



Wir und die Daten: Ein Beitrag aus der Perspektive der Biogeochemie & Erdystemanalyse

Martin Jung

Miguel Mahecha

Markus Reichstein

Aspekte

- Womit befassen wir uns?
- Wie sehen unsere Daten und Datenmanagementansätze aus?
- Womit haben wir zu kämpfen?



NOAH Visualization

System Erde



<https://www.mpg.de/mpforschung>

Max Planck FORSCHUNG

Das Wissenschaftsmagazin der Max-Planck-Gesellschaft 2.2013

B20396F

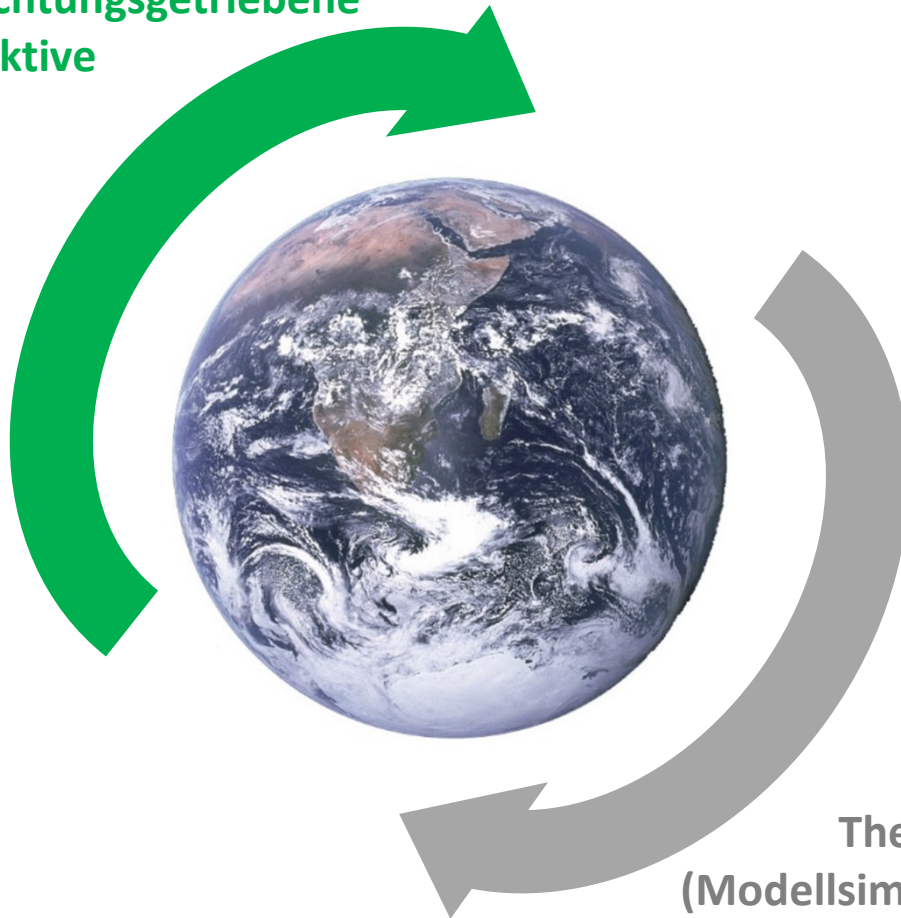


GEOWISSENSCHAFTEN

Verflochtene Erde

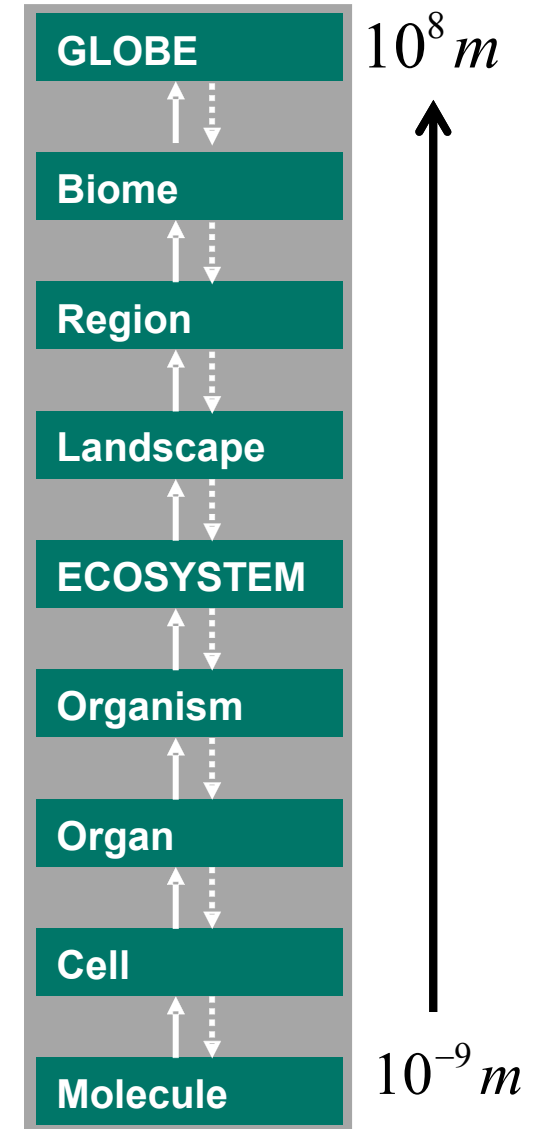
Wie funktioniert das Erdsystem?

