**PhylochipAnalyzer - A Program for Analysing Hierarchical Probe-Sets**

Katja Metfies[*], Philipp Borsutzki[‡], Christine Gescher[*], Linda K. Medlin[*],

Stephan Frickenhaus[*,+]


[*]Alfred-Wegener-Institute for Polar and Marine Research, Am Handelshafen 12,

27570 Bremerhaven, Germany

[‡] Urnenfelderstraße 9, 85051 Ingolstadt, Germany

[+]Corresponding author:

Stephan Frickenhaus, Alfred-Wegener-Institute for Polar and Marine Research,

Am Handelshafen 12, 27570 Bremerhaven

Tel. ++49471-4831 1179, Fax. ++49471-4831 1590, e-mail: Stephan.Frickenhaus@awi.de

1 **Abstract**

2 The recent introduction of phylochips that contain molecular probes facilitates environmental

3 microbial identification in a single experiment without previous cultivation. A set of probes

4 recognizing species at different taxonomic levels is denoted as a hierarchical set. Application

5 of hierarchical probe sets on a DNA-microarray allows the assessment of biodiversity with

6 different resolutions. It significantly increases the robustness of the results retrieved from

7 phylochip experiments because of the possible consistency checks of hybridization across

8 different taxonomic levels. Here, we present a computer program, PhylochipAnalyzer, for the

9 hierarchy editing and the evaluation of phylochip data generated from hierarchical probe-sets.

1   **Basic rationale:** Recently, more and more publications describe the application of DNA

2   microarrays for species identification (phylochips) from environmental samples (Guschin et

3   al. 1997; Loy et al. 2002; Metfies and Medlin 2004; Call 2005, Medlin et al. 2006).

4   Phylochips are DNA-microarrays containing molecular probes that bind to unique sequences

5   in a target. The target sequence is usually part of marker genes, e.g., the ribosomal RNA gene.

6   Ribosomal RNA-genes are particularly well suited for phylochip- and phylogenetic analysis,

7   because they are universal, found in all cellular organisms, are of relatively large size; and

8   contain both highly conserved and variable regions with no evidence for lateral gene transfer

9   (Woese 1987). The large number of published 18S rDNA-sequences, (e.g., RDP, Maidak et

10   al. 2001) makes it possible to design hierarchical probe sets that specifically target the 18S-

11   rDNA from higher taxa down to species level (Lange et al. 1996; Guillou et al. 1999; Groben

12   et al. 2004). Phylochips provide a promising tool to identify large numbers of microbial

13   species in complex environmental samples quickly without a cultivation step. Our phylochip

14   contains a hierarchical set of probes, which target phytoplankton species at different

15   taxonomic levels (Metfies and Medlin 2004; Medlin et al. 2006).  In a hierarchical probe-set,

16   a target species is only considered present, if all hierarchical probes for each species result in

17   a positive signal. Therefore, hierarchical probes add to the accuracy of molecular probe based

18   identification approaches.

19       In spite of the growing number of applications for phylochips, they represent only a

20   small proportion of all DNA-microarray related work. Most publications describe expression

21   studies (e.g., Lockhart et al. 2000; Stoughton 2005; Rensink 2005; Csako 2006).

22   Consequently, the majority of protocols are optimized for applications related to expression

23   analysis. However, the application of phylochips for species identification in environmental

24   samples presents technical challenges that are not encountered in gene expression studies of

25   laboratory samples (Peplies et al. 2003; Call et al. 2005; Metfies et al. 2006). There are

1 numerous commercial and non-commercial programs for the analysis of expression studies

2 (e.g. Dondrup et al. 2003; Vaquerizas et al. 2005) but few programs exist for phylochip

3 analysis. One example is the Unix-based program ChipChecker (Loy et al. 2002), which is

4 dedicated to data interpretation from phylochips. It calibrates signal to noise ratios to a set

5 threshold determined by the user and finds positive signals with respect to that threshold

6 based on the fact that a positive signal can only be located where there is a fully

7 complementary probe to its target. However, in a hierarchical probe set, a signal is only

