



Generic Research Data Infrastructure

Forschungsdatenkolloquium am 7. Juli 2017

Wilhelm Hasselbring

The GeRDI Team



Prof. Dr. Klaus Tochtermann



Prof. Dr. Wolfgang E. Nagel



Prof. Dr. Hans-Joachim Bungartz
Dr. Christian Grimm



Prof. Dr. Wilhelm Hasselbring

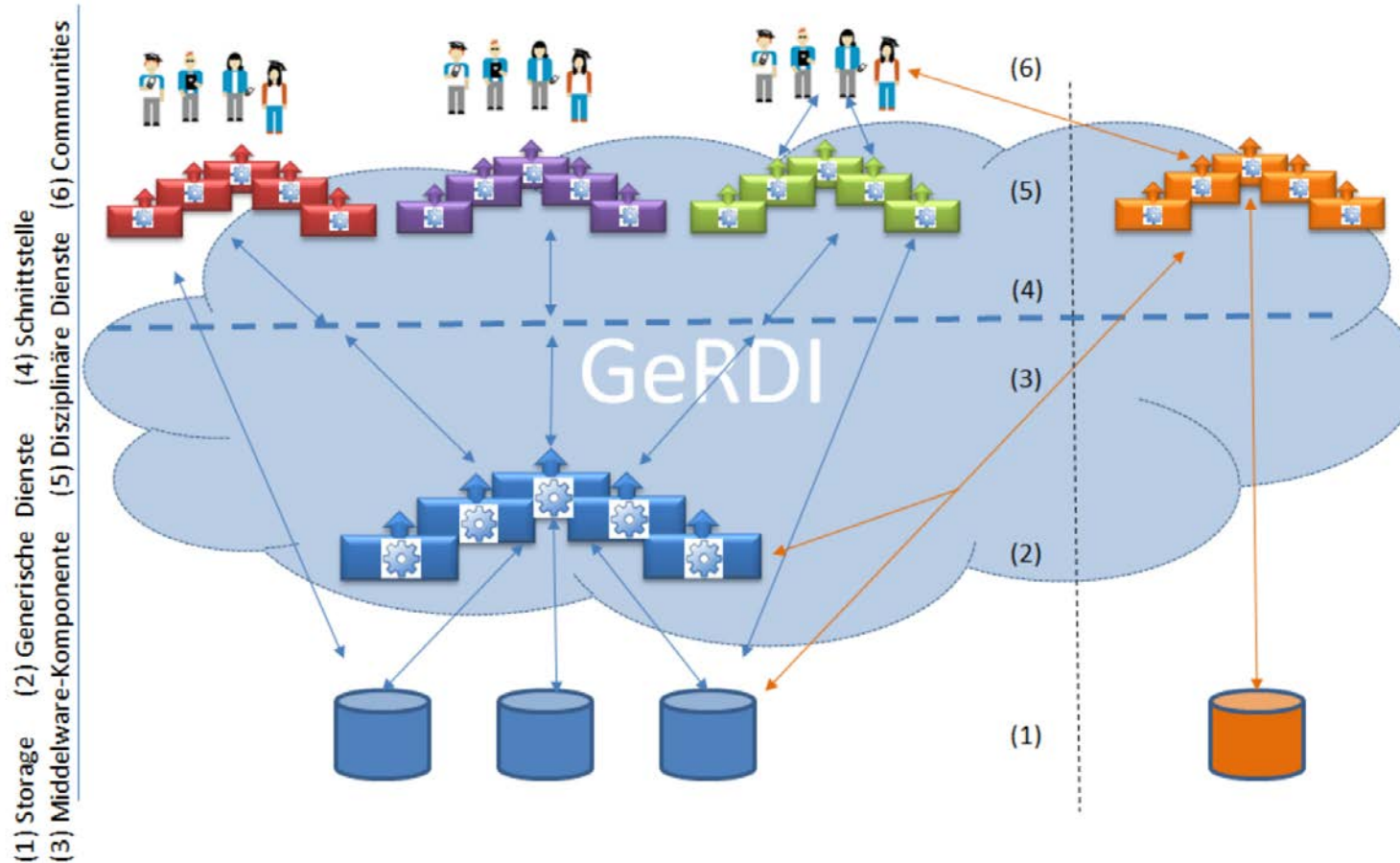


Prof. Dr. Arndt Bode
Prof. Dieter Kranzlmüller

GeRDI Vision and Mission

- Vision:
 - Multidisciplinary and FAIR research data management principles are widely accepted.
 - The common application of these principles results in great benefits for research, industry, society and our environment.
- Mission:
 - Contribution to the vision (not yet explicitly defined) with initial focus on the present evaluation communities.

Envisioned GeRDI Architecture (From the Proposal)



Evaluation Communities

GeRDI Evaluation Communities



<p>Economics</p>	<p>Life science, Humanities</p>
<p>Marine science</p>	<p>Environmental science</p>



Munich communities (LRZ)

Alpine Environment Data Analysis Center (AlpEnDac)

- Environmental Science and Medicine
- Geophysical health data from Climate, Glaciology, Radiology research

Hydrology and River Basin Management (HIOS)

- Environmental Science
- Geophysical data from Hydrology:
Water Resource and Decentralized Flood Management

UN International Strategy for Disaster Reduction (UNISDR)

- Socio-economic and geophysical data from research in:
Disaster Risk Reduction and Sociology of Disaster Modelling

Dresden communities (TUD)

Microscopy and Bioinformatics (CBG)

- Cell Biology and Genetics
- Images and sequence data from Gentic

Digital Humanities (Prof. Crane)

- Automatic text analysis and philology
- Textual data from image, text and speech analysis

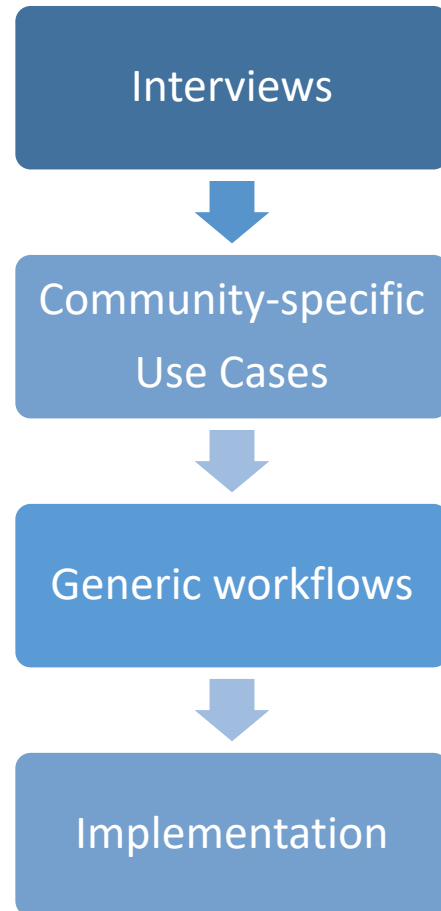
National Center for Tumor Diseases (NCT)

- Medial research in tumors
- Clinical and epidemiological studies and collection of biomaterial