Beobachtungsgetriebene
Perspektive



Focus: Biosphäre

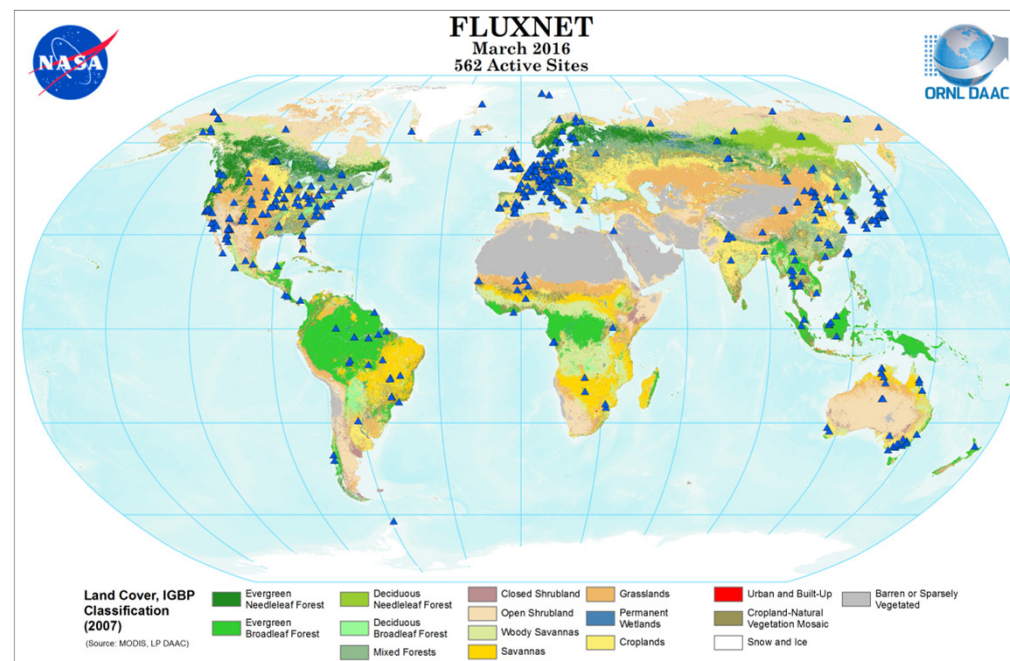
Theoretische
(Modellsimulations-)
Perspektive



Was für Daten?

„**Stationsdaten**“: (meist) Zeitserien verschiedener in-situ gemessener Variablen (z.B. Kohlenstoffaustausch zwischen Land und Atmosphäre) + diverse meta-Daten

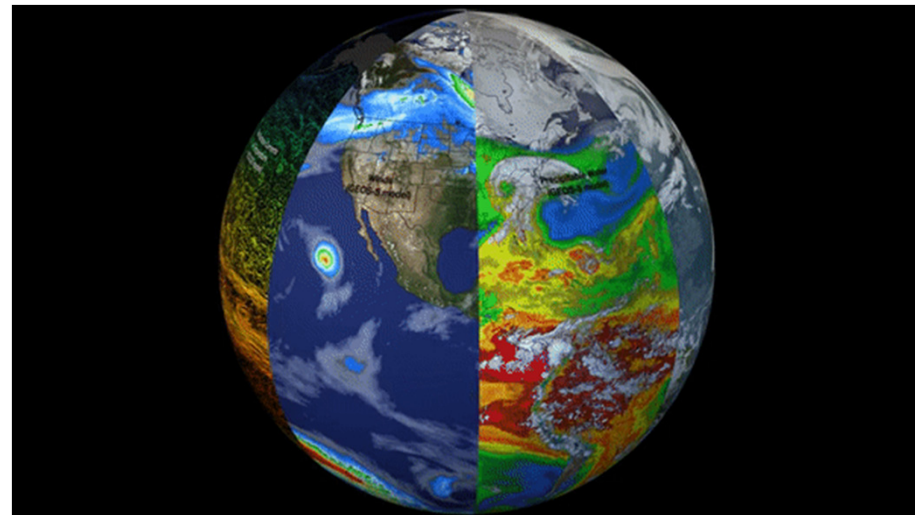
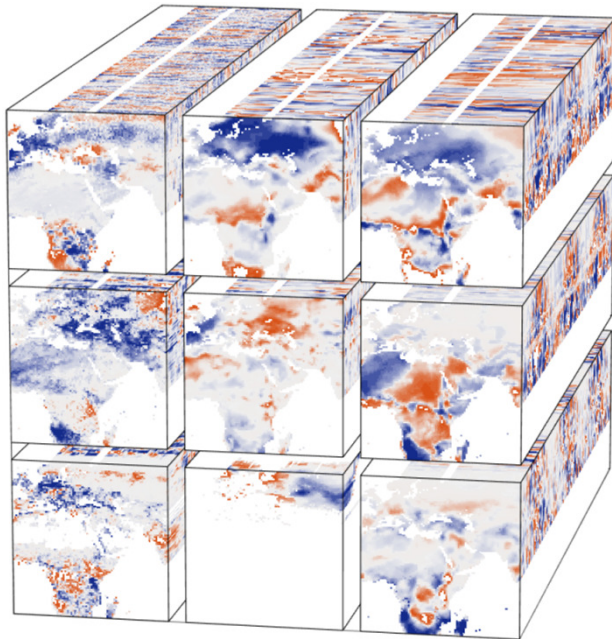
Quellen: „Kollaborative Netzwerke“



Was für Daten?

