P93 0.15 P141 0.59 Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42		Spec Prod -2 †	P42	0.32	P90	0.23	P138	0.26
Phenotype at time point of sampling, metabolome data of stationary samples Prod. P44 0.10 P92 0.14 P140 0.57 Phenotype at time point of sampling, metabolome data of stationary samples Spec.Act2 P45 0.19 P93 0.15 P141 0.59 Spec.Prod1 P46 0.18 P94 0.39 P142 0.55 Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42		Act	P/13	0.32	PQ1	0.18	P139	0.69
Phenotype at time point of sampling, metabolome data of stationary samples Spec.Act2 P45 0.19 P93 0.15 P141 0.59 Prod. P46 0.18 P94 0.39 P142 0.55 Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42		Spec.Act1	P44	0.20	P92	0.14	P140	0.57
metabolome data of stationary samples Prod. P46 0.18 P94 0.39 P142 0.55 Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42	Phenotype at time point of sampling.	Spec.Act2	P45	0.19	P93	0.15	P141	0.59
Spec.Prod1 P47 0.05 P95 0.45 P143 0.38 Spec.Prod2 P48 0.03 P96 0.45 P144 0.42	metabolome data of stationary samples	Prod.	P46	0.18	P94	0.39	P142	0.55
Spec.Prod2 P48 0.03 P96 0.45 P144 0.42		Spec.Prod1	P47	0.05	P95	0.45	P143	0.38
		Spec.Prod2	P48	0.03	P96	0.45	P144	0.42

* For a detailed description of how each phenotype (P) was defined, see Supplementary data file 1. P is used to indicate models generated without LN transformation; LN(P) is used to indicate models generated after LN transformation of the phenotype; LN(M) is used to indicate models generated after LN transformation of the metabolome data.

 \pm For these PLS models, the results for Spec.Act.-2 and Spec.Prod.-2 are identical to Spec.Act.-1 and Spec.Prod.-1, respectively. To calculate Spec.Act.-2 and Spec.Prod.-2 in principal DWT_{max} is used, except for samples collected before DWT_{max} was reached (as is the case for the mid log and late log samples). For these samples DWT at the time point of sampling was used, similar as for calculating Spec.Act.-1 and Spec.Prod.-1 (see also Supplementary data file 1).

Information content of the metabolomics data set with respect to the different quantitative phenotypes

About 44% of the PLS models generated for glucoamylase were considered good models ($R_{CV}^2 \ge 0.6$); for protease, this was 19% (see Table 1). When comparing Tables 1A (glucoamylase) and 1B (protease) with each other, one thing is obvious: the highest information content of the metabolomics data set was obtained with different quantitative phenotypes for the different products. For glucoamylase good models were especially obtained when based on metabolome data of the samples from the mid log growth phase, while most good PLS models for protease were based on inclusion of metabolome data from all three time samples. Furthermore, LN transformation of either the metabolome data or the phenotype data resulted in general in an increased number of PLS models with $R^2_{CV} \ge 0.6$. In addition, more good PLS models were generated with the quantitative phenotype based on the *maximum* activity or productivity instead of the phenotype based at the activity or productivity at the time point of sampling. Moreover, for glucoamylase productivity resulted in more models with a R^{2}_{CV} above the cut-off of 0.6, while for protease on average the selection of activity (i.e., amount of product formed) as phenotype resulted in a somewhat higher number of good models.

Identification of metabolites that correlate with the phenotype studied

Metabolites contributing the most to, for instance, protease activity or productivity can be identified by ordering the (relative) statistical importance of the metabolites by virtue of the weight factors (regression factors) as determined in the PLS models for all metabolites. In other words, by applying PLS, metabolites important for a specific phenotype can be identified and ranked based on the strength of their correlation with the phenotype of interest. For both products, one good PLS model was chosen as starting point for analysing the strongest correlating metabolites in more detail. Based on this analysis subsequently lists of correlating metabolites from other good PLS models were compared.