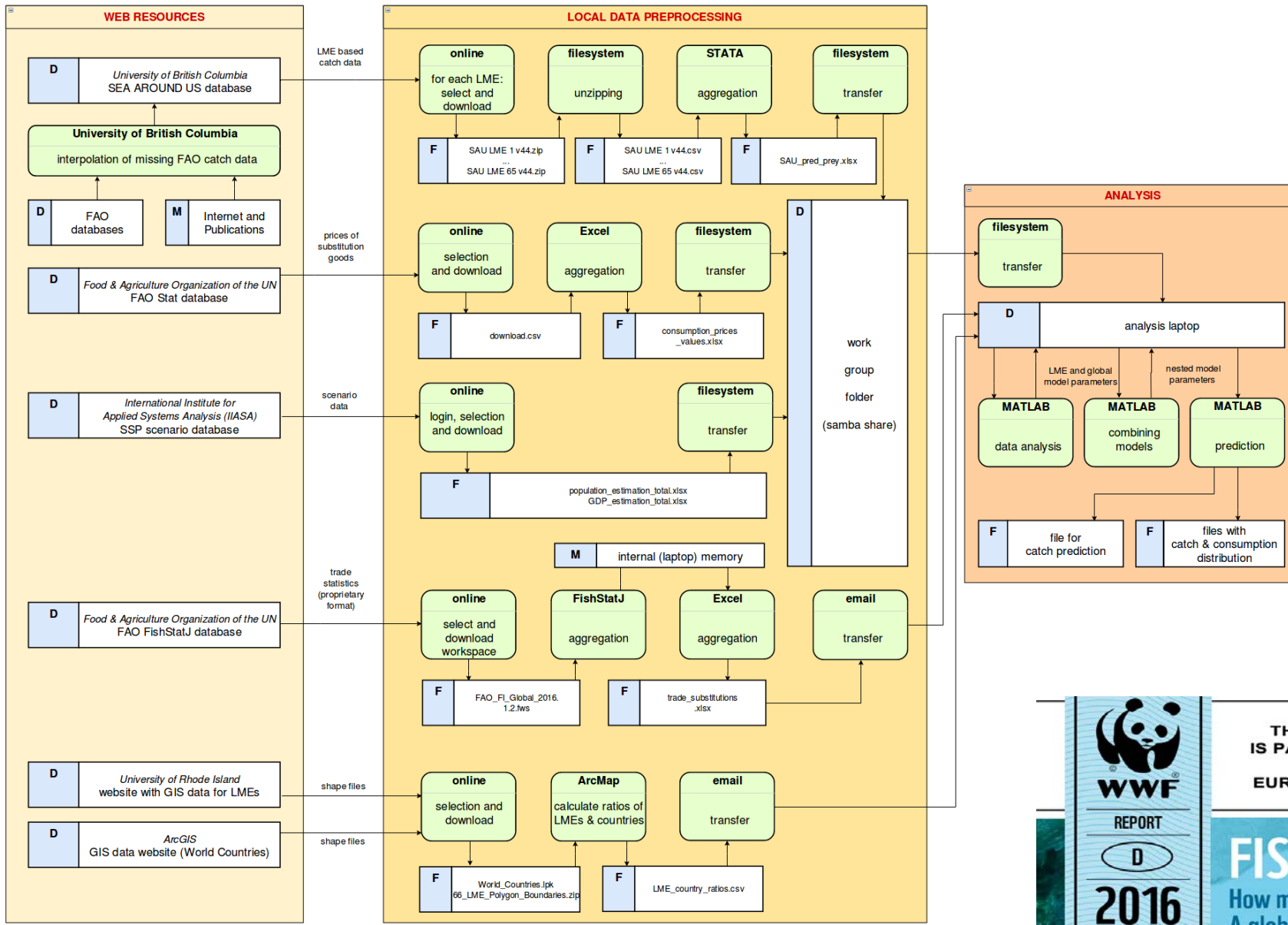
Kiel communities (ZBW, CAU)

- Socio-Economic Panel (SOEP):
Prof. Schupp, DIW
 - Socioeconomics
 - From wide ranging representative studies of private households:
Qualitative & quantitative data on composition, occupation, earnings, ...
- Enviromental, Resource and Ecological Ecomomics (EREE):
Prof. Quaas, Future Ocean
- Paleoceanography:
Prof. Dullo, GEOMAR

Requirements Engineering for Community Involvement



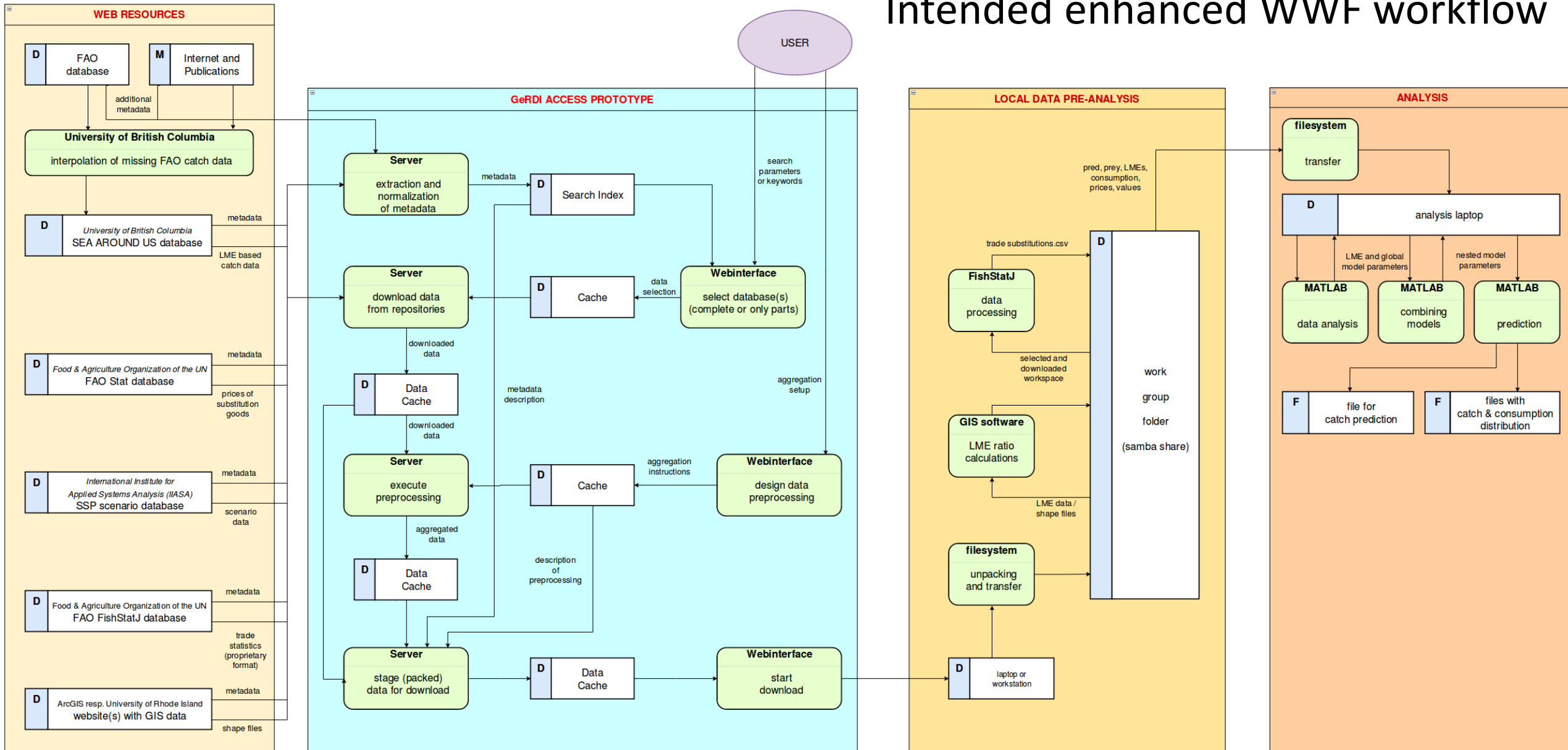
- Identification of ...
 - research workflows
 - relevant data repositories
- Elicitation of community-specific use cases
- Extraction of generic workflows