„Data Cubes“: Multivariate, räumlich und zeitlich aufgelöste „Gitter“ (grids) →
„Videos von Karten der Erde“

Quellen: Satellitenerdbeobachtung, Erdsystemmodelle



Wie sehen solche Daten aus?

```
uweber@pc027:~/data/grid/Global/0d50_daily/CERES/Ed3A/Data/Rn> ncdump -h Rn.CERES.Ed3A.720.360.2006.nc
netcdf Rn.CERES.Ed3A.720.360.2006 {
dimensions:
    lon = 720 ;
    lat = 360 ;
    time = 365 ;
variables:
    float Rn(time, lat, lon) ;
        Rn:long_name = "Net Radiation" ;
        Rn:units = "W m-2" ;
        Rn:scale_factor = 1.f ;
        Rn:add_offset = 0.f ;
        Rn:missing_value = -9999.f ;
        Rn:_FillValue = -9999.f ;
    double lon(lon) ;
        lon:long_name = "longitude" ;
        lon:standard_name = "longitude" ;
        lon:units = "degrees_east" ;
        lon:valid_range = -180.f, 180.f ;
    double lat(lat) ;
        lat:long_name = "latitude" ;
        lat:standard_name = "latitude" ;
        lat:units = "degrees_north" ;
        lat:valid_range = -90.f, 90.f ;
    double time(time) ;
        time:units = "days since 1582-10-14 00:00" ;
        time:calendar = "gregorian" ;

// global attributes:
    :title = "Net Radiation calculated from CERES input" ;
    :provided_by = "http://ceres.larc.nasa.gov/index.php" ;
    :created_by = "ulrich weber (uweber@bgc-jena.mpg.de)" ;
    :method = "Rn = (lwdown - lwup) + (swdown - swup)" ;
    :version = "Ed3A" ;
    :history = "Thu May  8 12:26:59 2014" ;
    :download = "http://ceres-tool.larc.nasa.gov/ord-tool/jsp/SYN1degSelection.jsp" ;
    :reference = "http://ceres.larc.nasa.gov/documents/DQ_summaries/CERES_SYN1deg_Ed3A_DQS.pdf" ;
}
```

BGI (\\Minerva) (M:) ▶ data ▶ DataStructureMDI ▶ DATA ▶ grid ▶ Global ▶ 0d50_daily ▶ CERES ▶ Ed3A ▶ Data ▶ Rn