Glucoamylase

For glucoamylase, most PLS models were above the threshold of $R_{CV}^2 = 0.6$ when using metabolome data of the samples collected during the mid log growth phase. From this group of models, the PLS model in relation to maximum activity (PLS model G7), was selected as starting point for target identification and comparison to other good PLS

models for glucoamylase. From this G7 PLS model, the 20 highest ranking metabolites are shown in Table 2. This top 20 included a relative high number of disaccharides and other sugar-derived compounds that were only present under glucoamylase inducing conditions (i.e. with glucose as carbon source). For all these dissacharides as well as some of the other compounds, such as DL-aminoadipic acid, 2,3-butanediol and xylitol the correlation is based on the absence of the compounds in all xylose samples and the presence in all glucose samples (Table 2). However, there is no clear correlation between their intracellular concentrations and maximum glucoamyase activity based on only the glucose samples (e.g. Fig. 3A). On the other hand, for putrescine, ornithine, glucose-6-phosphate, and fructose-6-phosphate there is a correlation between increasing intracellular levels of these compounds and maximum glucoamyase activity (e.g. see Fig. 3B).

When comparing the top 20's of other models with a $R^2_{CV} \ge 0.6$ with each other, especially the use of metabolome samples from particular sampling times was of influence on the resulting top 20 (see Supplementary data file 2A). When either the metabolome data of all time samples was used (e.g. model G49), or only the metabolome data of the mid log or late log samples (models G55 and G61, respectively), only four metabolites are present in all three resulting top 20's. These four metabolites were the compound tentatively identified as volemitol or perseitol, the compound tentatively identified as ribonic acid or xylonic acid, an unidentified disaccharide with a retention time of 42.02 min and another unidentified compound with ID AN 320-218 22.96 min (Supplementary data file 2A).



Fig. 3. Plot of the correlation between the metabolite tentatively identified as nigerose and maximum glucoamylase activity (A) and a similar plot for putrescine (B). O, Metabolome samples from xylose fermentations (n=11); **a**, metabolome samples from glucose fermentations (n=11).

rank	metabolite ID *	tentative identity	regression factor	visual correlation to phenotype †
1	dissacharide 39.13 min	nigerose	+	+ ‡
2	C5 sugar alcohol	xylitol	+	+ ‡
3	DL-aminoadipic acid		+	+ ‡
4	putrescine		+	+
5	disaccharide 319-361	kojibiose	+	+ ‡
6	ornithine		+	+
7	disaccharide 40.41 min	isomaltose	+	+ ‡
8	disaccharide 40.89 min	isomaltose §	+	+ ‡
9	xylose		-	- 11
10	histidine		+	0
11	glucose-6-phosphate		+	+
12	glucose		+	+ ‡
13	fructose-6-phosphate		+	+
14	AN 292-333 24.26 min	ribonic acid or xylonic acid	-	-
15	AN 201 26.51 min	unknown	+	+ ‡
16	spermidine		+	0
17	tryptophan		+	0
18	glutamine		+	+
19	2,3-butanediol		+	+ ‡
20	uric acid		+	+ ‡

Table 2. Twenty metabolites with the strongest correlation to glucoamylase as determined by PLS based on all mid log metabolome samples in relation to maximum activity (PLS model G7).

* All metabolites in this list were detected with the OS-GC-MS method.

†Visual correlation is indicated by + (positive correlation), – (negative correlation), or 0 (no apparent correlation); see also Supplementary data file 3A.
‡ Only or mainly high abundant on glucose, no apparent visual correlation within the glucose samples.

§ These are different mass fragments of the same compound.

| | Only high abundant on xylose.

The effect of LN transformation on the ranking of the potential targets was somewhat ambiguous. The effect of LN transformation of the phenotype or the metabolome data on the resulting top 20's was in several cases limited. For instance, for PLS models G7, G55 and G103 50% of the compounds were present in all three lists (for details, see Supplementary data file 2A). However, in other cases, i.e. PLS models G34, G82 and G130, this was only the case for 25% of the compounds (for details, see Supplementary data file 2A). The exact effect of LN transformation on the correlations of the metabolites with the phenotype was unclear; plotting of the peak areas of metabolites exclusively present in the top 20's after LN transformation against the

phenotype showed in some cases an improvement of the linear correlation, while in other cases the linear correlation deteriorated (data not shown).

