

Anforderungen an und Potenzial von eScience/GRID aus der Sicht der Naturwissenschaften

Wolfgang Voges
Max-Planck-Institut für extraterrestrische Physik

DINI e.V. Jahrestagung
10. Oktober 2007
Oldenburg

Dank an folgende Personen, die mir Material zur Verfügung gestellt haben bzw. deren OpenAccess-Präsentationen ich z.T. benutzt habe

Art Carlson (MPE)

Volker Springel (MPA)

Matthias Steinmetz (AIP, Astrogrid-D)

Hans-Falk Hoffmann (CERN)

Wilhelm Hasselbring (U-Oldenburg, WISENT-Grid)

Stephan Kindermann (DKRZ,C3-Grid)

Wolfgang Walkowiak (U-Siegen, HEP-Grid)

Gliederung

- Einführung in die Problematik
- Setting the scene: Beispiele aus den Naturwissenschaften (Astronomie, Hochenergie-Teilchenphysik, Klimaforschung, Wetterbeobachtungen, Geophysik)
- Open Access
- Zusammenfassung

Einführung in die Problematik

Information und Kommunikation
sind das
A und O
der Wissenschaft

Information ist gleichzusetzen mit:

Daten (Beobachtungen, Simulationen, Laborversuche, Bilder, Audio u. Videoaufzeichnungen, Publikationen, Zeitschriften, Bücher, „graue Literatur“, etc.)

Kommunikation wird ermöglicht durch:

Netzwerke (GRID), das Web, Portale... -> e-Infrastruktur

GRID bedeutet hier,

der Verbund von Hochleistungsrechnern (aber nicht nur!!), d.h. hohe CPU-Leistung, grosse Speicherplatz-Kapazitäten, schnelle Datenleitungen

>>>> e-Science

e-Science, Cyber-Infrastructure

- “e-Science is about *more* than networks, GRIDs, High Performance Computing...
- e-science is about global collaboration in key areas of science
- and the next generation of infrastructure that will enable it”
- John Taylor, Director Research Councils, UK, 2000
- NSF Cyber-infrastructure Initiative; “Blue Ribbon Panel on Cyber-infrastructure” 2002
- “e-Infrastrukturen Initiative” EU, FP 6,7: “fortgeschrittenste Wissensgesellschaft bis 2010” (Lissabon 2000)
- D-Grid – Neue Formen der Zusammenarbeit (BMBF) “VO”

Auf der German e-Science Conference 2007 in Baden-Baden wurden vier Bereiche des e-Science behandelt:

- Grid-Computing
- Wissensvernetzung
- e-Learning
- Open Access

Setting the scene:

Beispiele aus den Naturwissenschaften

(Astronomie, Hochenergie-Teilchenphysik,

Klimaforschung, Wetterbeobachtungen,

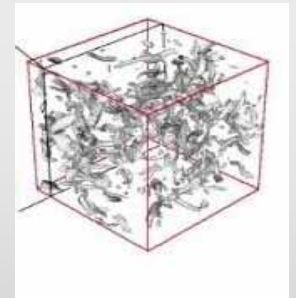
Geophysik)

Evolving Science

- Thousand years ago:
science was empirical
describing natural phenomena
- Last few hundred years:
theoretical branch
using models, generalizations
- Last few decades:
a computational branch
simulating complex phenomena
- Today:
data exploration (eScience)
synthesizing theory, experiment and computation with advanced data management, data mining, data visualisation, and statistics