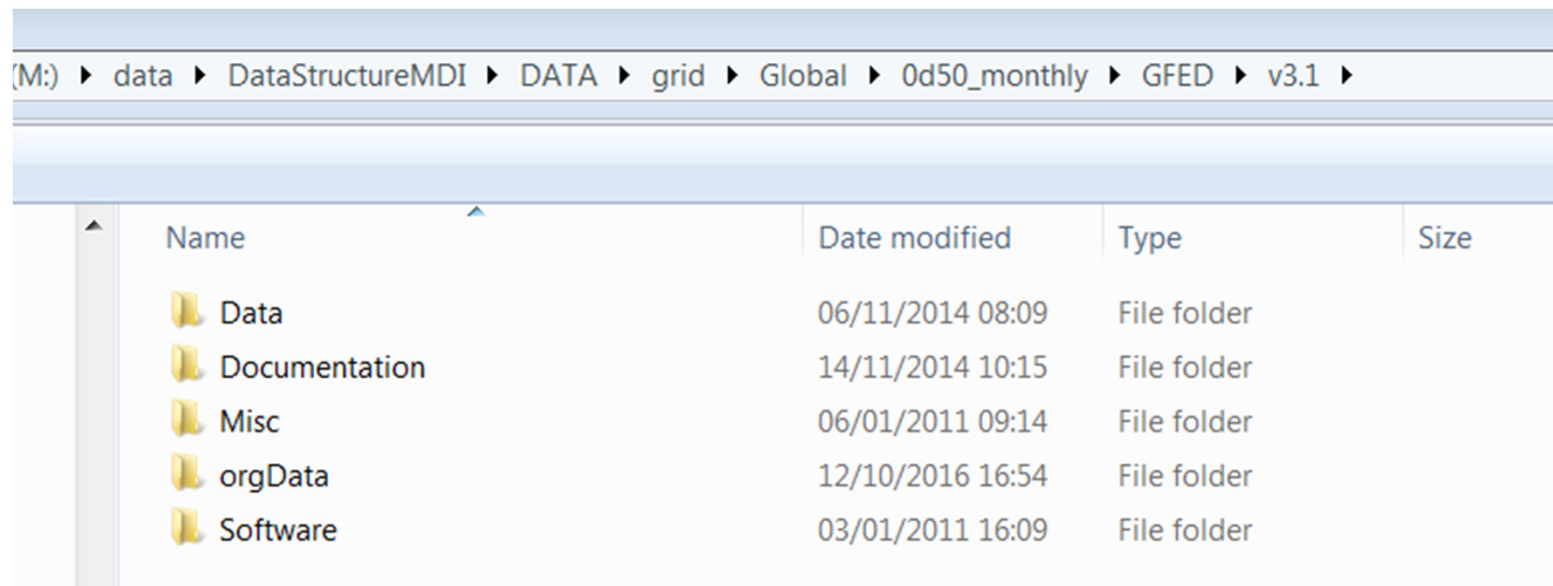
Folder

Name	Date modified	Type	Size
Rn.CERES.Ed3A.720.360.2000.nc	08/05/2014 12:26	NC File	309,838 KB
Rn.CERES.Ed3A.720.360.2001.nc	08/05/2014 12:26	NC File	369,576 KB
Rn.CERES.Ed3A.720.360.2002.nc	08/05/2014 12:26	NC File	369,576 KB
Rn.CERES.Ed3A.720.360.2003.nc	08/05/2014 12:26	NC File	369,576 KB
Rn.CERES.Ed3A.720.360.2004.nc	08/05/2014 12:26	NC File	370,588 KB

Wie verwalten wir unsere Daten?

Ordnerstruktur:

[räumlicher Ausschnitt] → [raum-zeitliche Auflösung] → [Daten-Produkt] → [Version]



Name	Date modified	Type	Size
Data	06/11/2014 08:09	File folder	
Documentation	14/11/2014 10:15	File folder	
Misc	06/01/2011 09:14	File folder	
orgData	12/10/2016 16:54	File folder	
Software	03/01/2011 16:09	File folder	

- Konsistentes Dateiformat (ncdf)
- Standardisierte Attributnamen
- Konventionen wie intern abgelegt

Datenstruktur Werkzeuge

BGI Intra : DataStructure : NetCdf Files : Search by Parameter

[Back to DataStructure](#)

Variables.VariableName = Tair

File-Path	Files
/Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Africa/0450_daily/BC_ERAIinterimV2/Data/Tair	Files
/Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Europe/0425_daily/DMI_HIRHAM_A1B_ECHAM5_CRU_DM_25KM_1951_2099V2.1/Data/Tair	Files
/Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Europe/0425_daily/KNMI_RACMO2_A1B_ECHAM5_R3_CRU_DM_25KM_1950_2100V2.1/Data/Tair	Files
/Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Europe/0425_daily/MPI-M-REMO_SCN_ECHAM5_CRU_DM_25KM_1951_2100V2.1/Data/Tair	Files
/Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Europe/0425_daily/WATCH_Interim_harmonized/DataV2.1/Tair_WFD	Files

Suchen

Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Global/0450_daily/BC_ERAIinterimV2/Data/Tair

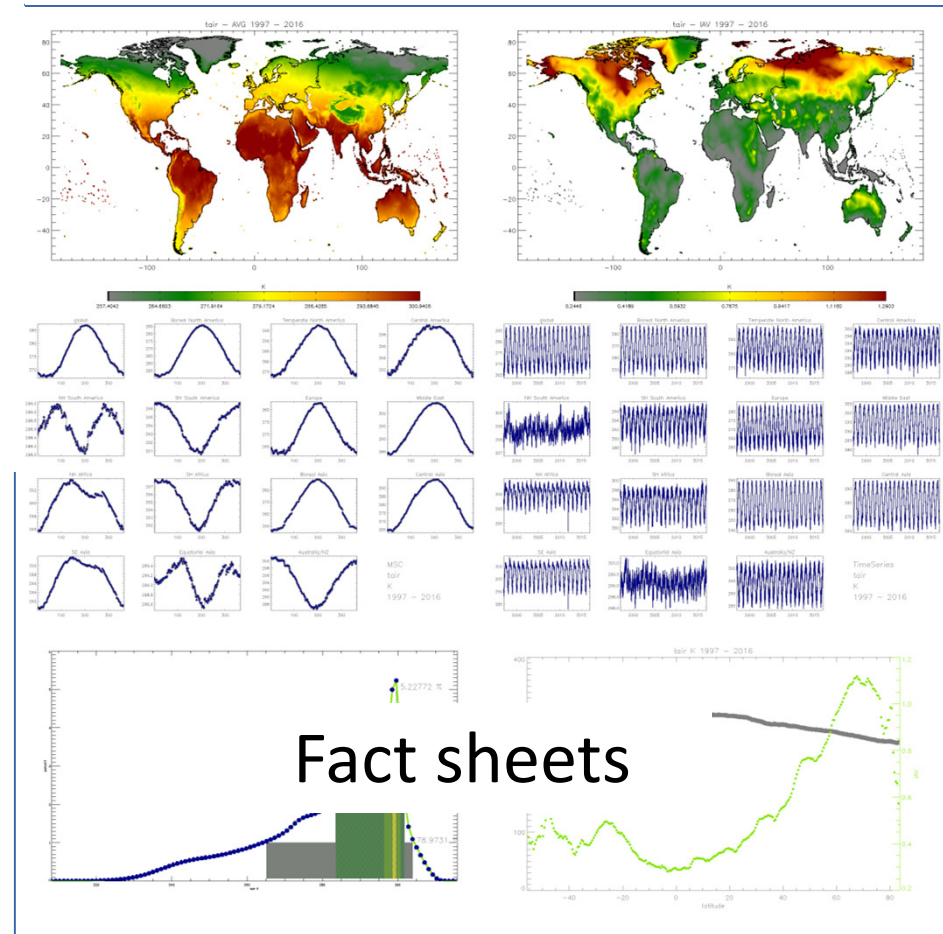
Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Global/0450_daily/CRUNCEPv4/Data/Tair

Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Global/0450_daily/CRUNCEPv5/Data/Tair

Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Global/0450_daily/CRUNCEPv6/Data/Tair

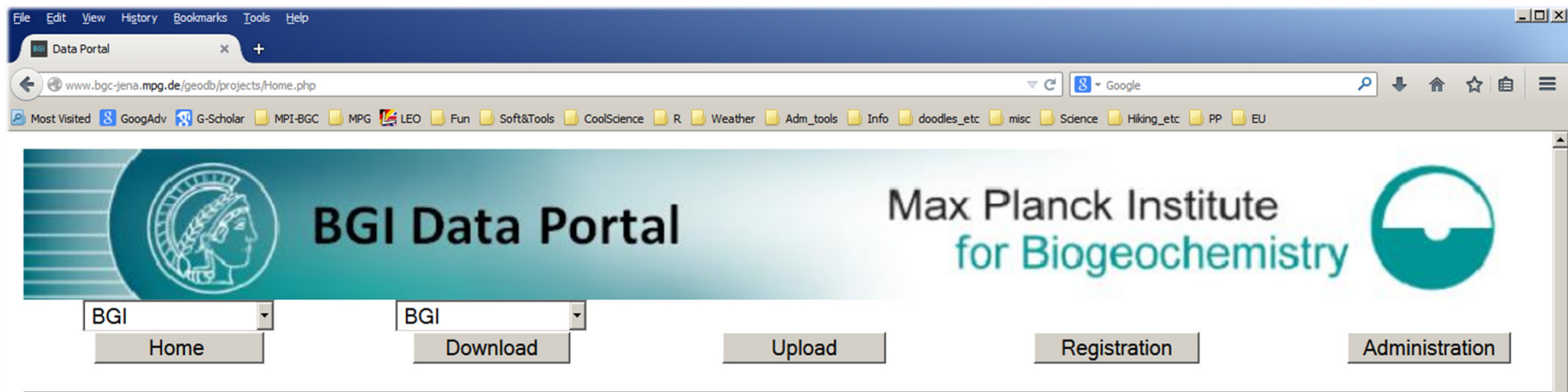
Net/Groups/BGI/data/DataStructureMDX/DATA/grid/Global/0450_daily/CRUNCEPv7/Data/Tair

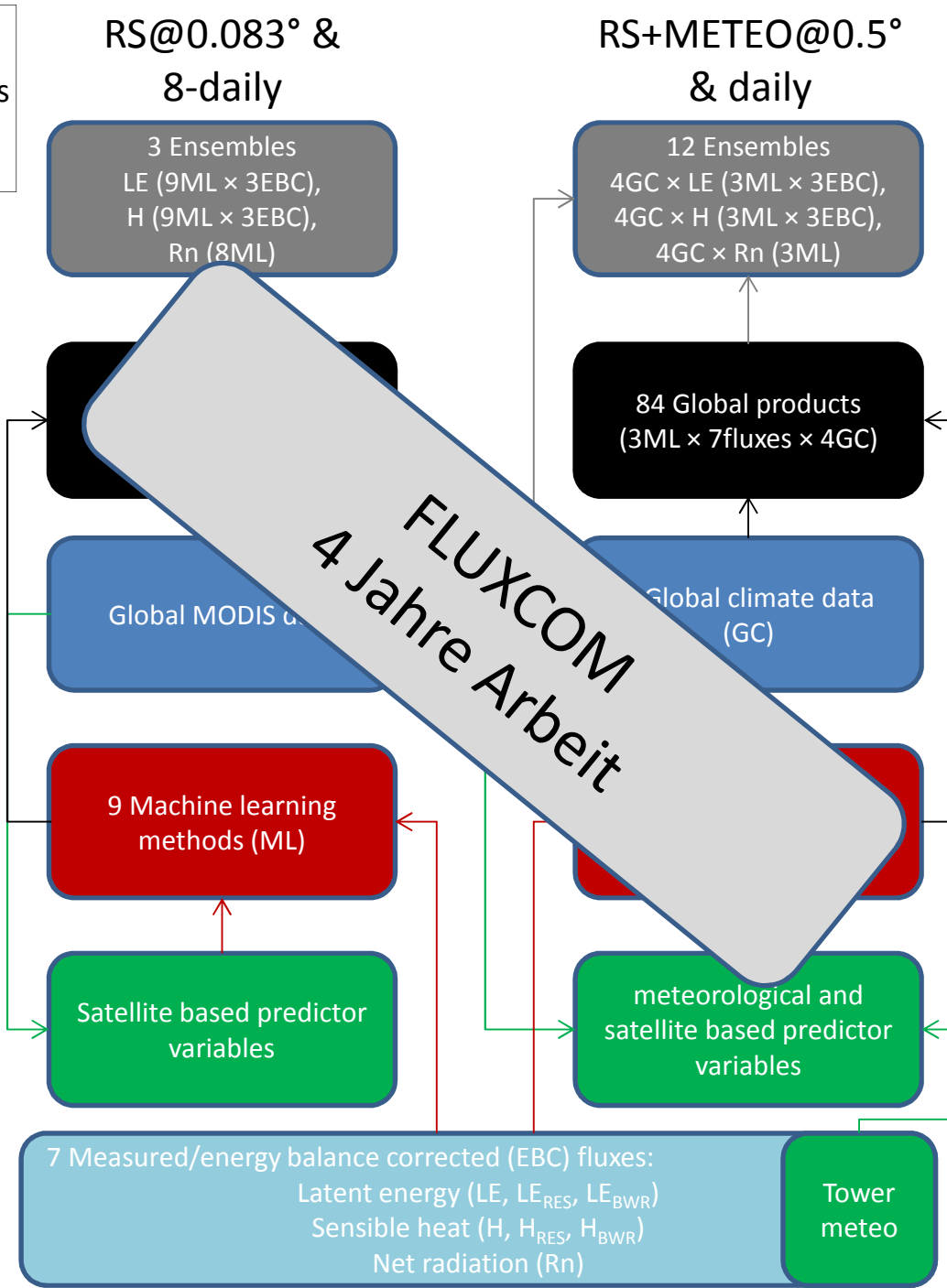
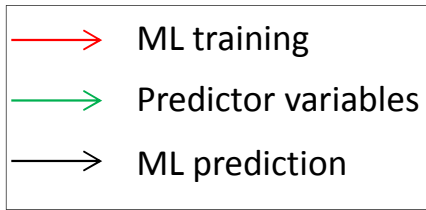
In Entwicklung: code für robuste und effiziente Verarbeitung



