# The PANGAEA® Data Warehouse

U. Schindler (1), M. Diepenbroek (1), H. Grobe (2), and R. Sieger (2)

(1) MARUM - University of Bremen, PANGAEA, Bremen, Germany (uschindler@pangaea.de, mdiepenbroek@pangaea.de),
(2) Foundation Alfred Wegener Institute for Polar and Marine Research (AWI), PANGAEA, Bremerhaven, Germany
(hannes.grobe@awi.de, rainer.sieger@awi.de)

PANGAEA® - Publishing Network for Geoscientific & Environmental Data (www.pangaea.de) currently provides published data entities for download that can be cited like publications using authors, year, title and a Digital Object Identifier (DOI). A lot of scientists, especially modellers, need compilations of various data sets for analyzing. On the other hand, data producers want to be cited for their work, which is nearly impossible with huge compilations containing thousands of distinct data sets, especially when data compilation is done manually outside of the PANGAEA data library.

We will present our recently introduced, AJAX-based web interface based on data warehouse technologies that can used for highly efficient retrievals and compilations of time slices or surface data matrixes on any measurement parameters out of the whole data continuum. The user is able to first select the original datasets using the current PANGAEA full text search engine (based on ISO 19115 metadata). After that he can select the configuration of measurement parameters and methods for the data compilation using an innovative scoring algorithm based on the original search query. The data matrix can be downloaded along with an ISO 19115 compatible metadata description about the compilation, referencing all original datasets. Additionally, each data point can be traced back to the original, citable data set.

# The PANGAEA® Data Warehouse

Uwe Schindler[1], Michael Diepenbroek[1], Hannes Grobe[2], Rainer Sieger[2]

## Abstract

PANGAEA® - Publishing Network for Geoscientific & Environmental Data (www.pangaea.de) currently provides published data entities for download that can be cited like publications using authors, year, title and a Digital Object Identifier (DOI). A lot of scientists, especially modellers, need compilations of various data sets for analyzing. On the other hand, data producers want to be cited for their work, which is nearly impossible with huge compilations containing thousands of distinct data sets, especially when data compilation is done manually outside of the PANGAEA data library.

We present our recently introduced, AJAX-based web interface based on data warehouse technologies that can used for highly efficient retrievals and compilations of time slices or surface data matrixes on any measurement parameters out of the whole data continuum. The user is able to first select the original datasets using the current PANGAEA full text search engine (based on ISO 19115 metadata). After that he can select the configuration of measurement parameters and methods for the data compilation using an innovative scoring algorithm based on the original search query. The data matrix can be downloaded along with an ISO 19115 compatible metadata description about the compilation, referencing all original datasets. Additionally, each data point can be traced back to the original, citable data set.

## About PANGAEA®

PANGAEA® (www.pangaea.de) is operated by the Alfred Wegener Institute for Polar and Marine Research (AWI), Bremerhaven, member of the Helmholtz Association of National Research Centres, funded by the Federal Ministry of Education and Research and the Center for Marine Environmental Sciences (marum) at the University of Bremen with support of the Department of Geosciences. Around 20 people are associated with WDC-MARE. The budget amounts approximately 1,2 Mio Euro per year for personnel, hard-, and software. Third party funds are about 70% of the total budget.The basic technical structure corresponds to a three tiered client/server architecture with a number of clients and middleware components controlling the information flow and quality. On the server side a relational database management system (RDBMS) is used for information storage. To ensure fast data access the data are mirrored in a data warehouse (SYBASE IQ) which is also used as interface to the german GRID community. All interfaces to the information system's metadata, datasets and data warehouse are based on web services including a simple map supported search engine.

**Database status (2/2009):**

578 931 data sets
4 669 156 961 data items

With its comprehensive user interfaces and the built in editorial system for import, and maintainance of information PANGAEA® is a highly efficient system for scientific data management and data publication.
The challenge of managing the heterogenous and dynamic data of environmental and geosciences was met in PANGAEA® through a flexible data model which reflects the information processing steps in the earth science fields and can handle any related analytical data.
Metadata are consistently stored in a relational database and can be served as ISO19115, DIF, FGDC, or DC compliant catalogues.
The data description of any data set includes the principle investigators (PI) name and email for contact, parameters, used methods as well as spatial and temporal coverage.
Except a few data from ongoing projects which are under moratorium all data are freely available and can be used by referencing the related publication or the data set citation.