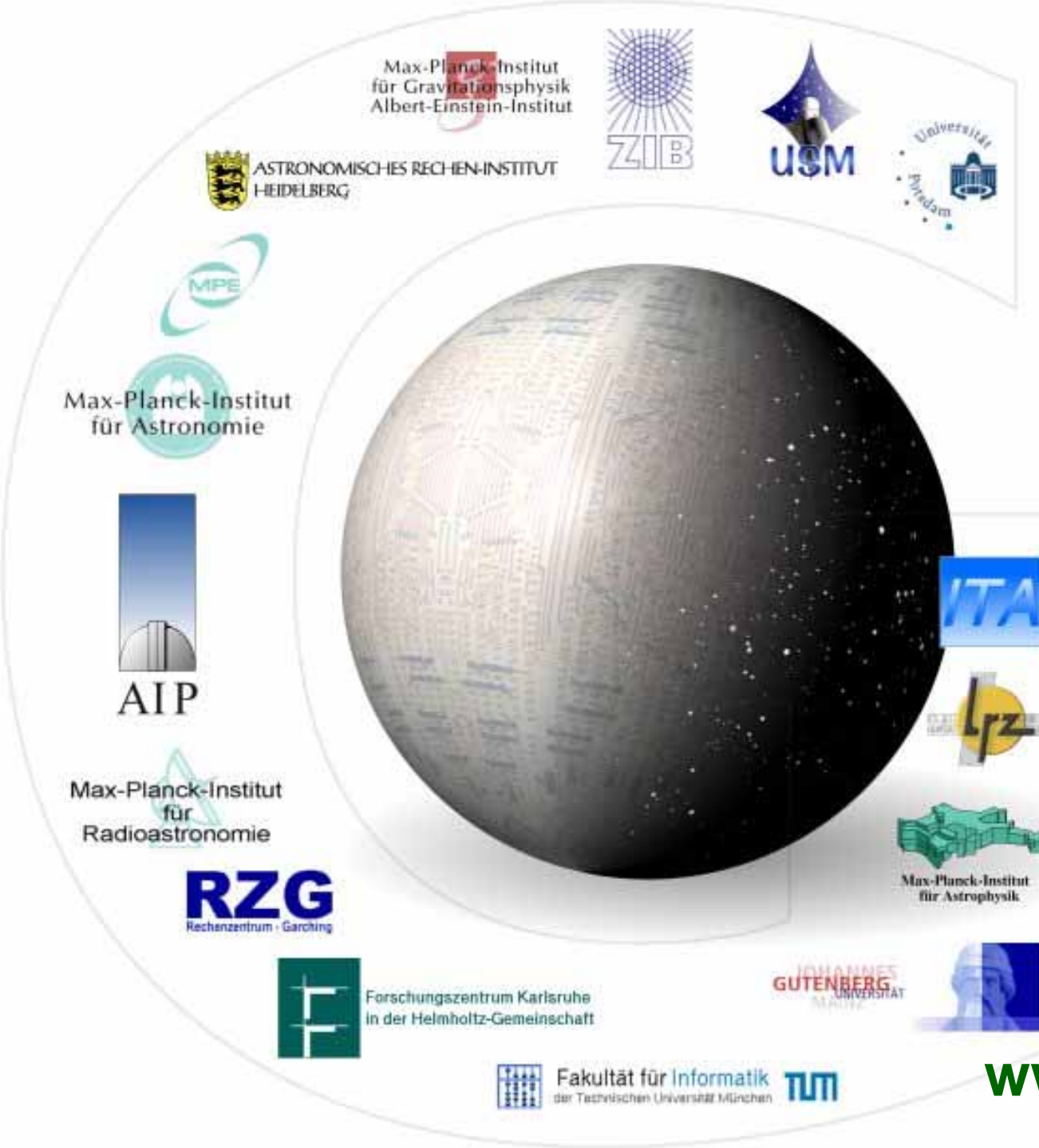


$$\left| \frac{a}{a} \right|^2 = \frac{4\pi G \rho}{3} - K \frac{c^2}{a^2}$$



German Astronomical Community GRID

AstroGrid-D



www.gac-grid.org

International Virtual Observatory Alliance

www.ivoa.org



www.g-vo.org

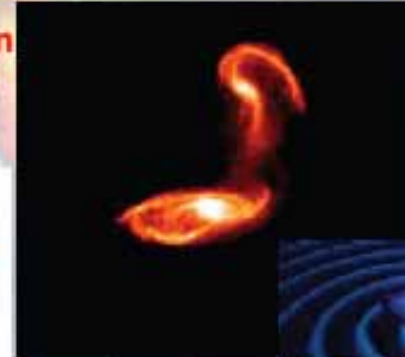


AstroGrid-D: Aufgaben

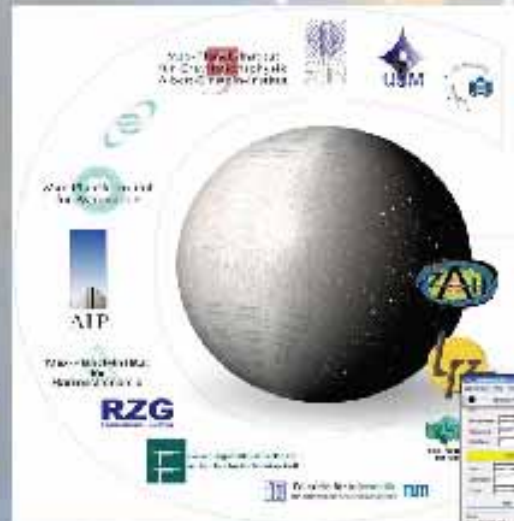


sanssouci.aip.de

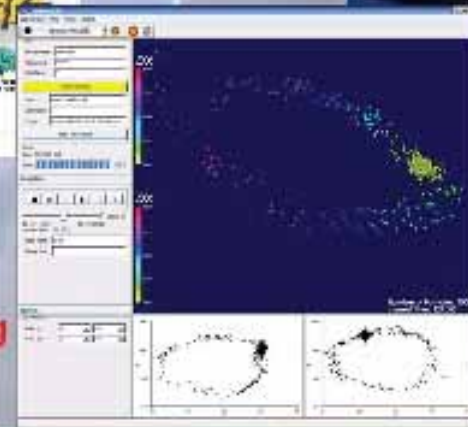
Galaxien-Kollision
Postprocessing



LISA



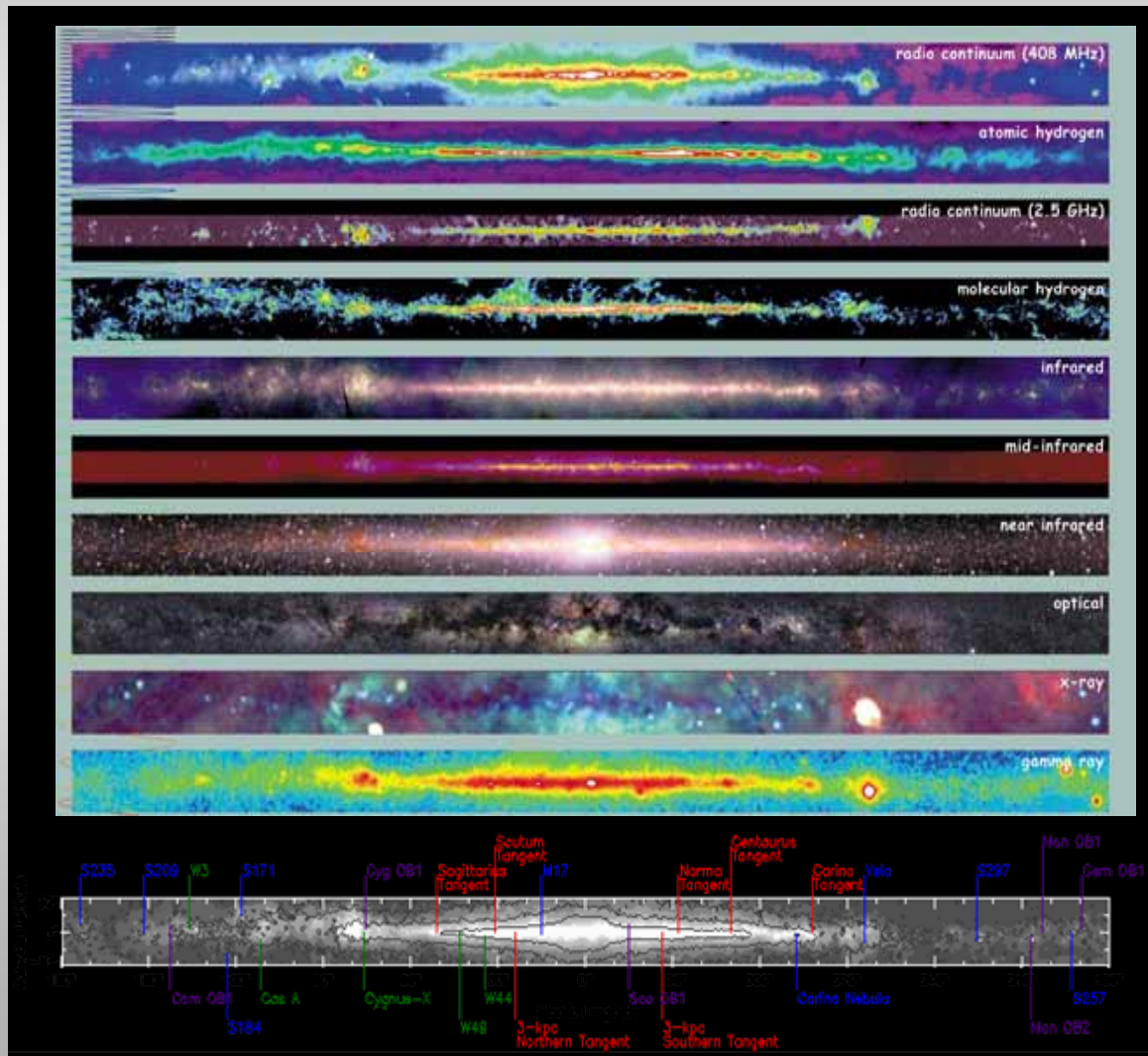
X-NBODY6++
Visualisierung



LOFAR

Combining images: a multi-wavelength view of the galactic plane

Zum besseren Verständnis der physikalischen Strahlungsmechanismen in den beobachteten Himmelsobjekten oder Regionen



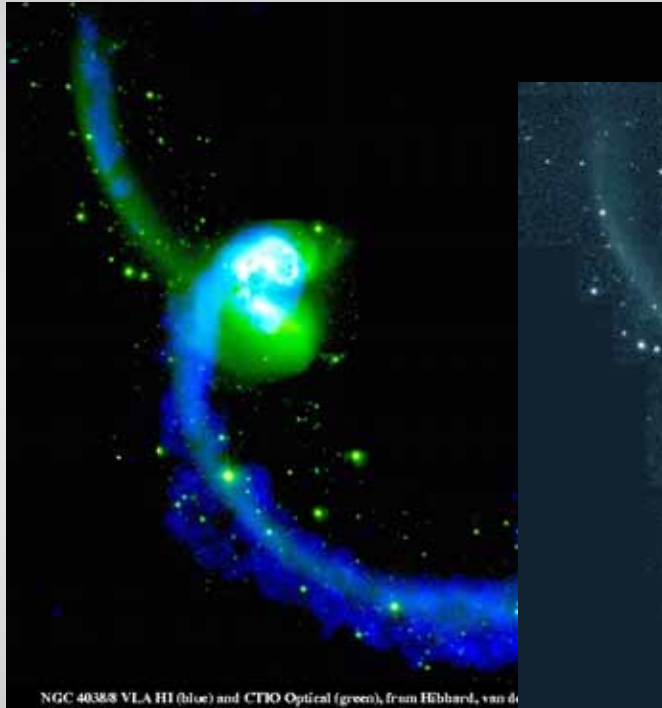
Dame, Hartmann, Thaddeus
2001, ApJ, 547,792