8 considered truly positive, if all probes in the hierarchy are positive. Therefore, the analysis of

9 hierarchically organized phylochips requires an additional step in comparison to the functions

10 provided by ChipChecker. The positive signals must be tested for their robustness in relation

11 to the hierarchy on the phylochip. In summary, a program for the analysis of hierarchically

12 organised phylochips has to provide an algorithm for the calculation of a signal to noise-value

13 and a tool that allows to set positive signals in relation to the hierarchy inherent in the design

14 of the probe-set. Here we present the program, PhylochipAnalyzer, that implements the

15 calculation of signal to noise ratios and the evaluation of phylochip-data with respect to probe

16 hierarchy.

17
18 **Funcionality and Implementation aspects of the Program**
19
20 PhylochipAnalyzer is a GUI-based Windows-program, developed under Borland-Delphi. The

21 program combines two strongly interconnected functions: hierarchy editing and data analysis.

22 The user starts editing interactively and graphically the hierarchy that is inherent in the

23 chip/probe design process. Editing is started by loading a spot description file in GAL-format

24 generated by the GenePix- software (Axon Instruments Inc., USA). A procedure to convert

25 other formats is described in the software documentation. Spot entries are shifted manually so

26 that a hierarchically structured tree-like layout appears, in correspondence to the hierarchical

27 probe design of the chip seen in Fig. 2A, upper part. Probes must not be placed  in a hierarchy

1    at all, e.g., positive or negative controls should be placed as stand-alone, i.e., with no parent

2    probes and no child probes. However, a positive control could be placed as the parent probe to

3    all other. The hierarchy is then saved as an XML-file that is used later for data analysis.

4    Whereas the XML-file stores the pure hierarchy information of the chip, spot-intensity data

5    are read from files with externally defined format, such as tab-delimited tables. The user may

6    include the probe sequence in the comment field. The hierarchy can be exported as a tree file

7    in Newick-format.

8        The second mode of operation is for the analysis of processed scanner data, i.e., tables

9    with data for foreground and background intensities of the individual spots. The presence or

10    absence of a hybridization signal is checked by a threshold criterion. The foreground-

11    background intensity contrasts are normalized with respect to intensities of the negative

12    control spots (Loy et al. 2002). Here intensity data of multiple copies (blocks) of the spots on

13    each chip are evaluated and means and standard deviations are computed. The results for the

14    blocks on the chip are shown independently (Fig. 2A, bottom right) such that entire blocks

15    can be excluded from the analysis. It is assumed that if some spots in a certain block are

16    identified as outliers or if positive controls fail, the user should exclude the whole block from

17    evaluation because of the questionable quality of hybridization. A false positive signal on a

18    higher hierarchical level has consequences for the validity of lower levels, down to the

19    species level: PhylochipAnalyzer marks all positive signals that are below the hierarchy level

20    of a spot showing a negative signal, i.e., corrected lines are crossed out. Because a signal is

21    marked positive when the majority of copies give signals above the threshold, a correction is

22    always contradicting. The user should inspect whether the underlying probe is correctly

23    designed or maybe placed in the wrong hierarchy level.

24    The user may export the evaluation results directly to an Excel-graph (Fig.2B) in which the

25    signals are given as bars, labelled with the probe identifier. The size of a bar indicates the

1  quality above the threshold, i.e., the longer the bar, the stronger the evidence for a positive

2  signal. All data are shown with error bars of the mean due to the variance over the different

3  blocks.