**Data for nearly 50 000 different measurement types, e.g.:**
Sediment and Ice Profiles, Seismic Profiles, Mineral distribution, Geological Maps, Atmospheric Profiles, Radiation Data, Distributed Samples, Oceanographic Profiles, Time Series, Plankton and Fish
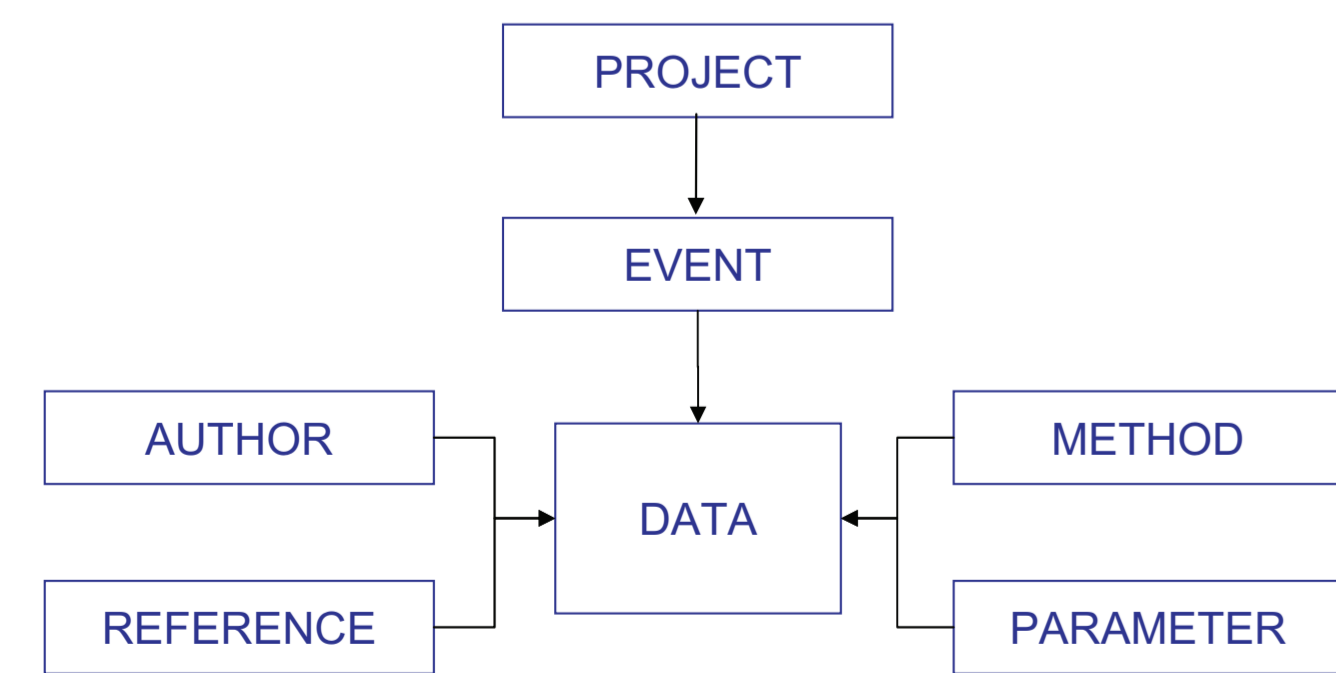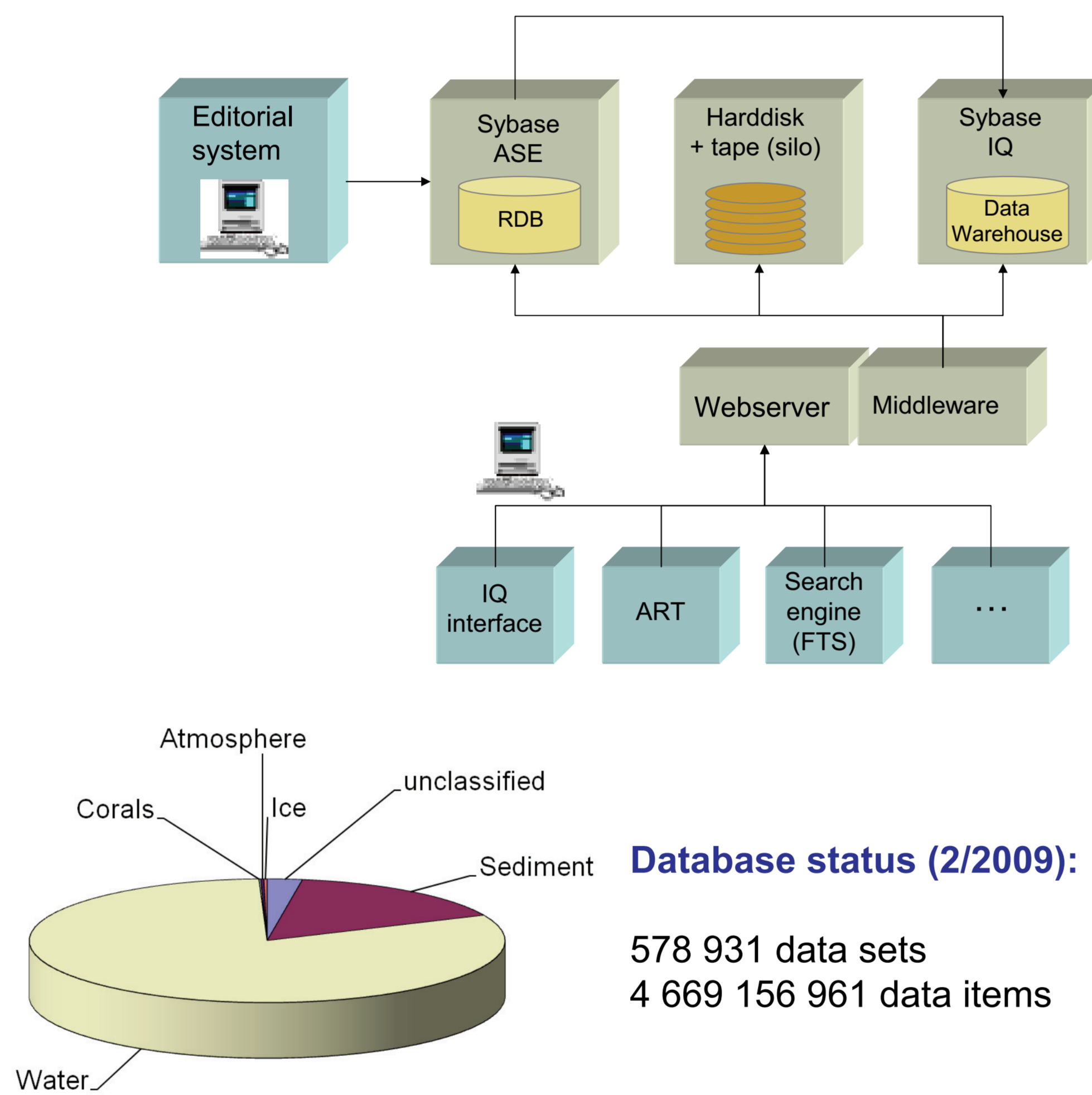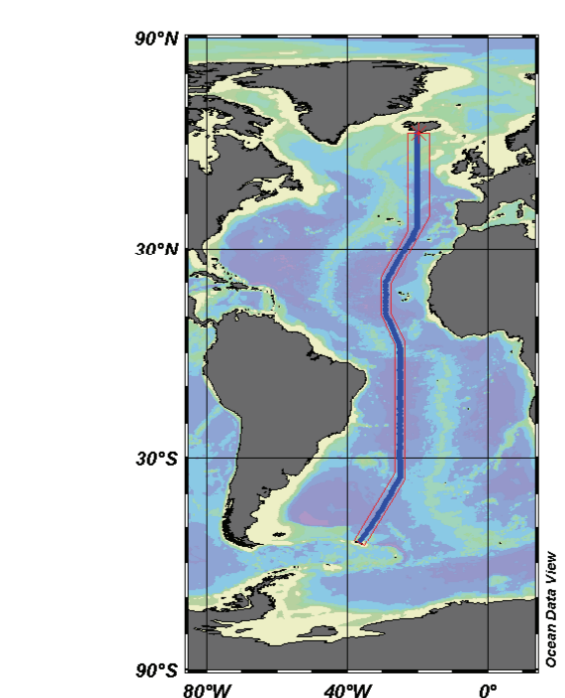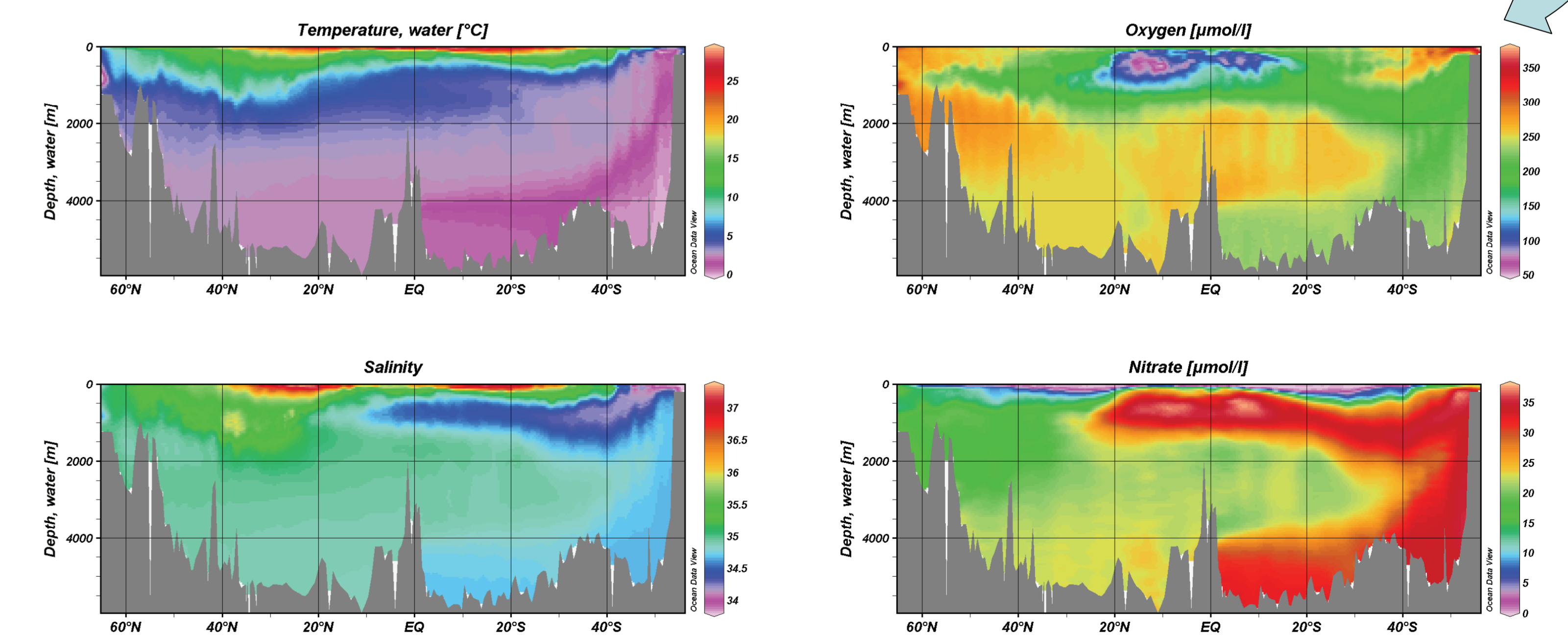
## Sample Workflow: WOCE A16 Profile