Price et al.
,2001,Astron.J.,121,2819

Snowden et al., 1997, ApJ,
485,125

Combining images: a multi-wavelength view of a galaxy merger (NGC 4038/39)

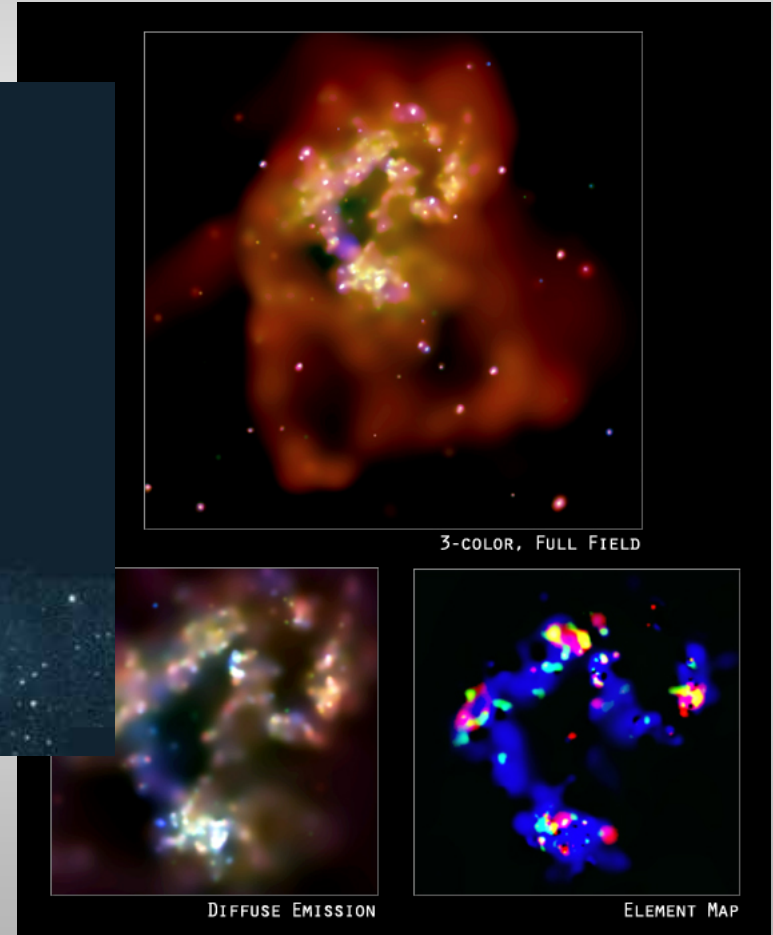
Radio



Optical



X-Ray

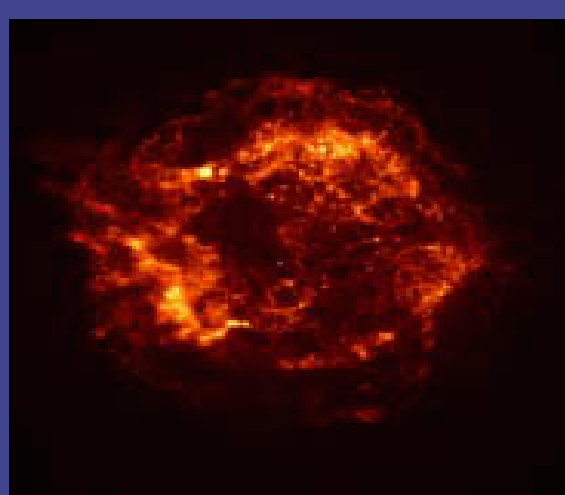


John Hibbard

<http://www.cv.nrao.edu/~jhibbard/n4038/n4038.html>

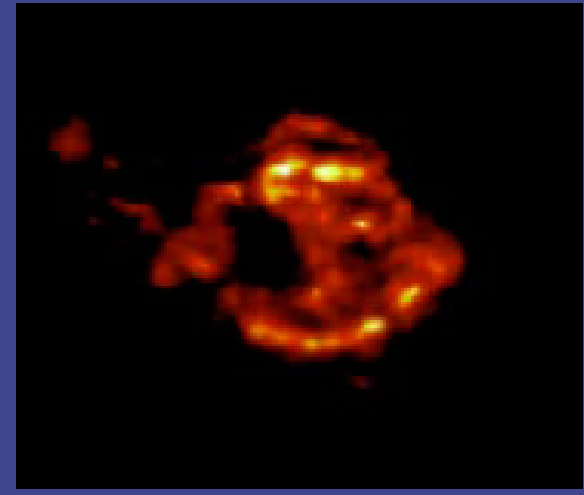
NASA/CXC/SAO/G. Fabbiano et al.

Virtual “all λ ” Observatory



Shocks seen in the X-ray

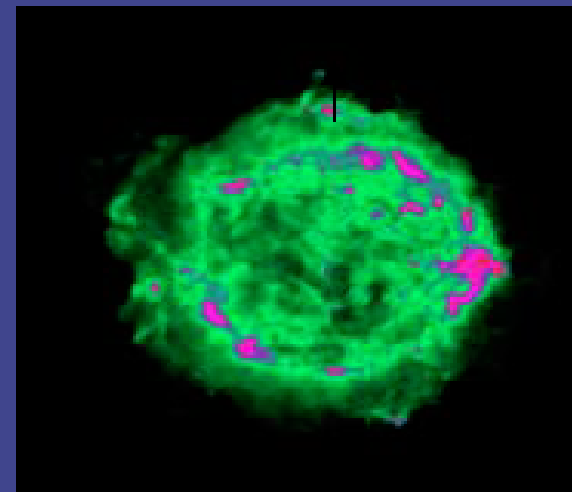
multi- λ
views of a
Supernova
Remnant



Dust seen in the IR



Heavy elements
seen in the optical



Relativistic electrons
seen in the radio

What can an OWL see ?



0.6 arcsec

OWL

GRID Simulations

‘**Millennium Run**’ of the Virgo consortium, parallel on 512 processors of the IBM p690 at RZG-Garching for 1 month

Simulation Set-up:

- Number of particles: $2160^3 = 10.077.696.000 = \sim 10^{10}$ particles
- Size of box: $L = 500 h^{-1} \text{ Mpc}$
- FFT: $2560^3 = 16.777.216.000 = \sim 17$ billion cells

Particle data
44 bytes / particle

Tree memory
40 bytes / particle

FFT workspace
24 bytes / mesh cell

Not needed concurrently!