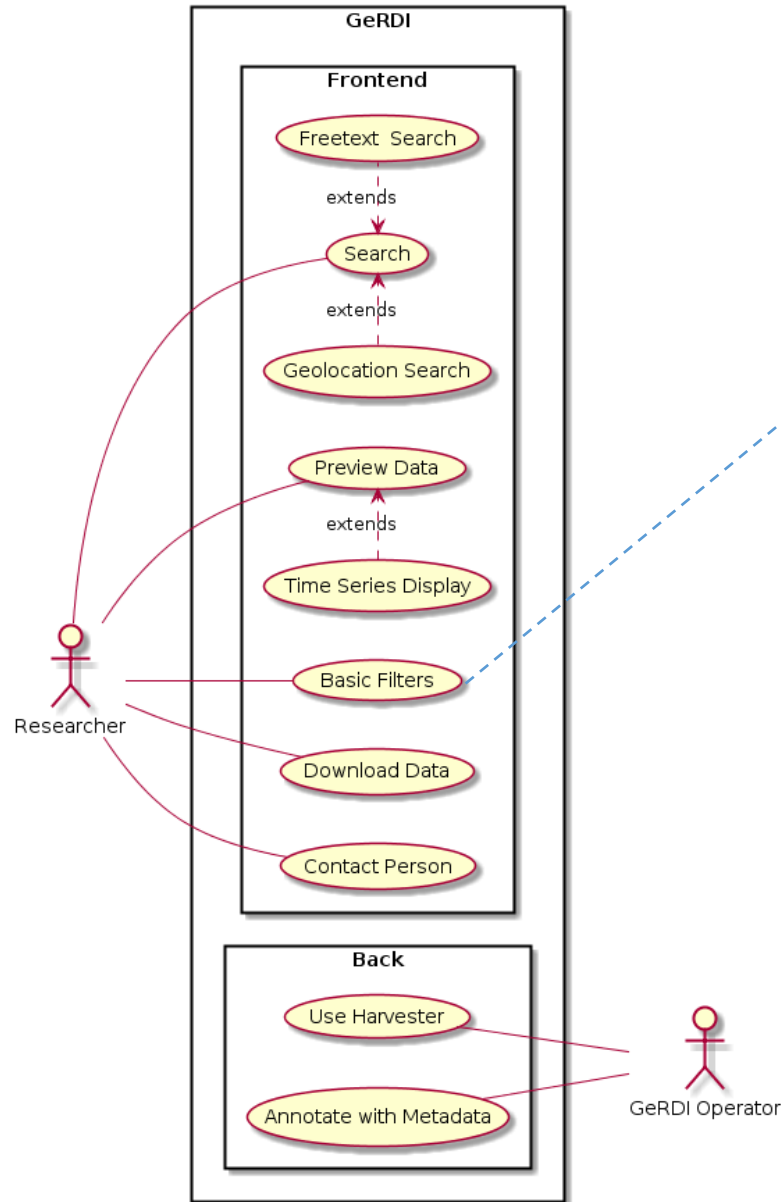
EREE
Scientific
Workflow,
as is.



Intended enhanced WWF workflow



Use cases



Use Case	Filter the results of the search
ID	SAI-121 - Filter the search results BACKLOG
Community	All communities
Precondition	There are relevant search results that can be filtered
Trigger	The researcher wants to use the filter function
Goal	The researcher get the filtered results of the search
Main Scenario	<ol style="list-style-type: none"> 1. The researcher click on filter he/she wants to apply 2. The researcher apply the filter, he/she has chosen 3. The results of the filtering are displayed
Extension	
Successful Postcondition	The researcher was able to limit the search results and so get the most relevant items
Failure Postcondition	The results of the filtering are not relevant or are empty

From requirements to implementation and automatic regression tests

- Use Case 2.0:
Combination of use cases with agile software development
- Use Case Slices:
Various flows through main and extended scenarios
- Described using acceptance tests formulated using BDD (test-driven development)

Use Case 2.0
<https://www.ivarjacobson.com/publications/white-papers/use-case-ebook>

Feature: Main Workflow

As a researcher

I want to get filtered results of the search

So I can find the relevant data easier

Scenario: The researcher clicks on filter he wants to apply

Given I am on the search page

And there are search results

When I click on filter button

Then the filter options are displayed

Scenario: Applying filter option [...]



Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Ecological Informatics

journal homepage: www.elsevier.com/locate/ecolinf

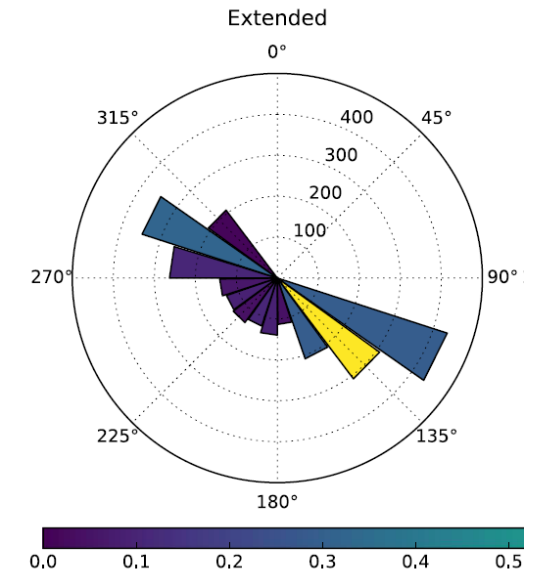
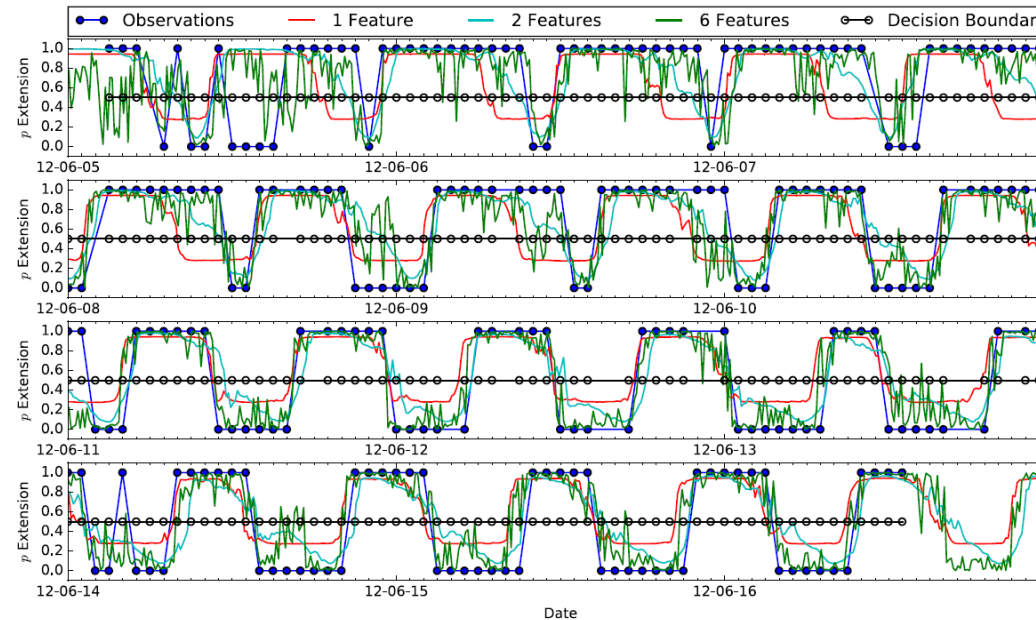
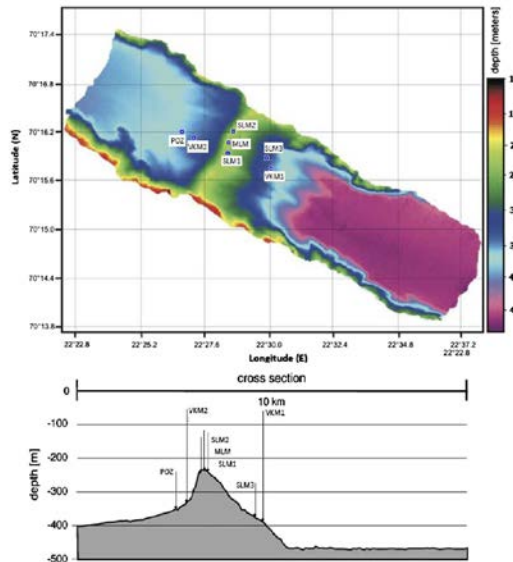