To start a data warehouse query, you must initially do a conventional full-text search engine query on the PANGAEA dataset metadata. Just enter some specific keywords for the data you are searching. In the example, we searched for data of the WOCE A16 profile. Optionally you can also add other constraints like a specific bounding box from where the data should come from or a date range for time series.
You get a list of citable datasets entities. In the past, you had to download all these single datasets as separate files and combine them using separate tools.
With the new data warehouse interface of PANGAEA, you can download the data as a single file containing only the relevant data points. From the result list, you can click on "Data Warehouse (BETA)" (currently only available for logged-in users) to open the new interface.

The result set of datasets is analyzed for measurement parameters, methods and used geocodes (corrdinates like latitude, longitude, date, time, elevation) and all of these are displayed in a scored list. An AJAX based GUI helps you to choose the interesting parameters and combine them in a "configuration" sheet for the resulting data file. Parameters and geocodes can be combined as columns in the resulting data matrix, methods can be choosen together with aggregation functions.
With one further click you get the tab-delimited result matrix as a single file.
As a lot of different datasets were involved when creating the result, each data point is referenced back to the citable dataset, the measurement value was taken from, helping to correctly reference the original data.
The result file can be visualized with common tools like Ocean Data View (ODV, http://odv.awi.de).

Data visualization with ODV (http://odv.awi.de)

## Author contacts:

[1] MARUM - University of Bremen, PANGAEA, Bremen, Germany (uschindler@pangaea.de, mdiepenbroek@pangaea.de)

[2] Foundation Alfred Wegener Institute for Polar and Marine Research (AWI), PANGAEA, Bremerhaven, Germany (hgrobe@pangaea.de, rsieger@pangaea.de)