4  **Validation**

5  The PhylochipAnalyzer was used to analyse data retrieved from a hybridization of PCR-

6  products of *Micromonas pusilla* 18S rDNA to a phylochip that contained 44 probes, including

7  a hierarchical probe-set for the Prasinophyte genus *Micromonas*. The hierarchical probe-set

8  consisted of six probes that bind, respectively, at the level of Kingdom (EUK 1209, EUK

9  328), Class (Chlo01, Chlo02), Clade or Order (Pras 04) and Genus (Micro01) to *Micromonas*

10  *pusilla*. The additional probes on the chip identified other phytoplankton taxa, a negative

11  control, and two positive controls. Fluorescence images of the hybridized phylochips were

12  taken with the Genepix 4000B Scanner (Axon Instruments Inc. USA). The signal intensities

13  were quantified using the GenePix 6.0 software (Axon Instruments Inc. USA). Raw data were

14  saved as a GPR-file and imported to the PhylochipAnalyzer-program. The computation of the

15  raw data with the PhylochipAnalyzer-program identified only positive signals for the

16  perfectly matching probes. For those probes, a signal/noise ratio was calculated that was

17  above the threshold. The complete hierarchical probe set resulted in positive signals, therefore

18  the signal for Micro01 can be considered truly positive (see Fig. 1 and Fig. 2B).

19  **Discussion**

20  The program simplifies tremendously the time consuming tasks of data processing of results

21  from hierarchical phylochips. This is from particular interest, if high-throughput data are

22  analyzed. The program is flexible with respect to configuration because the user can influence

23  the threshold criterion by modifying the code that is implemented as a Delphi-script. This

24  allows arbitrary modifications of the basic formula of data processing. Other formats of

25  intensity description can easily be converted into appropriate GAL-format. On screen, the

1    user may change the threshold value (default 2) interactively for sensitivity studies and

2    recalculation. The rather simple criterion for elimination of false positives could be extended

3    towards more quantitative measures. We plan to extend the program for quantitative analysis,

4    i.e., spots from higher hierarchical levels are expected to show stronger signals than the lower

5    hierarchical spots because they target more individuals. Multi-chip comparative analysis (e.g.,

6    clustering) for time-series analysis is also a desirable feature. The proposed XML-format for

7    hierarchy representation can be seen as a prototype for standardization in phylochip hierarchy

8    description. It is now necessary to introduce community standards for the representation of

9    both, chip description and data-processing details. For gene-expression analysis by means of

10   DNA-microarrays guidelines already exist (Brazma et al. 2001). Standards for phylochip

11   design and processing description are considered to be a prerequisite for permanent archiving

12   of publication supplemental data accompanied by catalogues of metadata in repositories.

18

1  **References**

2  Brazma A, Hingamp P, Quackenbush J, Sherlock G, Spellman P, Stoeckert C, Aach J,

3  Ansorge W, Ball WA, Causton HC, Gaasterland T, Glenisson P, Holstege FCP, Kim IF,

4  Markowitz V, Matese JC, Parkinson H, Robinson A, Sarkans U, Schulze-Kremer S, Stewart

5  J, Taylor R, Vilo J, Vingron M (2001) Minimum information about a microarray experiment

6  (MIAME)—toward standards for microarray data. Nature Genetics 29, 365 - 371.

7

8  Call DR (2005) Challenges and opportunities for pathogen detection using DNA microarrays.

9  Critical Reviews in Microbiology 31, 91-99.

10

11  Csako G (2006) Present and future of rapid and/or high-throughut methods for nulcleic acid

12  testing. Clinica Chimica acta 363, 6-31.

13

14  Dondrup M, Goesmann A, Bartels D, Kalinowski J, Krause L, Linke B, Rupp O, Sczyrba A,

15  Puhler A, Meyer F (2003) EMMA: a platform for consistent storage and efficient analysis of

16  microarray data. J Biotechnol 106:135-46.

17

18  Groben R, John U, Eller G, Lange M. and Medlin, L.K. (2004) Using fluorescently- labelled

19  rRNA probes for hierarchical estimation of phytoplankton diversity. Nova Hedwigia. 79, 313-

20  320.

21

22  Guschin DY, Mobarry BK, Proudnikov D, Stahl DA, Rittmann BE, Mirzabekov AD (1997)

23  Oligonucleotide microchips as genosensors for determinative and environmental studies in

24  microbiology. Appl. Environm. Microbiol. 63, 2397-2402.

25

1  Guillou L, Moon-van-der-Staay SY, Claustre H, Partensky F, Vaulot D (1999) Diversity and

2  abundance of Bolidophyceae (Heterokonta) in two oceanic regions. Appl. Environ. Microbiol.