Modeling polyp activity of *Paragorgia arborea* using supervised learning

Arne N. Johanson^{a,*}, Sascha Flögel^b, Wolf-Christian Dullo^b, Peter Linke^b, Wilhelm Hasselbring^a

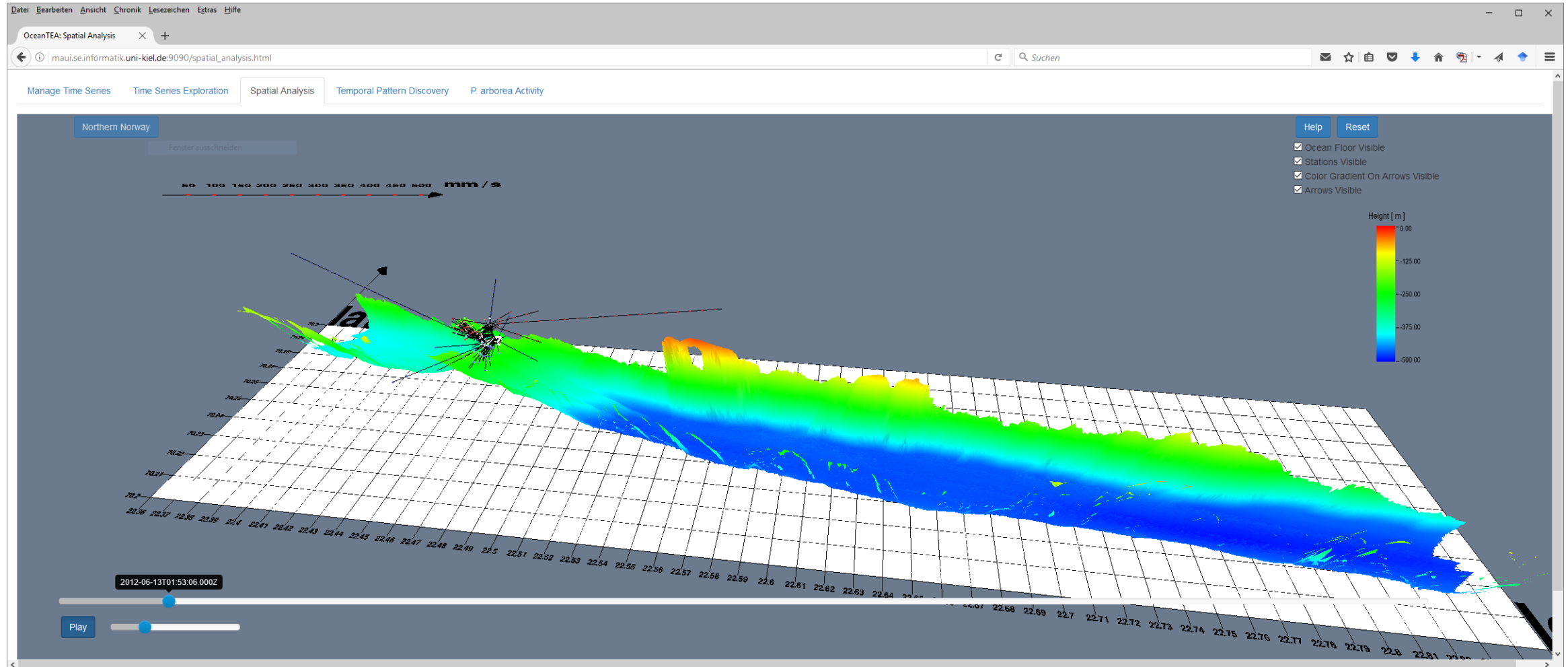
^a Software Engineering Group, Kiel University, Germany

^b GEOMAR Helmholtz Centre for Ocean Research, Kiel, Germany



Data available in OceanTEA

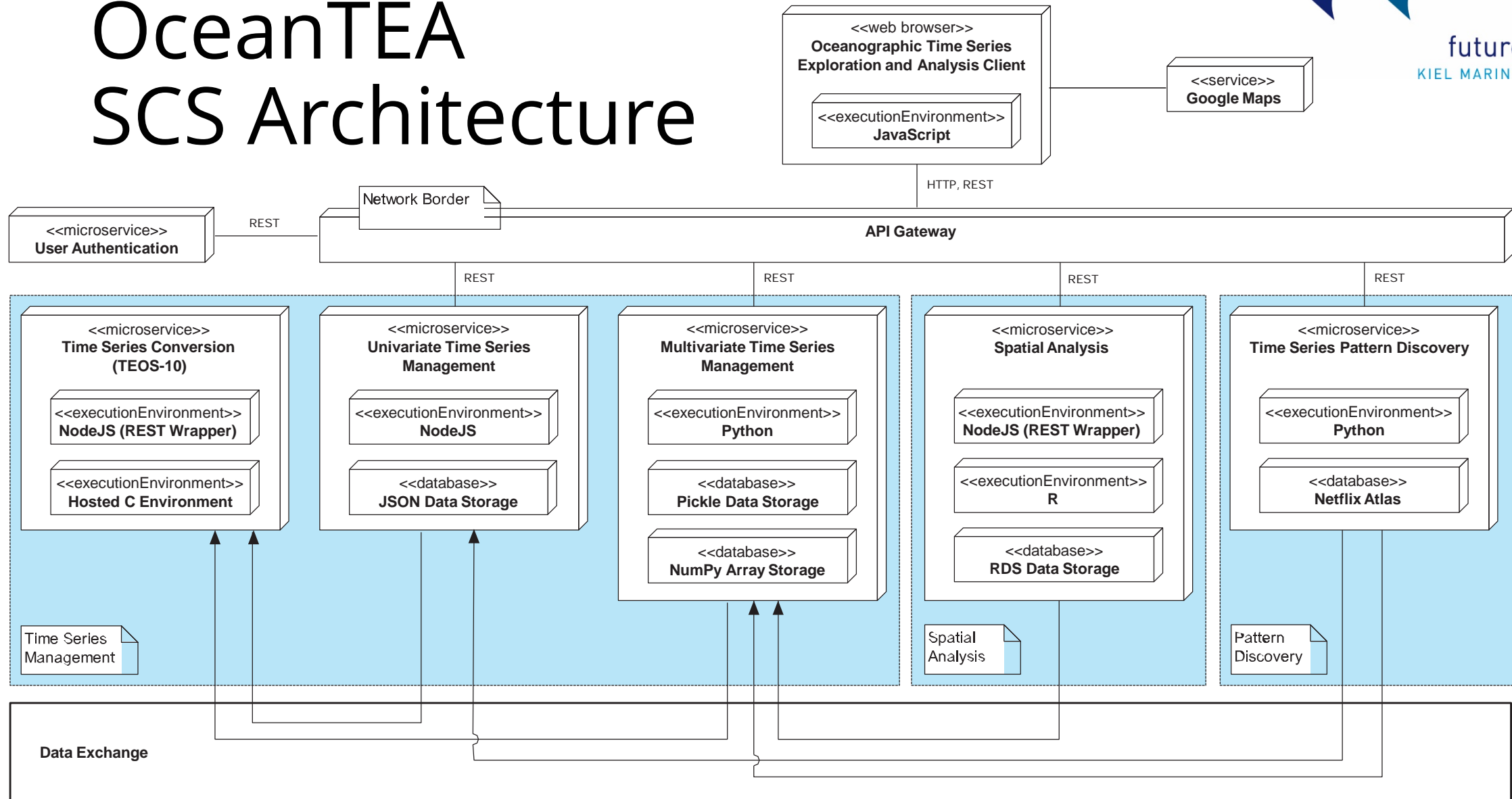
<http://maui.se.informatik.uni-kiel.de:9090/>