Protease

For protease, most PLS models were above the threshold of $R^{2}_{CV} = 0.6$ when using the metabolome data of all three samples collected during the fermentation. The PLS model in relation to maximum activity, model P1, was selected from this group of models as starting point for target identification and comparison to other good PLS models for protease. From this PLS model, the 20 highest ranking metabolites are shown in Table 3. This top 20 mainly consisted of unidentified compounds detected by LC-MS, making interpretation of the results difficult. Two of the metabolites were tentatively identified as 2,3-dihydroxy-3-methylpentanoic acid and 2,3-dihydroxy-3methylbutanoic acid, both known intermediates of the isoleucine and valine biosynthesis, respectively. A number of the compounds in the top 20 contained a phosphate-group; however, very little is known of a possible involvement of phosphorus sources on protease expression in aspergilli. In comparison to the glucoamylase results, the relative high contribution of compounds analyzed with the RP-LC-MS method was remarkable. Among others, RP-LC-MS is suitable for the detection of aromatic peptides and peptides larger than 4-5 amino acids, suggesting that at least some of the high ranked compounds could be peptide-derived. Unfortunately, for none of these compounds appropriate reference compounds are currently available to establish their exact identity.

When comparing the top 20's from good PLS models for protease with each other, the overall observations are in line with those for glucoamylase. Also for protease the largest differences between the top 20's were observed when comparing models which were based on different selections of the metabolome data, e.g. metabolome data of all time samples or only the metabolome data of mid log or late log samples (see Supplementary data file 2B for details). Furthermore, the influence on the resulting top 20's was very limited when using either activity or specific activity as phenotype. This is to be expected, given the strong correlation between activity and specific activity, or productivity and specific productivity. On the other hand, the effect of LN transformation of either the phenotype or the metabolome data was considerable, as the resulting top 20's showed 50% or less overlap with the top 20 without LN transformation (Supplementary data file 2B).

rank	metabolite ID *	tentative identity	regression factor	visual correlation to phenotype †
1	428.0417 (RP)	unknown	+	+
2	AN 110-336 13.53 min (GC)	unknown	+	+
3	phosphorylethanolamine related (GC)	unknown	+	0
4	712.1019 (RP)	unknown	+	+
5	AN 312 15.42 min (GC)	unknown	+	+
6	2,3-dihydroxy-3- methylpentanoic acid (GC)		+	+
7	223.0937 (IP)	monomethylphosphate	+	0
8	2,3-dihydroxy-3- methylbutanoic acid (GC)		+	+
9	AN 298-342 (GC)	unknown	+	+
10	AN 342-299 31.30 min (GC)	unknown	-	-
11	AN 211-283 20.80 min (GC)	unknown	+	0
12	446.0929 (IP)	monomethylphosphate ‡	+	0
13	monomethylphosphate (GC)		+	0
14	230.1734 (RP)	unknown	+	0
15	171.0420 (RP)	unknown	+	+
16	207.0929 (IP)	monomethylphosphate ‡	+	0
17	799.1182 (IP)	unknown	+	0
18	688.1035 (RP)	unknown	+	0
19	428.0743 (RP)	unknown	-	0
20	Adenosine (GC)		+	0

Table 3. Twenty metabolites with the strongest correlation to protease as determined by PLS based on all metabolome samples in relation to maximum activity (PLS model P1).

* The analytical method used to detect each metabolite is indicated in between brackets: GC, OS-GC-MS; IP, IP-LC-MS; and RP, RP-LC-MS.

* Visual correlation is indicated by + (positive correlation), – (negative correlation), or 0 (no apparent correlation); see also Supplementary data file 3B.
* These are different mass fragments of the same compound.

DISCUSSION

The choice for a certain quantitative phenotype in bioprocess optimization often seems rather random, but may have a major influence on the outcome of an optimization strategy. In this study, the information content of a metabolomics data set was determined with respect to different quantitative phenotypes related to the formation of two simple products, i.e. glucoamylase, and a more complex product, i.e. protease. When comparing the results of the two enzyme products glucoamylase and protease, it could be concluded that the information content of the metabolomics data set is higher for the simpler of these two products, i.e. glucoamylase. This is on the one hand remarkable, because the fermentation conditions from which the metabolome samples were collected in this study, were originally selected to result in large and evenly distributed variation in protease activity (Braaksma *et al.*, 2009).

Another important aspect influencing the information content of the metabolomics data set is the time point at which metabolome samples were collected. For instance, in this study the information content of the metabolome data from the mid log time samples was high in respect to glucoamylase (Table 1A), while it was low for protease (Table 1B). Based on this result, we conclude that data sets based on fewer experimental conditions but more metabolome samples in time may be more informative than a data set based on many experimental conditions and only one or a few time samples per condition. In addition, data sets based on more samples in time will allow the analysis of longitudinal effects in the data, i.e. metabolites whose correlation with product formation show a shift in time (Rubingh *et al.*, 2009).