Minimal RAM requirement

~840 GByte

Special code version
L-GADGET-II needs:

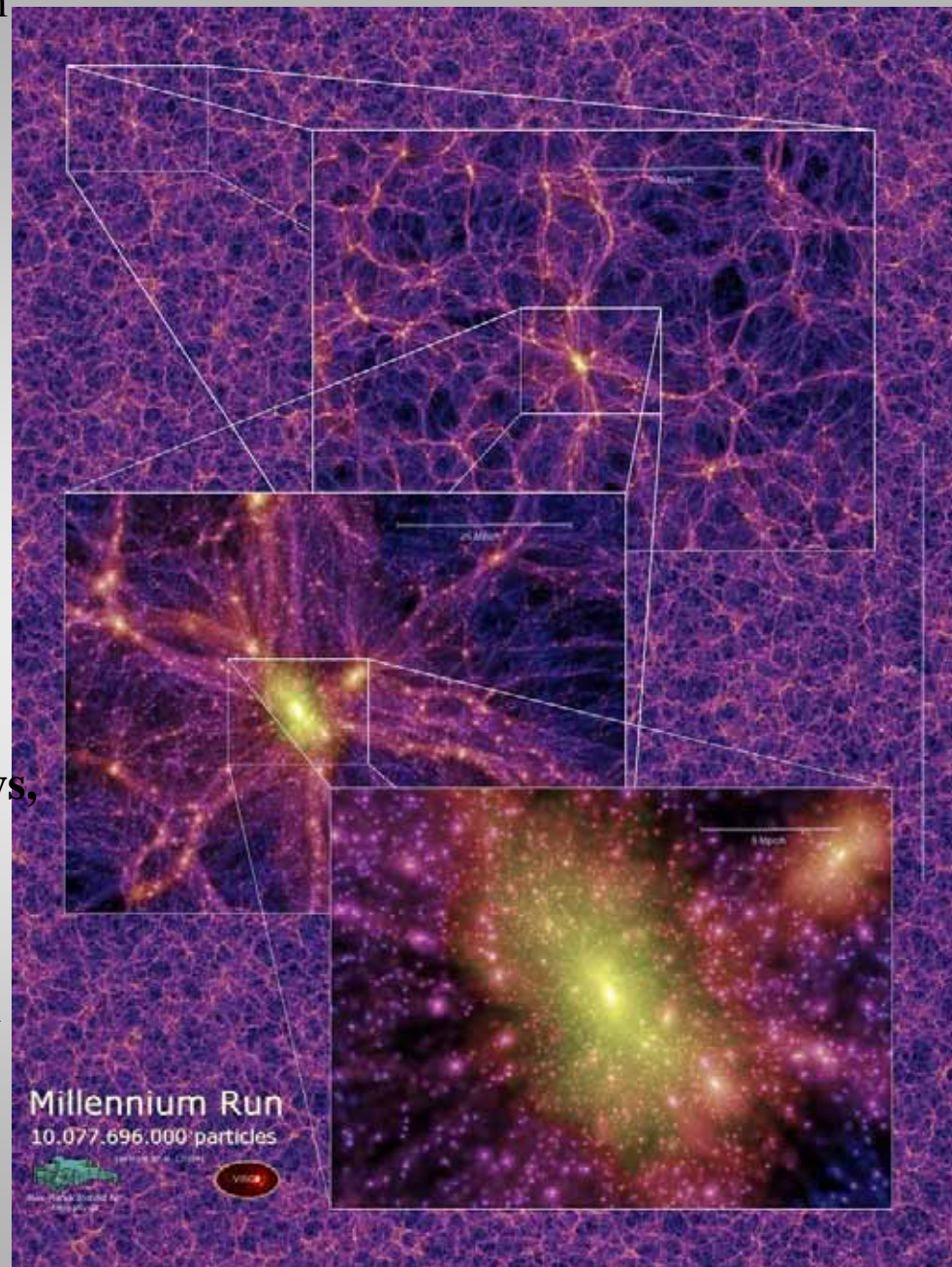
84 bytes / particle

(Assuming 1.5 mesh cells/particle)

38 processor years, ~30 Tbyte of data

The Millennium Run used more than 10 billion particles to trace the **evolution of the matter distribution** in a cubic region of the Universe over 2 billion light-years on a side. It kept the supercomputer at the Max Planck Society's Computer centre in Garching busy for **more than a month**. By applying sophisticated modeling techniques to the **25 Tbytes** of stored output, Virgo scientists have been able to create evolutionary histories both for the 20 million or so galaxies which populate this enormous volume and for the super-massive black holes which occasionally power quasars at their hearts. **By comparing such simulated data to large observational surveys, one can clarify the physical processes** underlying the buildup of real galaxies and black holes.

The illustration shows a projected density field for a 15Mpc/h thick slice of the redshift = 0 output, a massive cluster of galaxies. The overlaid panels zoom in by factors of 4 in each case, enlarging the regions indicated by the white rectangles.



New Science with the GRID



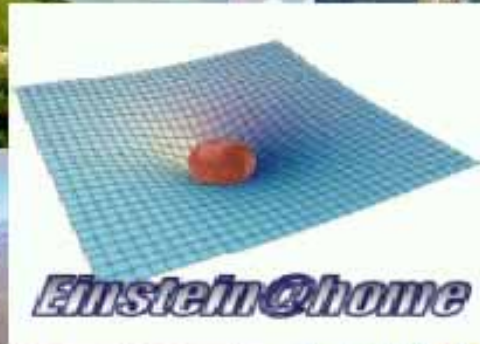
- Radioastronomy at $\lambda = 30\text{m}$
- Telescope size: 1000 km
- Solution: 10,000 small antennas distributed over $\text{Ø}1000\text{ km}$
- Propagation time correlations
- All-sky observation

➤ Available supercomputers, fast networks and intelligent software define the power of the telescope



AstroGrid-D: UseCase GEO600 / LIGO

Laser I nterferometer G ravitational Wave O bservatory





Gravitationswellen: Datenanalyse im Grid

- Die Datenauswertung kann in kleinen Datensets erfolgen, daher ist diese Auswertung “embarrassingly parallel”.
- **Einstein@Home** ist, ähnlich wie SETI@Home, auch auf einer Workstation zur Nutzung der Idle-Time geeignet.
- **Einstein@Home** ist somit auch eine ideal einfache Grid-Applikation, die viele Betriebssysteme unterstützt.
- Checkpointing und Recovery Verfahren erlauben eine sehr genaue Steuerung von CPU-Requirements und walltime-Vorgaben.
- Durch automatisches Deployment der Software, von Jobsubmission und Jobmanagement kann ein sehr gutes Skalieren der Grid-Applikation erzielt werden.
- 30000 CPU h pro Woche werden nach interner Statistik abgearbeitet