OceanTEA SCS Architecture

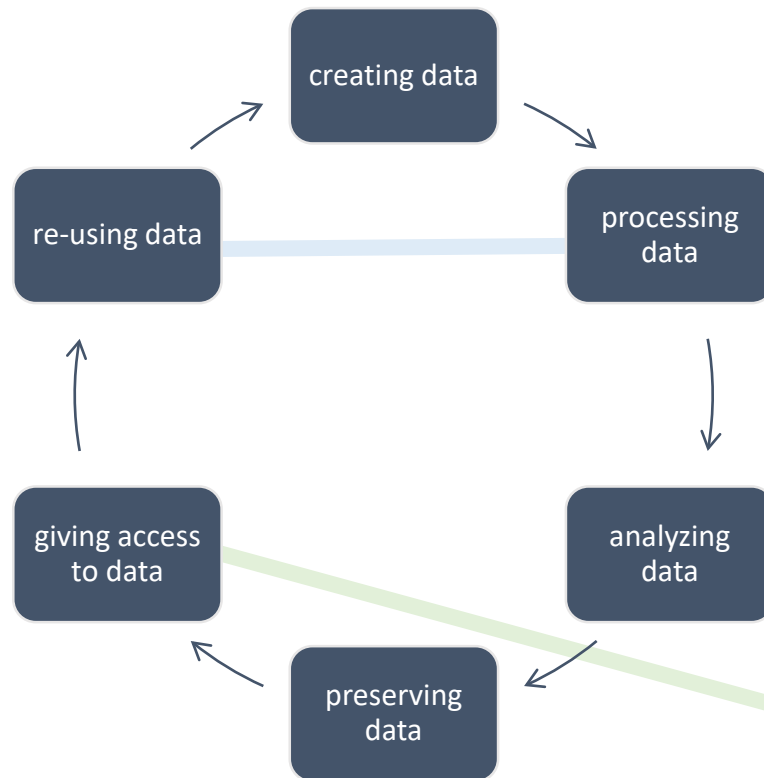


future ocean
KIEL MARINE SCIENCES



GeRDI services and workflow cycles

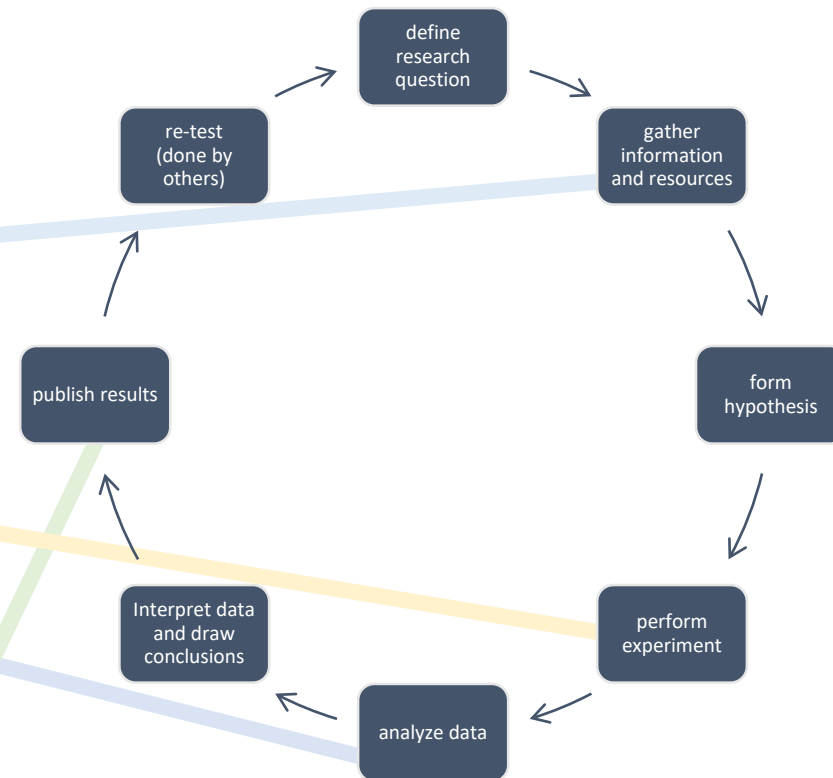
Research data lifecycle



Services

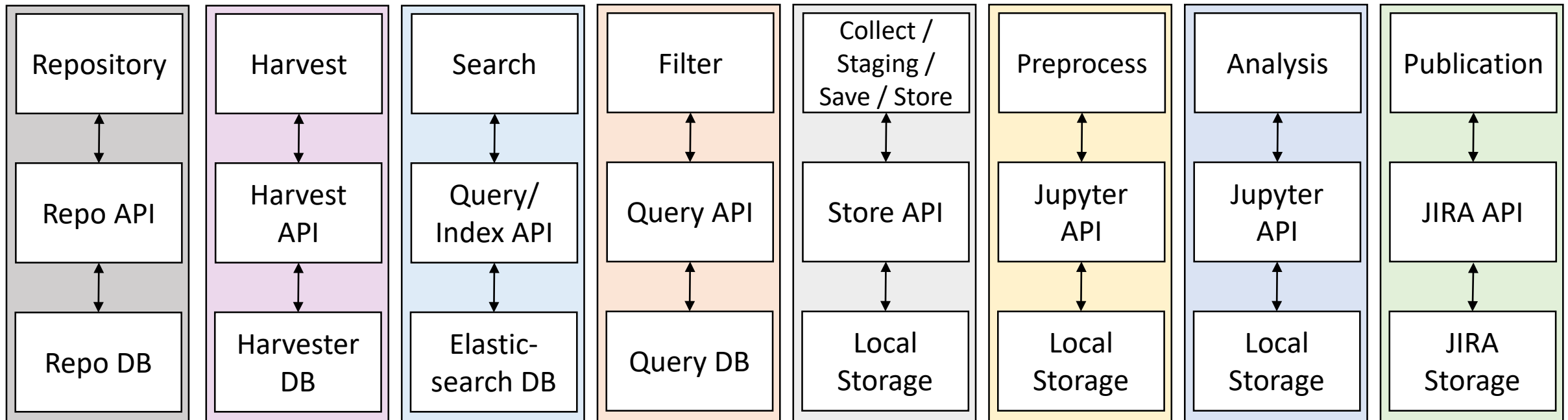
- Repository
- Harvest
- Search
- Filter
- Collect / Staging / Save / Store
- Preprocessing
- Analysis
- Publication

Research workflow



GeRDI's Architectural Style

API Gateway, Page Assembly Proxy, ...



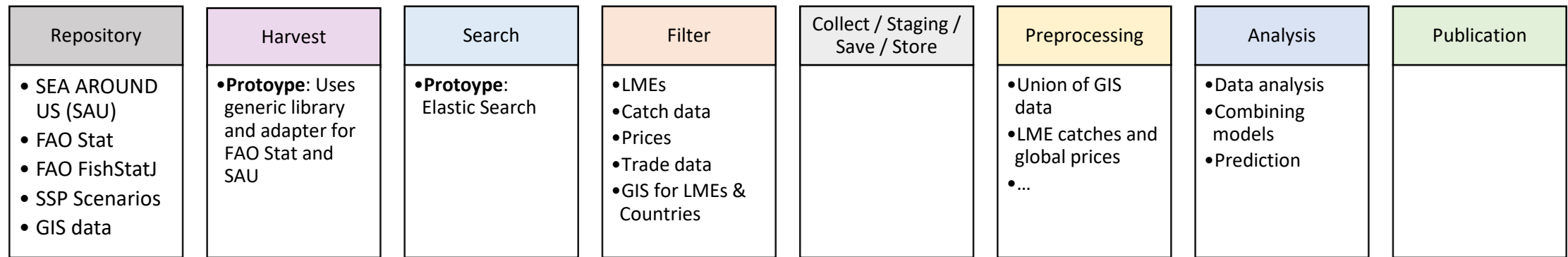
REST, Messaging, Prospector AAI, Scalability/Elasticity, Monitoring, Control Center, ...