Our results show that the effect of different ways to calculate the quantitative phenotype on the information content and resulting targets is much smaller than the effect of the time point of sampling. In general, the number of PLS models with a R^2_{CV} above the threshold value was higher when quantitative phenotypes were used that were based on the maximum activity or productivity instead of the activity or productivity at the time point of sampling (Table 1). A possible explanation for this is the more distinct variation in phenotypic values for the maximum phenotype. This may correlate better to the variation in the metabolome data present at a time point when phenotypic differences are perhaps not yet that clearly visible. Nevertheless, the effect of either maximum phenotype or phenotype at time point of sampling on the resulting top 20's is limited (Supplementary data file 2). This holds for the different description of the phenotype (e.g. activity versus productivity, or activity versus specific activity) as well. Conversely, the effect of LN transformation was considerable. Not only did the number of PLS models with a R²_{CV} above the threshold value increase with LN transformation of the phenotype or the metabolome data, the resulting top 20's were often considerably different from the top 20 based on the data without LN transformation. However, it should be noted that it is difficult to interpret the effect of LN transformation, especially as it is not clear how LN transformation and data pretreatment methods (e.g. scaling methods such as range scaling) influence each other with regard to complex metabolome data (van den Berg *et al.*, 2006).

With the MVDA tool PLS the quantifiable phenotype of interest can be related to the metabolome data set as a whole and at the same time take into account the

relationship between the metabolites (van der Werf *et al.*, 2007). Without this, it would be necessary to plot the metabolite concentrations of each metabolite against the phenotype in order to investigating the relation between individual metabolites and the quantifiable phenotype of interest. However, in case of a large number of metabolites, and as in our case a large number of phenotypes as well, this approach will result in an extremely large number of plots to analyze. Moreover, in such plots the intrinsic interdependency of the metabolites is neglected. However, despite these advantages of MVDA over a univariate approach, interpretation of the relation of the metabolites ranked by PLS to the quantifiable phenotype of interest is not straightforward. Several aspects, as listed below, have to be taken into account when interpreting the results of a PLS model.

(1) The positive or negative regression factors that are a measure for the contribution of a metabolite to the phenotype cannot be directly translated into how a metabolite actually correlates to the phenotype. These regression factors are not only a measure for the correlation of a single metabolite to the phenotype, but also for the correlation of this metabolite to other metabolites. Therefore, for a more detailed biological interpretation it is recommended to plot the concentrations of highly correlated metabolites against the quantifiable phenotype.

(2) Not all metabolites found to be correlating to the phenotype of interest are involved in the production of this product, either as inducer/inhibitor or precursor/side-product. With MVDA no distinction can be made between metabolites that correlate to the phenotype due to either a cause or an effect relation. For instance, one may conclude that the disaccharides found to be correlating to high glucoamylase activity (Table 2) induce glucoamylase secretion and thus cause the high activities. However, it is also possible that the identified disaccharides were formed from glucose by transglucosylation activity from glucoamylase (Nikolov *et al.*, 1989), and thus are an effect of glucoamylase activity ('effect correlation'). For strain improvement in particular cause relations are of importance.

(3) Related to the previous subject is the occurrence of confounding effects, i.e. the situation that an extraneous factor correlates with both the phenotype and a metabolite. This can result in the false conclusion that there is a causal relationship between the phenotype and that specific metabolite. For example, there is only significant glucoamylase activity when *A. niger* is cultured on glucose instead of on xylose. Also several metabolites, such as uric acid and xylitol, are mainly present when *A. niger* is cultured on glucose. Therefore, one may conclude that there is a direct correlation between these metabolites and glucoamylase production. However, these

compounds may not be directly linked to glucoamylase production per se, but perhaps both glucoamylase and these metabolites independently correlate to growth on a specific carbon source.

(4) With the comprehensive analytical methods used in this study not only known compounds are analyzed, but also all peaks of unknown identity are included in the data set. One last aspect hampering the interpretation of the results of the data analysis is the correlation of these unidentified metabolites with the phenotype.