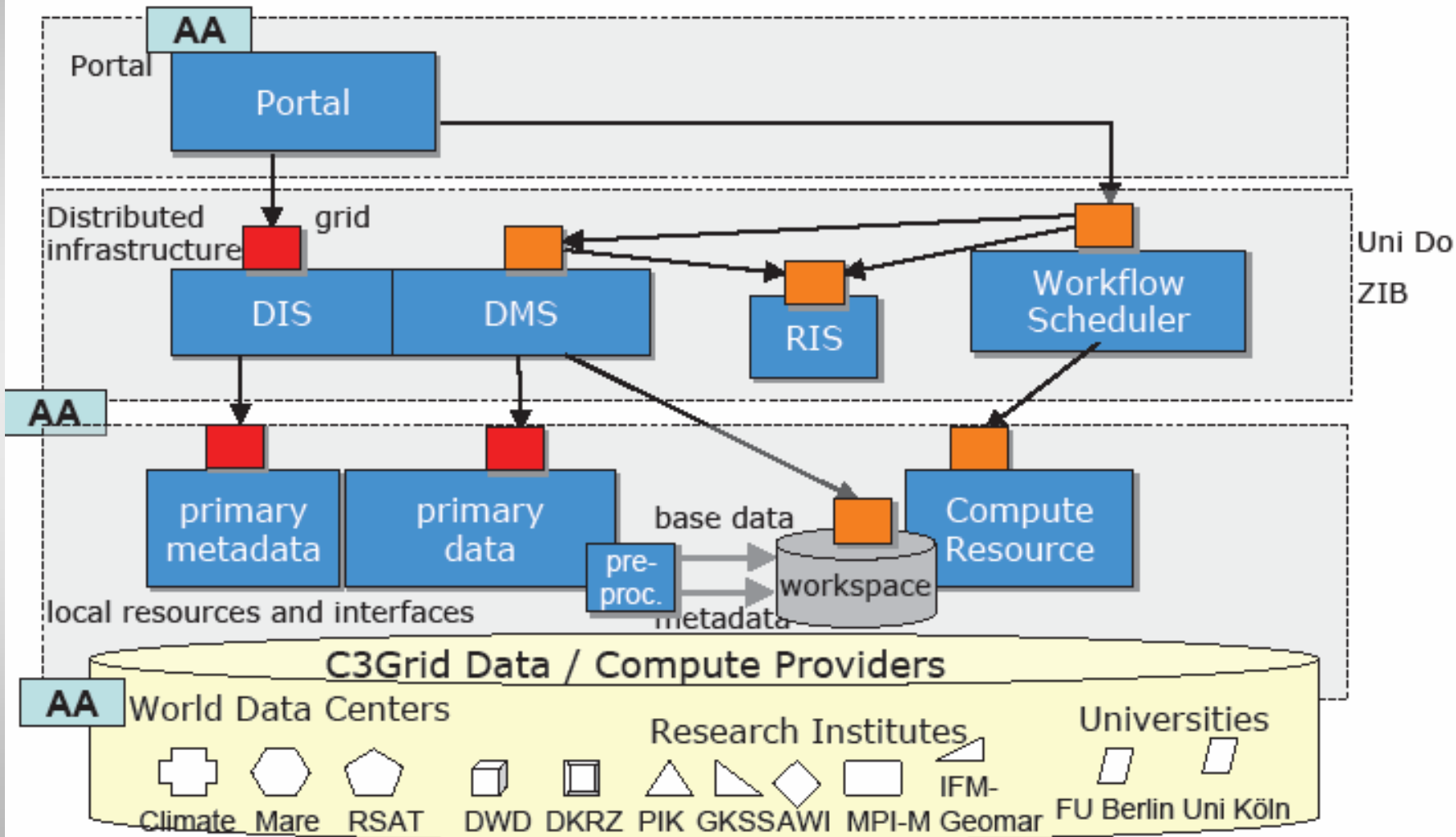
Collaborative Climate Community Data and Processing Grid - C3Grid -

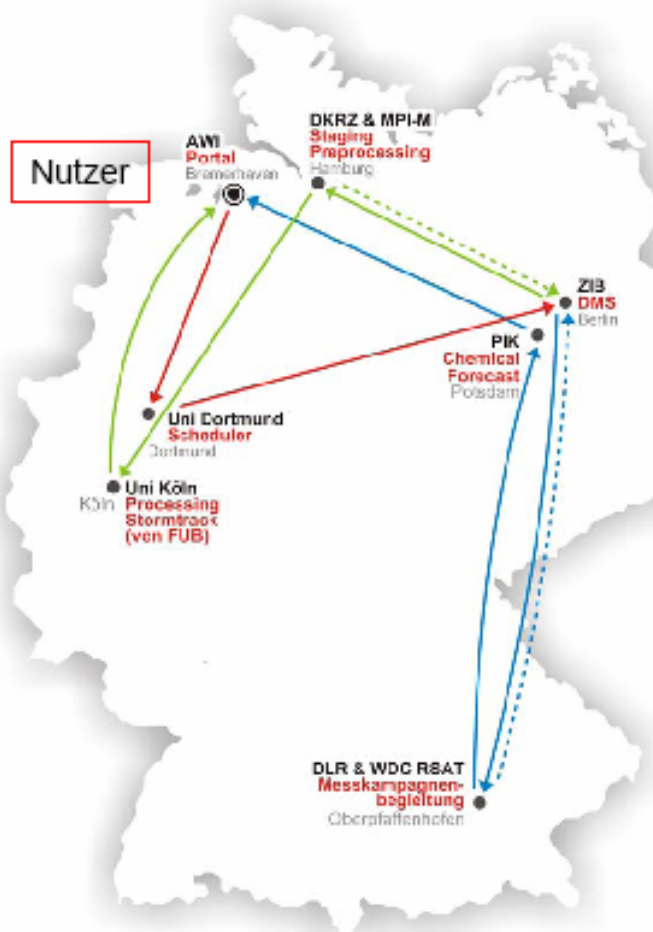
Stephan Kindermann, DKRZ
For C3-Team

Model Output Data + Observation Data + Analysis Data

Needed: A collaborative environment
for **data access, data exchange**
and **data processing**







1. Login at C3 Portal (→)
2. Query for datasets (DIS) (→)
3. Submitting job (Scheduler) (→)
4. DMS (→)
5. Data staging and pre-processing (→, →)
6. GridFTP from data provider to compute provider (→, →)
7. Analysis job execution (→, →)
8. Publication of results (→, →)



GEFÖRDERT VOM



Bundesministerium
für Bildung
und Forschung

WISENT Wissensnetz Energiemeteorologie

Prof. Dr. W. Hasselbring

OFFIS

Institut für Informatik

FuE-Bereich Betriebliches Informationsmanagement (BI)