Design Rationale

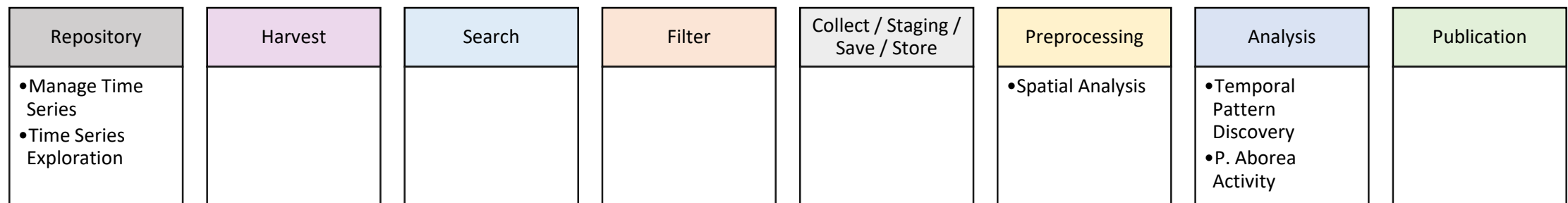
- Entry / Exit Options for users at various points in the workflow
- Self-contained Systems (SCS)
 - Also known as microservices (<http://scs-architecture.org>)
 - Each SCS is an autonomous (web) application:
For its domain all data, the logic to process that data and all code to render the web interface is contained within the SCS.
 - Communication with each other or 3rd party should be asynchronous:
This decouples the systems, reduces the effects of failure, and thus supports autonomy.
 - SCS should not share business code to avoid tight coupling.

GeRDI services and community workflows (current snapshot at CAU)

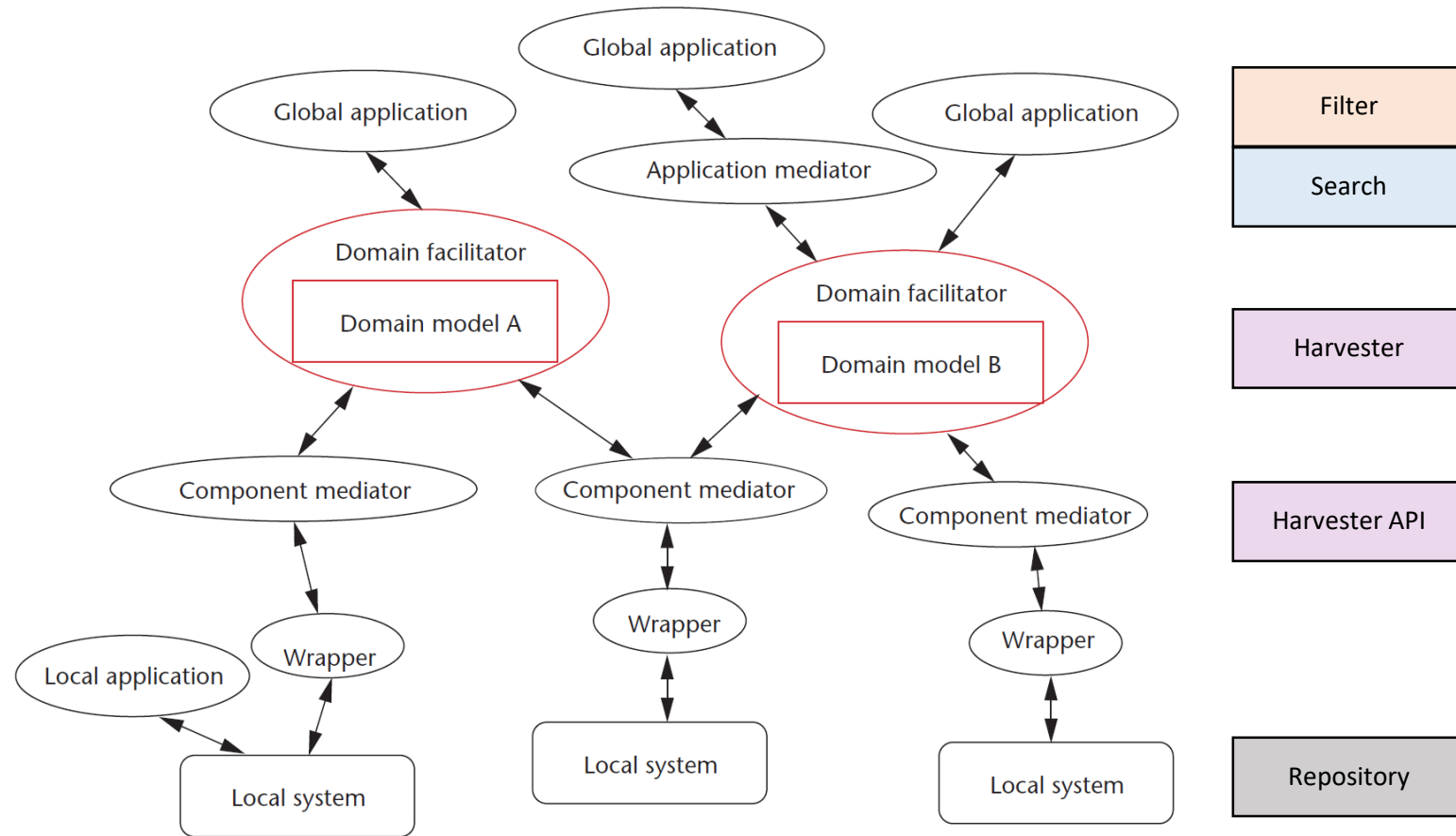
WWF Paper (EREE, Prof. Quaas)



OceanTEA (Bio & Environment, Prof. Dullo)



Mediator Architecture and Metadata



Source:
 Hasselbring, W.
 (2002) "Web Data
 Integration for E-
 Commerce
 Applications", IEEE
 Multimedia, 9 (1).
 pp. 16-25. DOI
 10.1109/93.978351.



GeRDI as a software product line

- “A software product line (SPL) is a set of software-intensive systems that share a common, managed set of **features** satisfying the specific needs of a particular market segment or mission and that are developed from a common set of core assets in a prescribed way.”
[Software Engineering Institute, CMU]
- Software product lines or application families are applications with **generic** functionality that can be adapted and configured for use in a specific context.

Demo

Live demo our search prototype

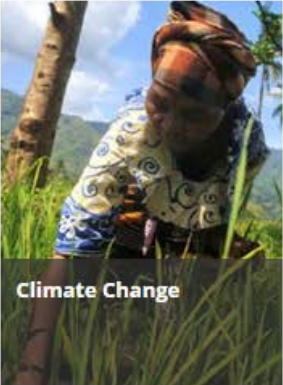
Without Gerdi (1 / 5)



- About FAO
- In Action
- Countries
- Themes
- Media
- Publications
- Statistics
- Partnerships



Making the international trade in plants and seeds a safer venture
International plant health body adopts new global standards



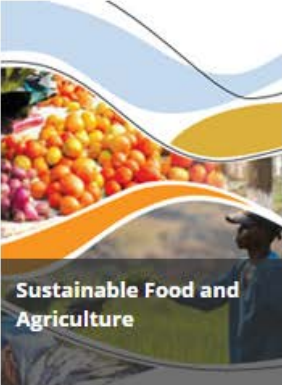
Climate Change



FACING FAMINE



FAO and the Sustainable Development Goals



Sustainable Food and Agriculture



Resilience

Without Gerdi (2 / 5)



International Conference on Agricultural Statistics VII

Modernization of Agricultural Statistics in
Support of the Sustainable Development
Agenda



Statistics at FAO

FAO develops methods and standards for food and agriculture statistics, provides technical assistance services and disseminates data for global monitoring. Statistical activities at FAO include the development and implementation of methodologies and standards for data collection, validation, processing and analysis. FAO also plays a vital part in the global compilation, processing and dissemination of food and agriculture statistics, and provides essential statistical capacity development to member countries.