Taking into account the various aspects that influence the interpretation of the PLS results, as discussed above, specific metabolites identified as important to the question under study can be distilled from the initial list of potential targets that result from PLS. For optimization of glucoamylase production glucose-6-phosphate and fructose-6-phosphate are among the most likely targets. The enzyme glucose-6phosphate isomerase catalyzes the conversion of glucose-6-phosphate into fructose-6phosphate. The ratio between the concentrations of glucose-6-phosphate and fructose-6-phosphate is approximately a factor seven higher than expected based on the equilibrium constant for glucose-6-phosphate isomerase (data not shown). The relative accumulation of glucose-6-phosphate may on the one hand suggest that the activity of this enzyme is a bottleneck in the flux through the glycolysis. On the other hand, this aberration of the equilibrium may be required to obtain a sufficient flux in the direction of the pentose phosphate pathway (PPP), in order to generate sufficient NADPH. Melzer et al. (2007) also observed that under glucoamylase-producing conditions the flux of glucose through the PPP was higher than for non-producing conditions. However, in our study even under non-producing conditions the ratio between glucose-6-phosphate and fructose-6-phosphate concentrations is approximately a factor seven higher than expected. This weakens the hypothesis that the flux through the PPP may only be insufficient under glucoamylase-producing conditions, although when glucose was used as carbon source the absolute concentrations of both metabolites are higher. Alternatively, also absolute metabolite concentrations could be involved in regulation of metabolite fluxes (e.g. allosteric effects). All in all, in view of the crucial position of glucose-6-phosphate isomerase at the branch point between the glycolysis and the PPP, the regulation of the activity of this enzyme may be a means to regulate the fluxes through these two pathways and thus optimize glucoamylase production. Putrescine and ornithine are the two other most likely targets for optimization of glucoamylase production. Ornithine is the starting point for the synthesis of polyamines such as putrescine. Little is known about the actual function of putrescine and other polyamides in *A. niaer*. In *A. nidulans*, there is an absolute requirement of polyamides in growth and development (Tabor &

Tabor, 1985; Jin *et al.*, 2002). The positive correlation between glucoamylase production and putrescine suggests that glucoamylase production may be stimulated by either addition of this polyamine to the medium or overexpression of the gene encoding ornithine decarboxylase, the enzyme responsible for the conversion of ornithine into putrescine.

No obvious targets were found in relation to protease production. Moreover, the majority of the compounds correlating to protease activity are unidentified compounds (Table 3). The presence of several compounds analyzed with the RP-LC-MS method in Table 3 suggests the possible involvement of small peptides in protease induction. Unfortunately, identification of peptides with the RP-LC-MS method has proved to be quite difficult, also because of the lack of appropriate reference compounds. Therefore, in order to further investigate the possible role of peptides in protease induction, additional methods will have to be deployed that offer more detailed information on the (partial) identity of peptides.

It was anticipated that the relation between intracellular metabolite concentrations and extracellular protease activity would not be straightforward, because extracellular protease activity is a complex phenotype, consisting of multiple enzyme activities. Recent analysis of the secretome of *A. niger* has indicated the presence of up to 20 different secreted proteases in the medium (Tsang *et al.*, 2009; Braaksma *et al.*, 2010b). Possibly, an approach with metabolomics alone is not sufficient for identifying targets for such a complex phenotype and an integrated systems biology approach is required.

Besides glucoamylase and protease production, also citric acid production was analysed as a phenotype. Although the experimental design of our data set was not optimally suited for this product, resulting in very few reliable PLS models (Supplementary data file 4), several TCA cycle intermediates (isocitrate, α -ketoglutarate) were identified as correlating with citric acid production (results not shown). Altogether, this study illustrates that with a combined metabolomics/MVDA approach relevant targets for strain and process improvement can be identified, as the relevance of several of the identified leads seem confirmed by what already is known in literature (e.g. the role of glucose-6-phosphate isomerase in glucoamylase production). Moreover, this study demonstrates the importance of experimental design in top-down systems biology studies, not only with regard to the fermentation conditions, but also with respect to the time point of sampling and the selection and calculation of the quantitative phenotype to be pursued.

ACKNOWLEDGEMENTS

The authors thank Karin Overkamp for her helpful input on the sample work-up, Maud Koek and Maarten Hekman for the GC-MS and LC-MS analysis, Bas Muilwijk and Richard Bas for identification of the metabolites, and Carina de Jong-Rubingh for the smoothing algorithm analyses and calculating the quantitative phenotypes. This project was carried out within the research programme of the Kluyver Centre for Genomics of Industrial Fermentation, which is part of the Netherlands Genomics Initiative / Netherlands Organization for Scientific Research.