Was machen wir mit den Daten?

- Finden
 - Akquirieren
 - Vorprozessieren für Datenstruktur
 - Verwenden
 - neue Datenprodukte → publizieren & verteilen (DOIs via Bibliothek und MPDL)
- } Data Manager





LETTER

doi:10.1038/nature09396

Recent decline in the global land evapotranspiration trend due to limited moisture supply

Martin Jung¹, Markus Reichstein¹, Philippe Ciais², Sonia I. Seneviratne³, Justin Sheffield⁴, Michael L. Goulden⁵, Gordon Bonan⁶, Alessandro Cescatti⁷, Jiquan Chen⁸, Richard de Jeu⁹, A. Johannes Dolman⁹, Werner Eugster¹⁰, Dieter Gerten¹¹, Damiano Gianelle¹², Nadine Gobron¹³, Jens Heinke¹⁴, John Kimball¹⁴, Beverly E. Law¹⁵, Leonardo Montagnani¹⁶, Qiaozhen Mu¹⁷, Brigitte Mueller³, Keith Oleson⁶, Dario Papale¹⁸, Andrew D. Richardson¹⁹, Olivier Roupsard²⁰, Steve Running¹⁷, Enrico Tomelleri¹, Nicolas Viovy², Ulrich Weber¹, Christopher Williams²¹, Eric Wood⁴, Sönke Zaehle¹ & Ke Zhang¹⁴

LETTER

doi:10.1038/nature20780

Compensatory water effects link yearly global land CO₂ sink changes to temperature

Martin Jung¹, Markus Reichstein^{1,2}, Christopher R. Schwalm³, Chris Huntingford⁴, Stephen Sitch⁵, Anders Ahlström^{6,7}, Almut Arneth⁸, Gustau Camps-Valls⁹, Philippe Ciais¹⁰, Pierre Friedlingstein¹¹, Fabian Gans¹, Kazuhito Ichii^{12,13}, Atul K. Jain¹⁴, Etsushi Kato¹⁵, Dario Papale¹⁶, Ben Poulter¹⁷, Botond Raduly^{16,18}, Christian Rödenbeck¹⁹, Gianluca Tramontana¹⁶, Nicolas Viovy¹⁰, Ying-Ping Wang²⁰, Ulrich Weber¹, Sönke Zaehle^{1,2} & Ning Zeng^{21,22}

LETTER

doi:10.1038/nature13731

Global covariation of carbon turnover times with climate in terrestrial ecosystems

Nuno Carvalhais^{1,2}, Matthias Forkel¹, Myroslava Khomik^{1,3}, Jessica Bellarby^{4,5}, Martin Jung¹, Mirco Migliavacca^{1,6}, Mingquan Mu⁷, Sassan Saatchi⁸, Maurizio Santoro⁹, Martin Thurner¹, Ulrich Weber¹, Bernhard Ahrens¹, Christian Beer^{1,10}, Ales James T. Randerson⁷ & Markus Reichstein¹

Terrestrial Gross Carbon Dioxide Uptake: Global Distribution and Covariation with Climate