FAO has a decentralized statistical system and statistical activities cover the areas of agriculture, forestry and fisheries, land and water resources and use, climate, environment, population, gender, nutrition, poverty, rural development, education and health as well as many others.

FAO Statistical Programme of Work

The FAO Statistical Programme of Work is a collaborative effort that is overseen by the Chief Statistician and supported by the Inter-Departmental Working Group on Statistics. These two mechanisms ensure strengthened coordination and cooperation on statistical matters and guarantee the high quality of FAO data.

The Statistical Programme of Work provides a summary of all of the principal statistical activities at FAO, and a detailed description of all the individual statistical activities carried out by FAO Divisions active in the field of statistics. It presents the organization's operational activities according to different statistical categories and domains.



Get FAO statistics

FAOSTAT: the world's la...



Related links

- [FAO statistics website](#)
- [Food security statistics](#)
- [Forestry statistics](#)
- [Fisheries and aquaculture](#)
- [Trade and markets](#)
- [Water information and statistics](#)
- [Coordinating Working Party on Fishery Statistics](#)

Data products

Without Gerdi (3 / 5)

FAOSTAT



FAOSTAT provides free access to food and agriculture data for over 245 countries and territories and covers all FAO regions from 1961 to the most recent year available.

[Explore Data](#)



Database Updates

[Macro Indicators \(Macro-Statistics\)](#)
March 29, 2017

[Deflators \(Prices\)](#)
March 24, 2017



FAO Statistical Yearbooks

The FAO Statistical Yearbook provides a selection of indicators on food and agriculture by country.

The first part of the book includes thematic spreads with data visualizations (graphs, charts, and maps) with basic text. The second part has country-level tables for a selected number of

Bulk Download

All FAOSTAT Data

Updated on Mar 29, 2017

Tweets by @FAOstat



#didyouknow? Net #forest countries has increased on bit.ly/2bUPJzP #GFRA

Without Gerdi (4 / 5)

FAOSTAT

Home Data Country Indicators Compare Data Definitions and Standards FAQ Search an

Data

DOMAINS

Filter the domain list e.g. crops, food security, fertilizers



Production

- Crops
- Crops Processed
- Live Animals
- Livestock Primary
- Livestock Processed
- Production Indices
- Value of Agricultural Production



Inputs

- Fertilizers
- Fertilizers archive
- Fertilizers - Trade Value
- Pesticides Use
- Pesticides Trade
- Land
- Employment Indicators



Emissions

- Agriculture Total
- Enteric Fermentation
- Manure Management
- Rice Cultivation
- Synthetic Fertilizer
- Manure applied to land
- Manure left on pastures
- Crop Residues
- Cultivation of Organic Soils
- Burning - Savanna
- Burning - Crop Residues
- Energy Use



Trade

- Crops and livestock products
- Live animals
- Detailed trade matrix
- Trade Indices



Population

- Annual population



Investment


- Machinery
- Machinery Archive
- Government Expenditure
- Credit to Agriculture



Emissions

- Land Use Total
- Forest Land
- Cropland

Without Gerdi (5 / 5)

COUNTRIES REGIONS SPECIAL GROUPS 

Filter results e.g. afghanistan

- Germany
- Ghana
- Gibraltar
- Greece
- Greenland
- Grenada
- Guadeloupe

Select All Clear All

Greece x


ELEMENTS

Filter results e.g. area harvested

- Area harvested
- Yield
- Production Quantity
- Seed

Select All Clear All

Production Quantity x

ITEMS ITEMS AGGREGATED 

Filter results e.g. agave fibres nes

- Fruit, pome nes
- Fruit, stone nes
- Fruit, tropical fresh nes
- Garlic
- Ginger
- Gooseberries
- Grain, mixed

Select All Clear All

Garlic x

YEARS

Filter results e.g. 2014

- 2014
- 2013
- 2012
- 2011
- 2010
- 2009




Select All Clear All

2010 x

Output Type: Table Pivot

Thousand Separator in 'Show Data': None Comma Period

Output Formatting Options: Flags Codes Units Null Values

 Show Data  Download Data 

Integrated search in multiple repositories with GeRDI

fish AND Norway

6 results found in 25 ms

Commodity Balances - Livestock and Fish Primary Equivalent
<http://www.fao.org/faostat/en/#data/BL>
Food supply data is some of the most important data in FAOSTAT. In fact, this data is for the basis for estimation of global and national undernourishment assessment, when it is combined with parameters and other data sets. This data has been the foundation of food balance sheets ever since they were first constructed. The data is accessed by both business and governments for economic analysis and policy setting, as well as being used by the academic community.

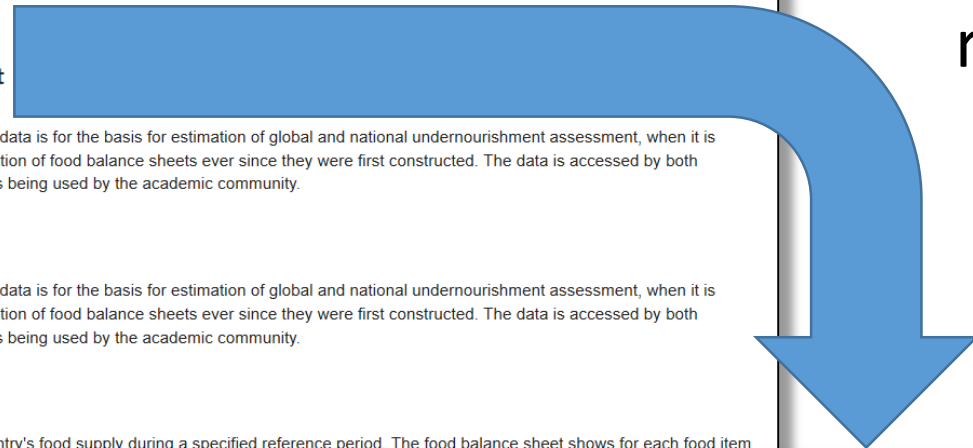
Food Supply - Livestock and Fish Primary Equivalent
<http://www.fao.org/faostat/en/#data/CL>
Food supply data is some of the most important data in FAOSTAT. In fact, this data is for the basis for estimation of global and national undernourishment assessment, when it is combined with parameters and other data sets. This data has been the foundation of food balance sheets ever since they were first constructed. The data is accessed by both business and governments for economic analysis and policy setting, as well as being used by the academic community.

Food Balance Sheets
<http://www.fao.org/faostat/en/#data/FBS>
Food Balance Sheet presents a comprehensive picture of the pattern of a country's food supply during a specified reference period. The food balance sheet shows for each food item - i.e. each primary commodity and a number of processed commodities potentially available for human consumption - the quantity of each foodstuff produced in a country added to the total quantity imported and adjusted to any change in the supply available during that period. On the utilization side a distinction is made between the quantity used for food and non-food uses, losses during storage and transportation, and food supplies available for human consumption.

Live animals
<http://www.fao.org/faostat/en/#data/TA>
The food and agricultural trade dataset is collected, processed and disseminated by FAO according to the standard International Merchandise Trade Statistics Methodology. The data is mainly provided by UNSD, Eurostat, and other national authorities as needed. This source data is checked for outliers, trade partner data is used for non-reporting countries or missing cells, and data on food aid is added to take into account total cross-border trade flows. The trade database includes the following variables: export quantity, export value, import quantity and import value. The trade database includes all food and agricultural products imported/exported annually by all the countries in the world.