3  65, 4528-4536.

4

5  Lange M, Guillou L, Vaulot D, Simon N, Amann RI, Ludwig W, Medlin LK (1996)

6  Identification of the class Prymnesiophyceae and the genus *Phaeocystis* with ribosomal RNA-

7  target nucleic acid probes detected by flow cytometry.  J. Phycol. 32, 858-868

8

9  Lockhart DJ, Winzeler EA (2000) Genomics, gene expression and DNA arrays. Nature 405,

10  827-836.

11

12  Loy A, Lehner A, Lee N, Adamczyk J, Meier H, Ernst J, Schleifer KH, Wagner M (2002).

13  Oligonucleotide Microarray for 16S rRNA Gene-Based Detection of All Recognized

14  Lineages of Sulfate-Reducing Prokaryotes in the Environment., Appl. Environm. Microbiol.
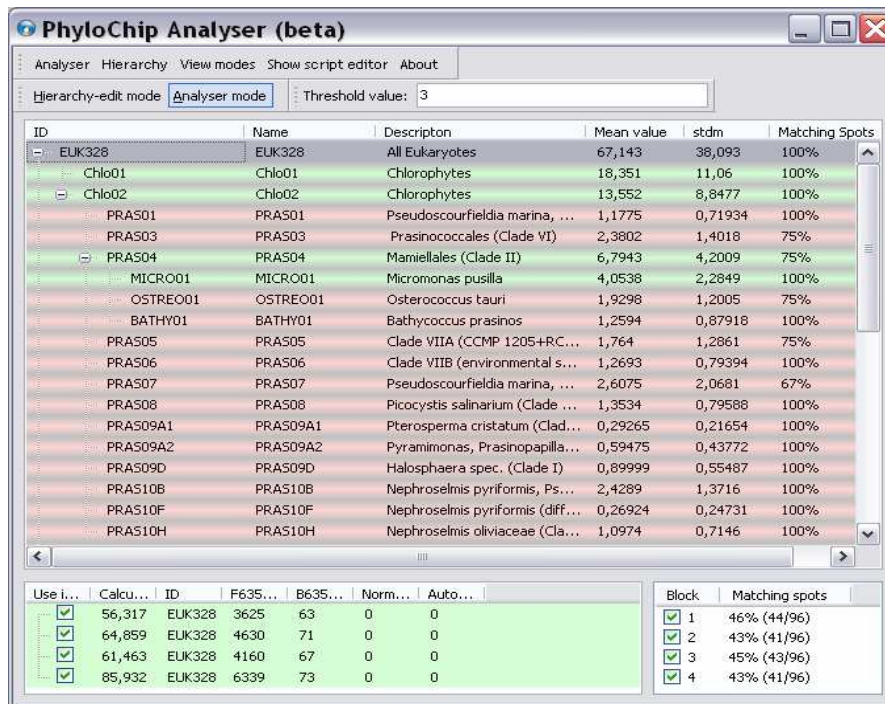
15  68, 5064-5081.

16

17  Maidak BL, Cole JR, Lilburn T G, Parker CTJ, Saxman PR, Farris RJ, Garrity, GM, Olson

18  GJ, Schmidt TM, Tiedje JM (2001) The RDP-II (Ribosomal Database Project). Nucleic.

19  Acids Res. 29, 173-174.

20

21  Medlin LK, Metfies K, Mehl H, Wiltshire K, Valentin K (2006) Picoeukaryotic plankton

22  diversity at the Helgoland time series site as assessed by three molecular methods. Microbial

23  Ecology, 52, 53-71.

24

1    Metfies K, Medlin L (2004). DNA Microchips for Phytoplankton: The Fluorescent Wave of

2    the Future. Nova Hedwigia, 79, 321-327.

3

4    Peplies J, Glöckner FO, Amann R (2003) Optimization strategies for DNA microarray-based

5    detection of bacteria with 16S rRNA-targeting oligonucleotide probes. Appl. Environ.