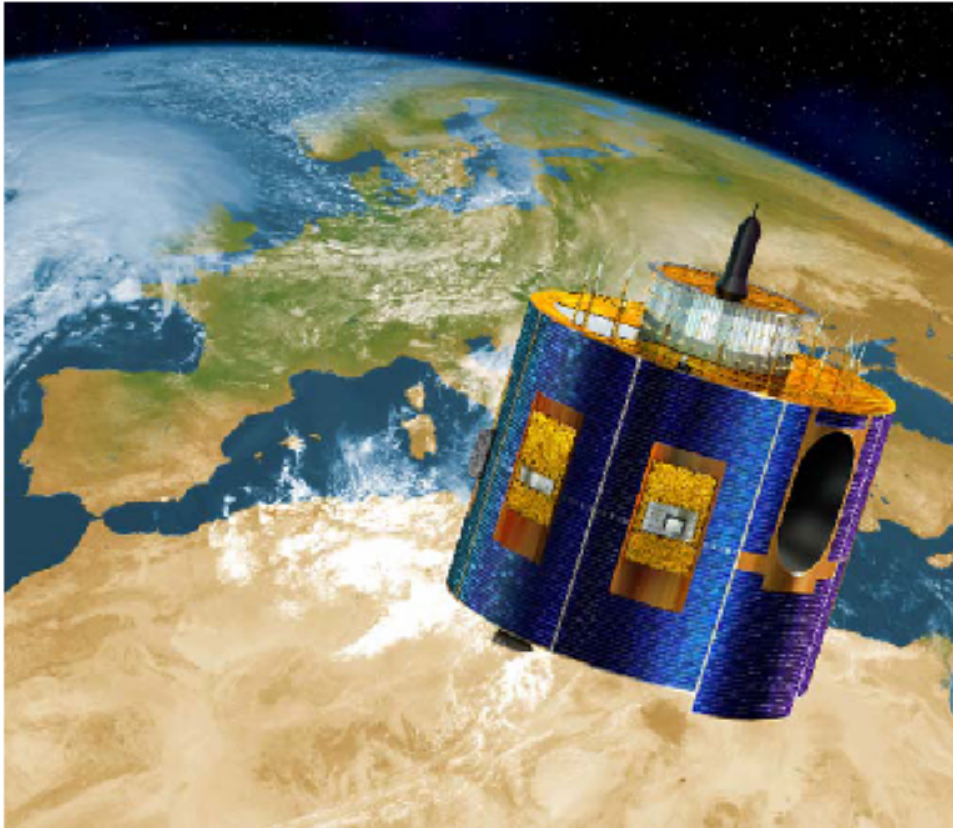
WISENT 

WISSENSNETZ
ENERGIEMETEOROLOGIE

D-Grid AHM, 10. September 2007
W. Hasselbring
OFFIS

Seite 1

Energiemeteorologie



- Gewinnung von Informationen zur Charakterisierung der fluktuierenden Energieerzeugung aus Solar- und Windenergie
- Erkenntnisgewinn durch interdisziplinäre Kooperation
 - Klima und Wetter (Meteorologie)
 - Umwandlung von Wind und Sonnenenergie (Physik)
 - Leistungsfähige, flexible verteilte Systeme (Informatik)

Wettersatellit MSG-1: Die Forscher der Universität Oldenburg und des DLR können damit alle 15 Minuten flächendeckend und präzise den verfügbaren "Brennstoff" für Solaranlagen bestimmen.

Quelle: EUMETSAT

Daten zur Optimierung der fluktuierenden Energieerzeugung aus Solar- und Windenergie

Herausforderungen:

- Sehr große Datenmengen (viele Terabyte)
- Komplexe Prozessketten (teilweise in Echtzeit)
- Sicherung der zukünftigen Energieversorgung



Photovoltaik-Anlage in Berlin-Brandenburg
Quelle: WZL/UBA



Windkraftanlagen in Ostpreußen
Quelle: WZL/UBA

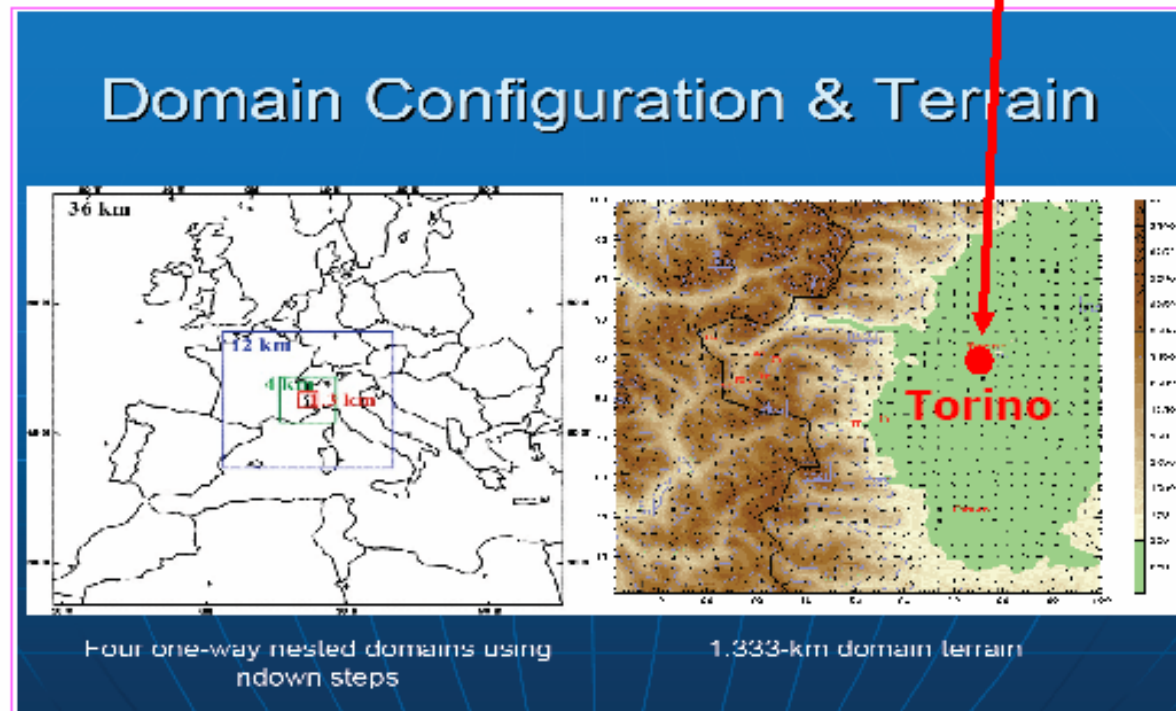
Eine der Anwendungen:

Weather **Research & Forecasting** Model (WRF)

- Wird sowohl zur meteorologischen Grundlagenforschung als auch zur operativen Wettervorhersage genutzt.
- Nutzung in WISENT zur mesoskaligen Wettervorhersage.

Nutzung für die olympischen Winterspiele 2006, hier zur Illustration der Möglichkeiten.

24-Stunden Vorhersage benötigte 4 h auf einem Parallelrechner mit 42 Prozessoren.