Trade Indices
<http://www.fao.org/faostat/en/#data/TI>
The food and agricultural trade dataset is collected, processed and disseminated by FAO according to the standard International Merchandise Trade Statistics Methodology. The data is mainly provided by UNSD, Eurostat, and other national authorities as needed. This source data is checked for outliers, trade partner data is used for non-reporting countries or missing cells, and data on food aid is added to take into account total cross-border trade flows. The trade database includes the following variables: export quantity, export value, import quantity and import value. The trade database includes all food and agricultural products imported/exported annually by all the countries in the world.

Crops and livestock products
<http://www.fao.org/faostat/en/#data/TP>
The food and agricultural trade dataset is collected, processed and disseminated by FAO according to the standard International Merchandise Trade Statistics Methodology. The data is mainly provided by UNSD, Eurostat, and other national authorities as needed. This source data is checked for outliers, trade partner data is used for non-reporting countries or missing cells, and data on food aid is added to take into account total cross-border trade flows. The trade database includes the following variables: export quantity, export value, import quantity, and import value. The trade database includes all food and agricultural products imported/exported annually by all the countries in the world.



Commodity Balances - Livestock and Fish Primary Equivalent

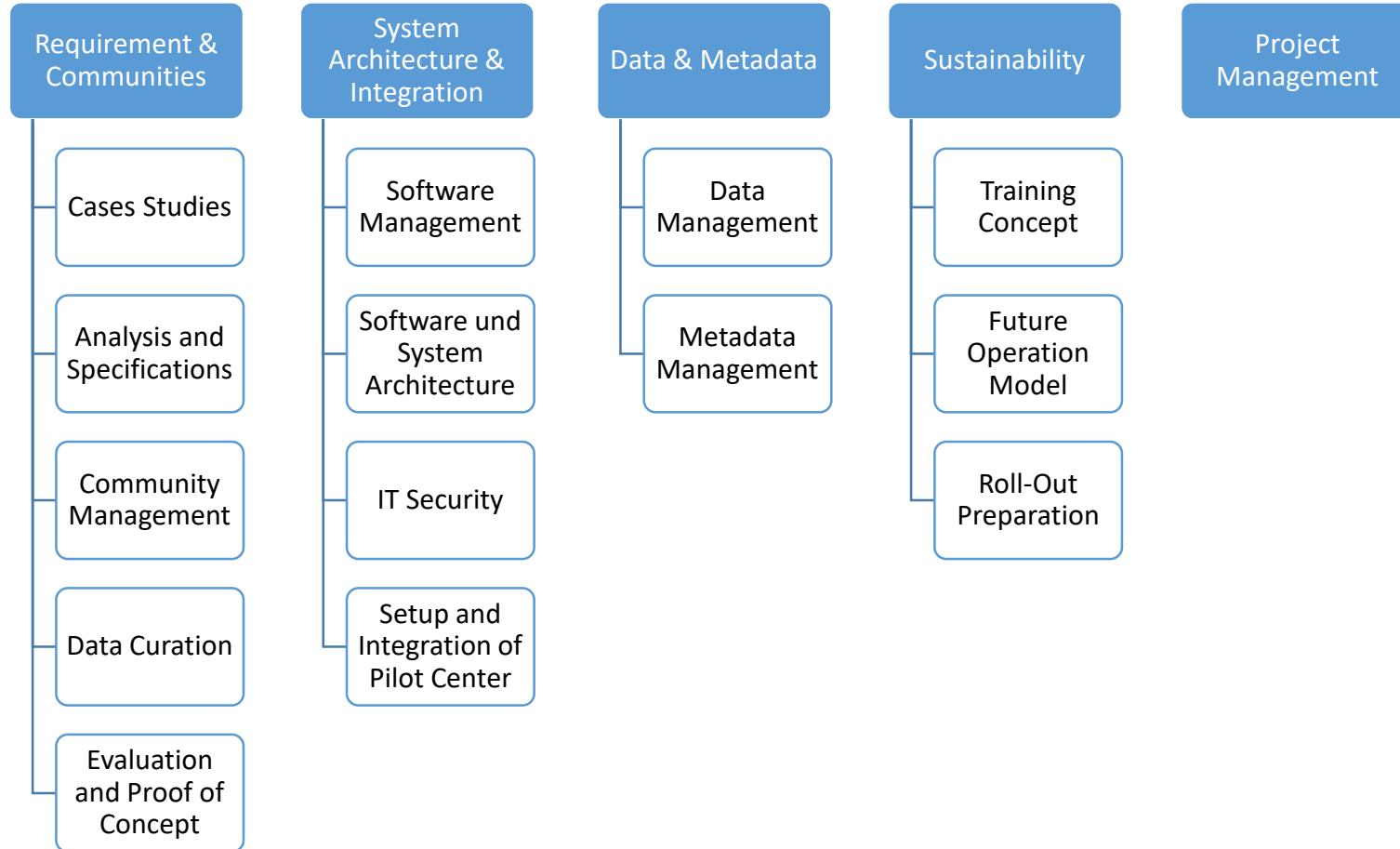


DESCRIPTION

Food supply data is some of the most important data in FAOSTAT. In fact, this data is for the basis for estimation of global and national undernourishment assessment, when it is combined with parameters and other data sets. This data has been the foundation of food balance sheets ever since they were first constructed. The data is accessed by both business and governments for economic analysis and policy setting, as well as being used by the academic community.

ADDITIONAL METADATA

Work Packages



Summary & Outlook

- GeRDI as Contribution to the European Open Science Cloud
<https://www.youtube.com/watch?v=SC4-O8BmI4I>
- Funding: 3 M€ for three years (first phase)
- KOLab
 - Kiel Open Data and Software Lab
- Digital Ocean 2017
 - 20. September 2017, 12-18 h at ZMB
https://www.kms.uni-kiel.de/de/veranstaltungen/digital_ocean2017
- More to come at
 - <http://www.gerdi-project.de/>



References

- (1) Grunzke, R., Adolph, T., Biardzki, C., Bode, A., Borst, T., Bungartz, H. J., Busch, A., Frank, A., Grimm, C., Hasselbring, W., Kazakova, A., Latif, A., Limani, F., Neumann, M., de Sousa, N. T., Tendel, J., Thomsen, I., Tochtermann, K., Müller-Pfefferkorn, R. and Nagel, W. E. (2017): “Challenges in Creating a Sustainable Generic Research Data Infrastructure” In: Softwaretechnik-Trends, 37 (2).
- (2) Hasselbring, W. and Steinacker, G. (2017): “Microservice Architectures for Scalability, Agility and Reliability in E-Commerce” In: IEEE International Conference on Software Architecture 2017, April 03-07, 2017, Gothenburg, Sweden.
- (3) Hasselbring, W. (2002): “Web Data Integration for E-Commerce Applications”, IEEE Multimedia, 9 (1). pp. 16-25. DOI 10.1109/93.978351.
- (4) Hasselbring, W. (2016): “Microservices for Scalability” In: International Conference on Performance Engineering (ICPE 2016), March 2016, Delft, Netherlands.
- (5) Johanson, A., Flögel, S., Dullo, W. C., Linke, P. and Hasselbring, W. (2017): “Modeling Polyp Activity of Paragorgia arborea using Supervised Learning” In: Ecological Informatics, 39 . pp. 109-118. DOI 10.1016/j.ecoinf.2017.02.007.
- (6) Johanson, A., Flögel, S., Dullo, C. und Hasselbring, W. (2016): “OceanTEA: Exploring Ocean-Derived Climate Data Using Microservices” In: Sixth International Workshop on Climate Informatics (CI 2016), September 22-23. 2016, Boulder, Colorado.