Christian Beer,^{1*} Markus Reichstein,¹ Enrico Tomelleri,¹ Philippe Ciais,² Martin Jung,¹ Nuno Carvalhais,^{1,3} Christian Rödenbeck,⁴ M. Altaf Arain,⁵ Dennis Baldocchi,⁶ Gordon B. Bonan,⁷ Alberte Bondeau,⁸ Alessandro Cescatti,⁹ Gitta Lasslop,¹ Anders Lindroth,¹⁰ Mark Lomas,¹¹ Sebastiaan Luyssaert,¹² Hank Margolis,¹³ Keith W. Oleson,⁷ Olivier Roupsard,^{14,15} Elmar Veenendaal,¹⁶ Nicolas Viovy,² Christopher Williams,¹⁷ F. Ian Woodward,¹¹ Dario Papale¹⁸

CARBON CYCLE

The dominant role of semi-arid ecosystems in the trend and variability of the land CO₂ sink

Anders Ahlström,^{1,2*} Michael R. Raupach,^{3†} Guy Schurgers,⁴ Benjamin Smith,¹ Almut Arneth,⁵ Martin Jung,⁶ Markus Reichstein,⁶ Josep G. Canadell,⁷ Pierre Friedlingstein,⁸ Atul K. Jain,⁹ Etsushi Kato,¹⁰ Benjamin Poulter,¹¹ Stephen Sitch,¹² Benjamin D. Stocker,^{13,14} Nicolas Viovy,¹⁵ Ying Ping Wang,¹⁶ Andy Willshire,¹⁷ Sönke Zaehle,⁶ Ning Zeng¹⁸



click for updates

PMAS PLUS

Global and time-resolved monitoring of crop photosynthesis with chlorophyll fluorescence

Luis Guanter^{a,1,2}, Yongguang Zhang^{a,1}, Martin Jung^b, Joanna Joiner^c, Maximilian Voigt^d, Joseph A. Berry^d, Christian Frankenberg^e, Alfredo R. Huete^f, Pablo Zarco-Tejada^g, Jung-Eun Lee^h, M. Susan Moranⁱ, Guillermo Ponce-Campos^j, Christian Beer¹, Gustavo Camps-Valls^k, Nina Buchmann^l, Damiano Gianelle^m, Katja Klumppⁿ, Alessandro Cescatti^o, John M. Baker^p, and Timothy J. Griffis^q

Probleme und neue Konzepte

Erstellung, Dokumentation, Publikation und Teilen von Prozessierungs- und Analyseketten (workflows)

- Für Transparenz und Reproduzierbarkeit
- Für Wiederverwendbarkeit
- Für effektives und effizientes Arbeiten („Wissenschaft“ vs Softwareentwicklung und Datamanagement)

Warum wird oft nicht der code eines Projektes publiziert?

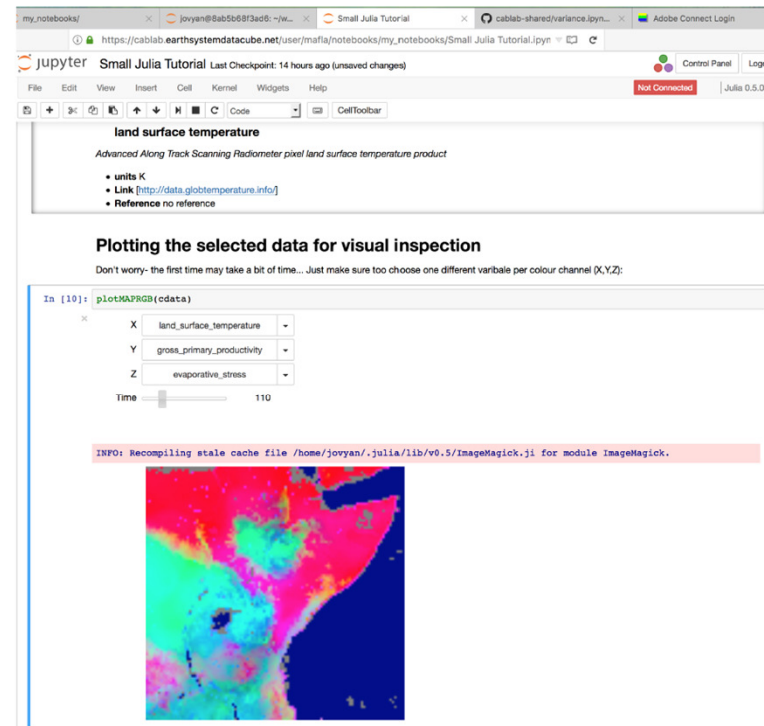
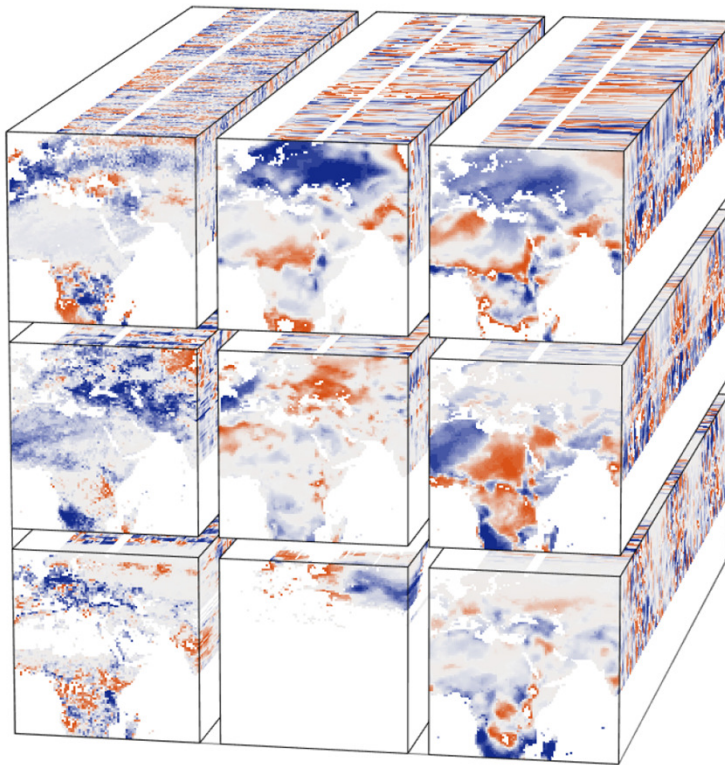
- Oft schwer lesbar und wiederverwendbar – großer Anteil betrifft performantes Verarbeiten und nicht „Wissenschaft“ (z.B. Parallelisierung, Kopplung an Hardware)
- Oft schlecht dokumentiert