HEP CG

High Energy Physics Community Grid

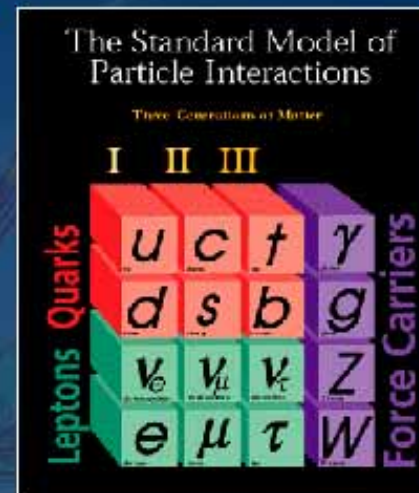
Wolfgang Walkowiak
Universität Siegen

D-Grid All-Hands-Meeting
Göttingen
10. September 2007



Grundlegende Fragen der Teilchenphysik

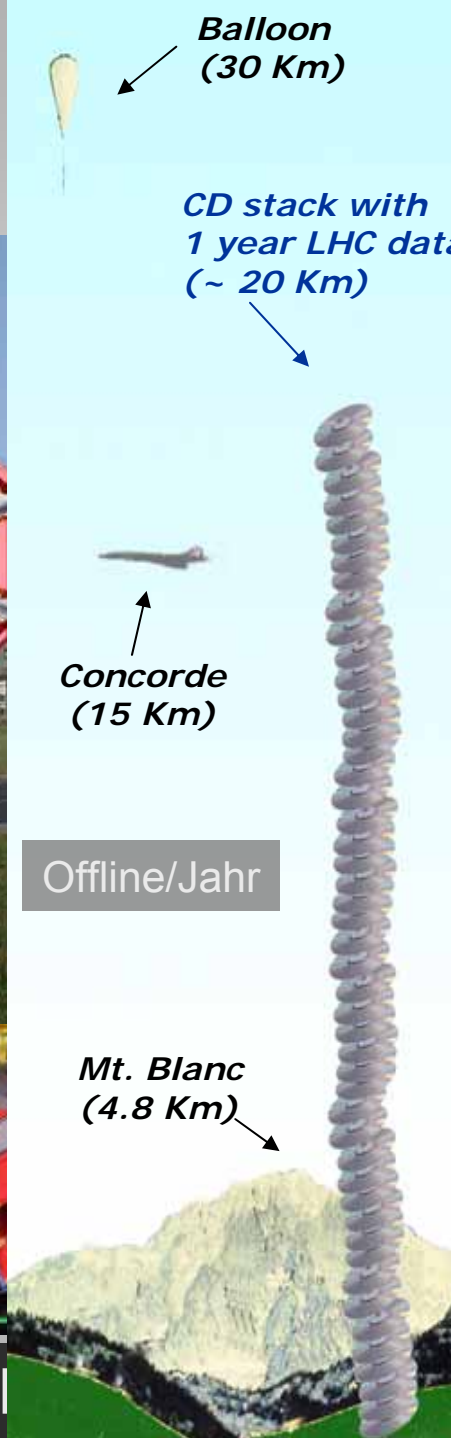
- Woraus besteht Materie?
- Welche Kräfte wirken zwischen Elementarteilchen?
- Woher kommt Materie?
- Worin besteht der Ursprung der Masse?



- Riesige Beschleuniger an Forschungszentren (CERN, ...)
- ~ 10 000 Wissenschaftler in grossen, internationalen Kollaborationen.



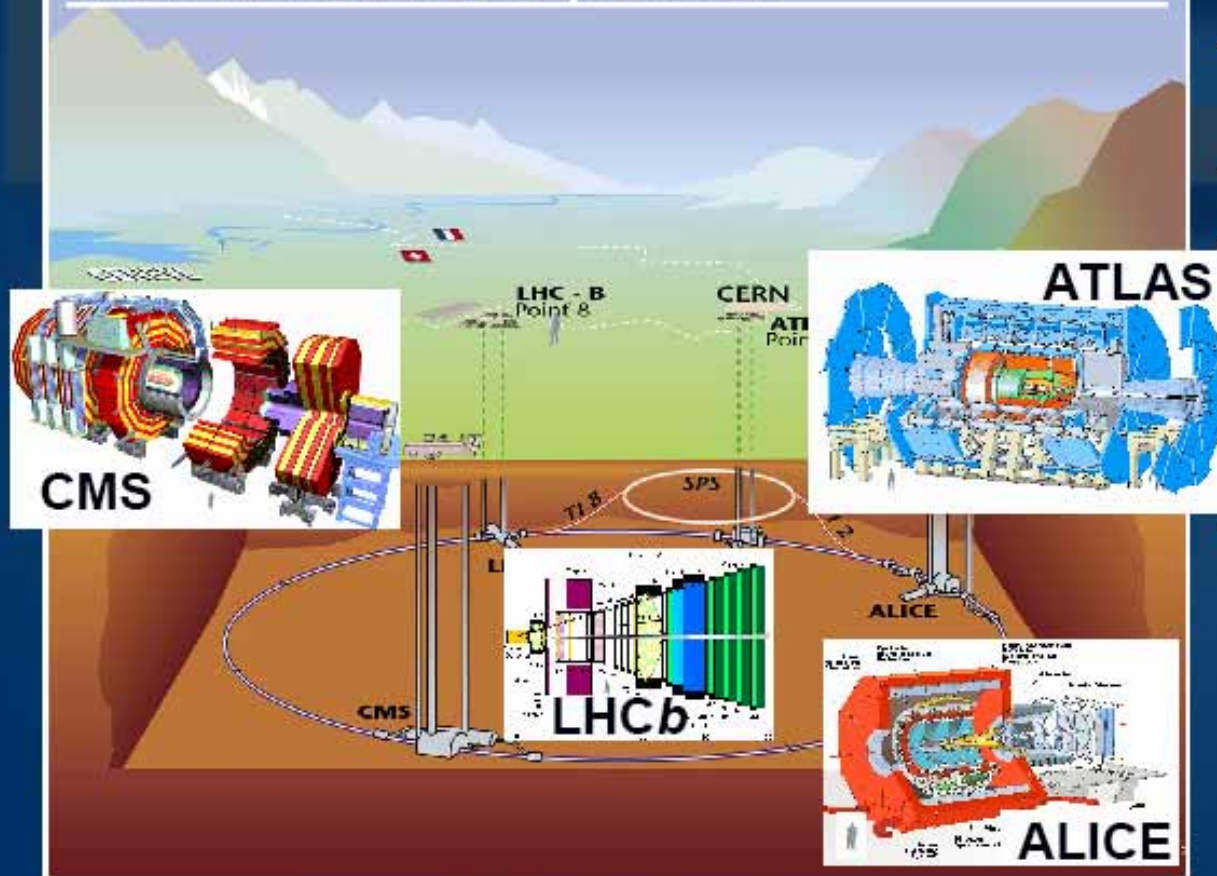
Teilchenphysik “made by CERN, ... Wuppertal,



Aber: 1 Petabyte/s Informationsfluss in die LHC-

Large Hadron Collider startet nächstes Jahr!