6    Microbiol., 69(3), 1397-407.

7

8    Stoughton RB (2005). Applications of DNA Microarrays in Biology. Annual Review of

9    Biochemistry, 74, 53-82.

10

11   Rensink WA, Buell CR (2005) Microarray expression profiling resources for plant genomics.

12   Trends in Plant Sciences 10, 603-609.

13

14   Vaquerizas JM, Conde L, Yankilevich P, Cabezon A, Minguez P, Diaz-Ulriarte R, Al-

15   Shhrour F, Herrero J, Dopazo J (2005) GEPAS, an experiment-oriented pipeline for the

16   analysis of microarray gene expression data. Nucleic Acids Res. 33(Web Server issue):W616-

17   20.

18

19   Woese C.R. (1987). Bacterial evolution. Microbiological Reviews. **51**,  221-271.
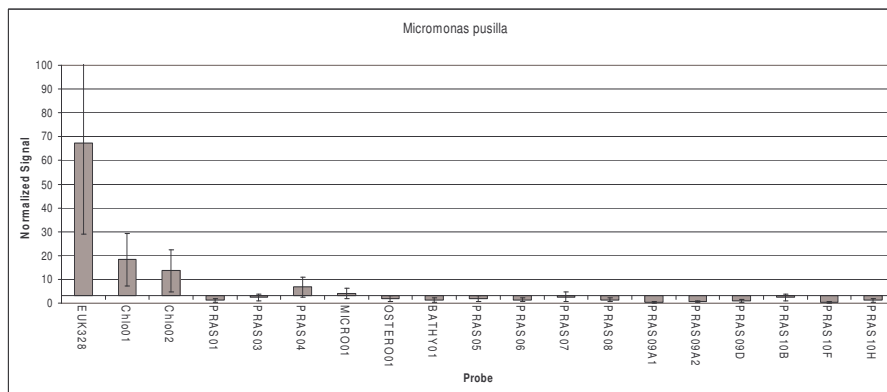
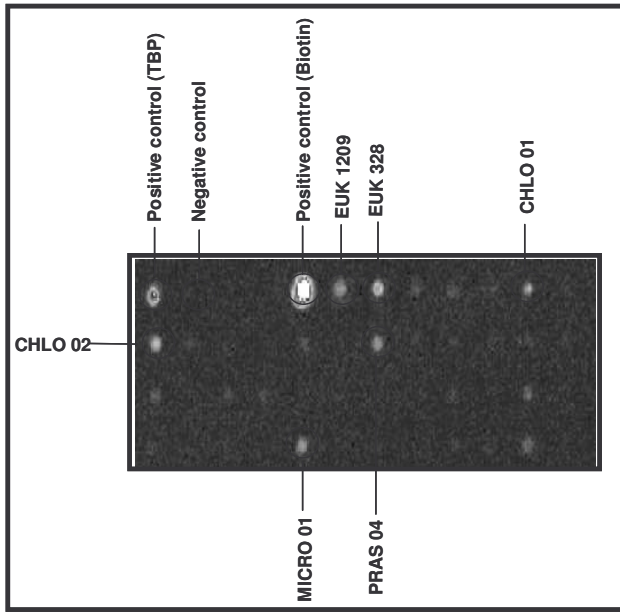20

1    **Figure 2:**

2    **A**



3

4    **B**



5

**Figure 2. A:** Screenshot of the analyser-mode. Any set of molecular probes can be organized as a user defined phylogenetic tree by a drag and drop function in editor-mode. The screenshot displays a tree of probes that bind to Prasinophytes at different hierarchical levels. The bottom part shows an individual probe result for the selected probe (EUK328, top part). **B**: Output of signal-noise values in graphical form.

1 **Figure 1:**



**Figure 1:** The 18S rDNA of *Micromonas pusilla* was hybridized to a set of 44 probes. The set of probes contained a hierarchical set that binds to the 18S rDNA of *M. pusilla* at four different taxonomic levels (EUK 1209, EUK 328, Chlo01, Chlo02, Pras04 and Micro01).