Earth System Data Cube

Harmonisierte Daten



Analysewerkzeug



<http://earthsystemdatacube.net/>

Das data cube Analysewerkzeug

- Viele Standardoperationen sind effizient implementiert (high-performance computing läuft im Hintergrund)
- Neue Funktionen können dazu gefügt werden
- Workflows können einfach dokumentiert und geteilt werden (z.B. via Jupiter notebooks)
- Julia, Python, R interfaces

```
In [4]: @time cube_filled = map(gapFillMSC,cdata,46,max_cache=1e7);
```

```
18.017533 seconds (14.40 M allocations: 877.965 MB, 1.70% gc time)
```

```
In [5]: @time cube_anomalies = map(removeMSC,cube_filled,46,max_cache=1e7);
```

```
3.923924 seconds (5.05 M allocations: 198.965 MB, 3.20% gc time)
```

```
In [6]: @time cube_normalized = map(normalize,cube_anomalies,max_cache=1e7);
```

```
3.708907 seconds (2.90 M allocations: 154.698 MB, 1.22% gc time)
```

```
In [7]: @time scores = map(recurrences,cube_normalized,5.0,5,zeros(Float32,506,506),max_cache=1e7);
```

```
61.147698 seconds (47.40 M allocations: 1.997 GB, 0.92% gc time)
```

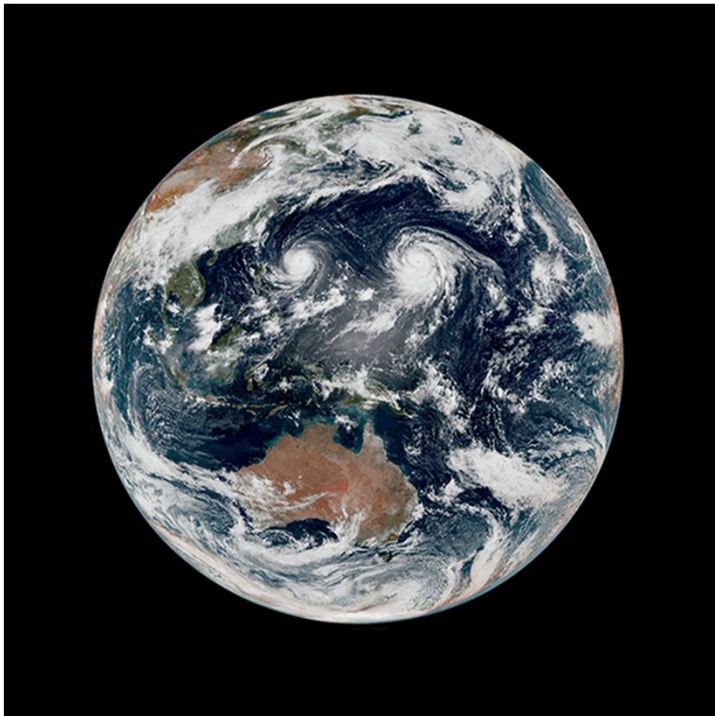

Probleme und neue Konzepte



Too Big Data

Mehrere hundert Terrabytes an Erdsystemdaten werden täglich von Satelliten aufgezeichnet

z.B. Himawari-8 produziert ca.
350 GB/Tag



Hohe zeitliche Auflösung (Minuten)



Hohe räumliche Auflösung (Meter)

Cloud computing – aber wie und wo?

Klassische Ansatz: Datenakquise → Vorprozessierung → Analyse ... funktioniert nicht mehr (bzw. sehr uneffektiv)

„Neue“ Ansatz: Cloud computing → Implementierung von workflows in big data clouds (z.B. Google Earth Engine, Amazon ...)

- enthalten meist nicht alle benötigten (großen) Daten
- keine Garantie für Reproduzierbarkeit (Daten „veralten“ und werden nicht für immer vorgehalten)
- erfordert neues „know-how“
- kostenintensiv (?)
- IT Sicherheit (?)

**"Sollen wir alles auf
Google-Earth Engine
rechnen, oder gibt es eine
MPG Lösung..?"**



Vielen Dank für Ihre Aufmerksamkeit!

mjung

mmahecha

mreichstein

@bgc-jena.mpg.de