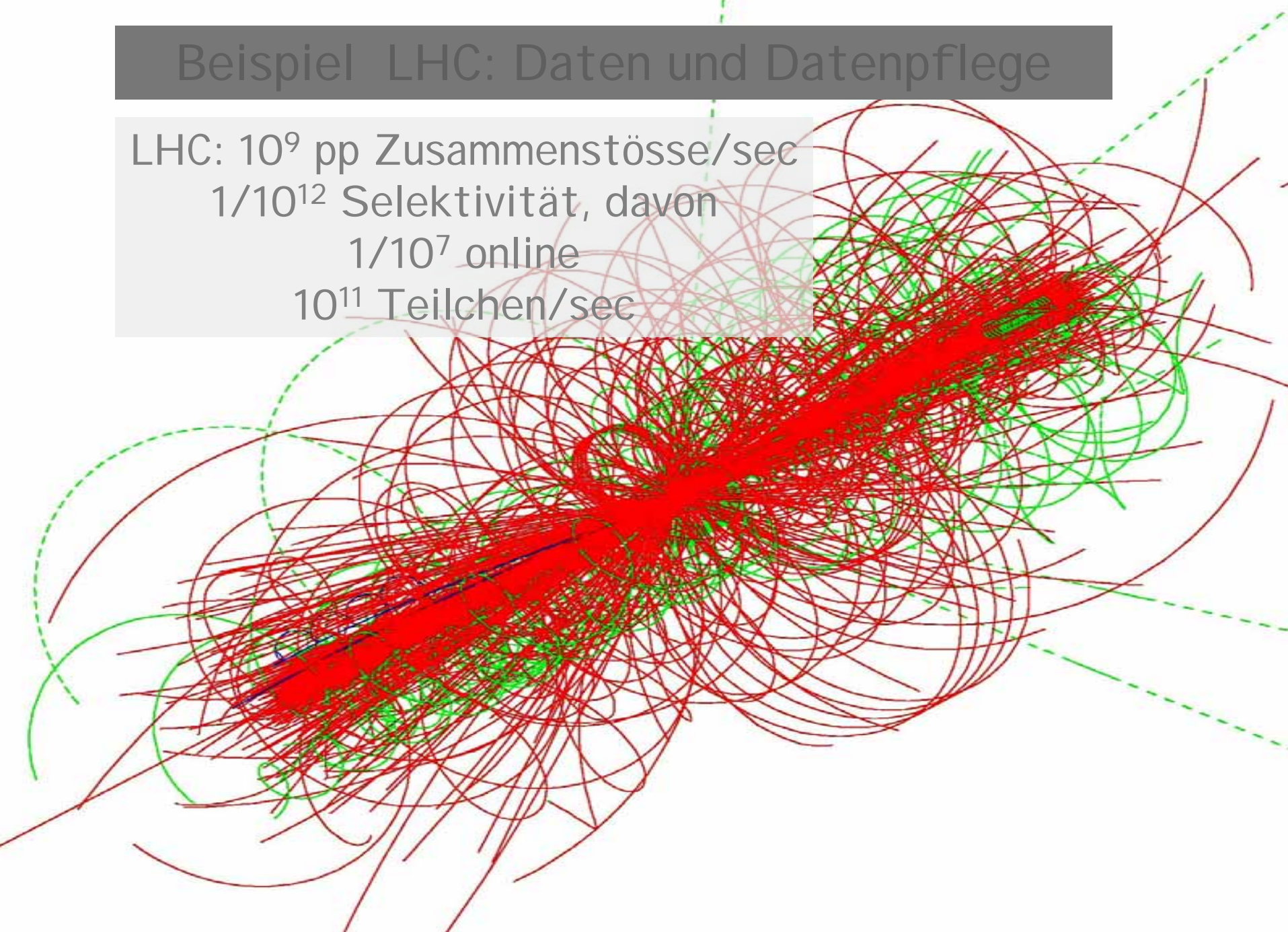
Overall view of the LHC experiments.



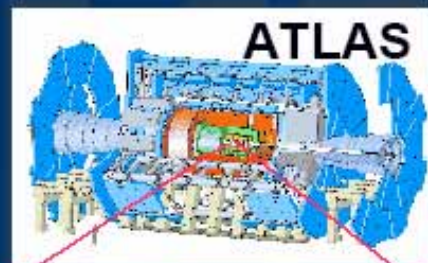
- LHC-Umfang: 27.6 km
- Proton-Proton
- 4 Großexperimente mit ~ 2 500 Physikern/ Experiment
- Entwicklung und Bau über > 10 Jahre

Beispiel LHC: Daten und Datenpflege

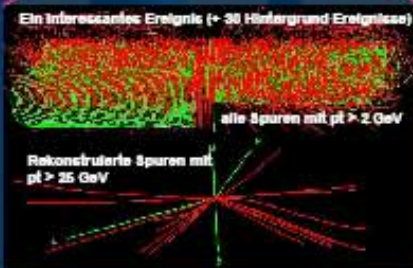
LHC: 10^9 pp Zusammenstöße/sec
 $1/10^{12}$ Selektivität, davon
 $1/10^7$ online
 10^{11} Teilchen/sec



Die Herausforderung



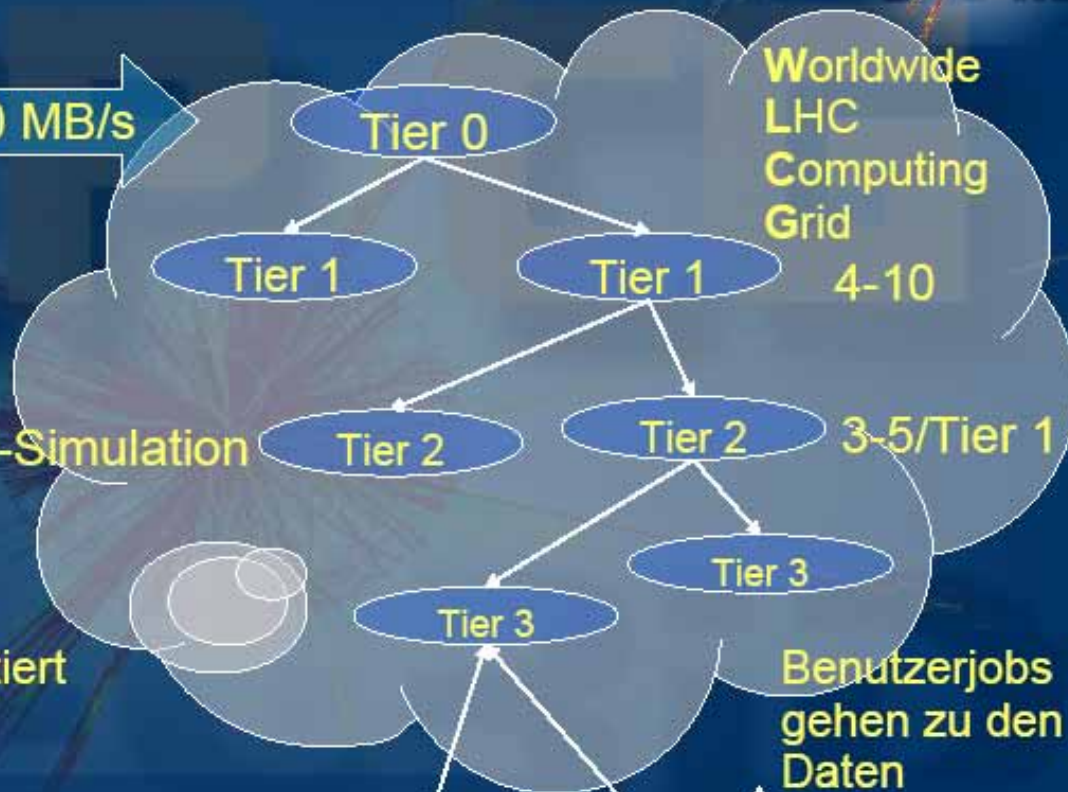
320 MB/s



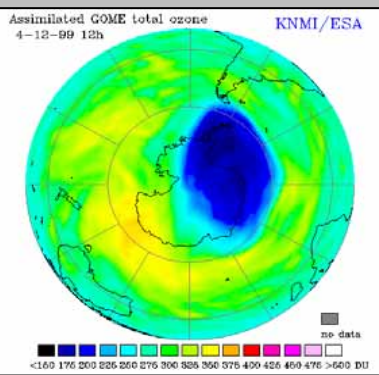
Kollisionrate 40 MHz
~ 200 Ereignisse/s selektiert
eine 2 GB Datei in 5 s

Pro Experiment:
~10 PetaByte / Jahr
~ 10^9 Ereignisse / Jahr
~ 10^3 Batch- und interaktive Nutzer

MC-Simulation



ENVISAT is providing very large data flows...



Altitude 0 to 100 km: GOMOS, MIPAS and SCIAMACHY are building a three-dimensional profile of ozone concentrations in the atmosphere.

Altitude 0 to 20 km: MIPAS and SCIAMACHY are detecting low levels of gases from industry, power generation and agriculture.

Altitude 0 to 10 km: MERIS obtains an image in which the clouds you see are but a part of a complex map of the concentration of water vapour.

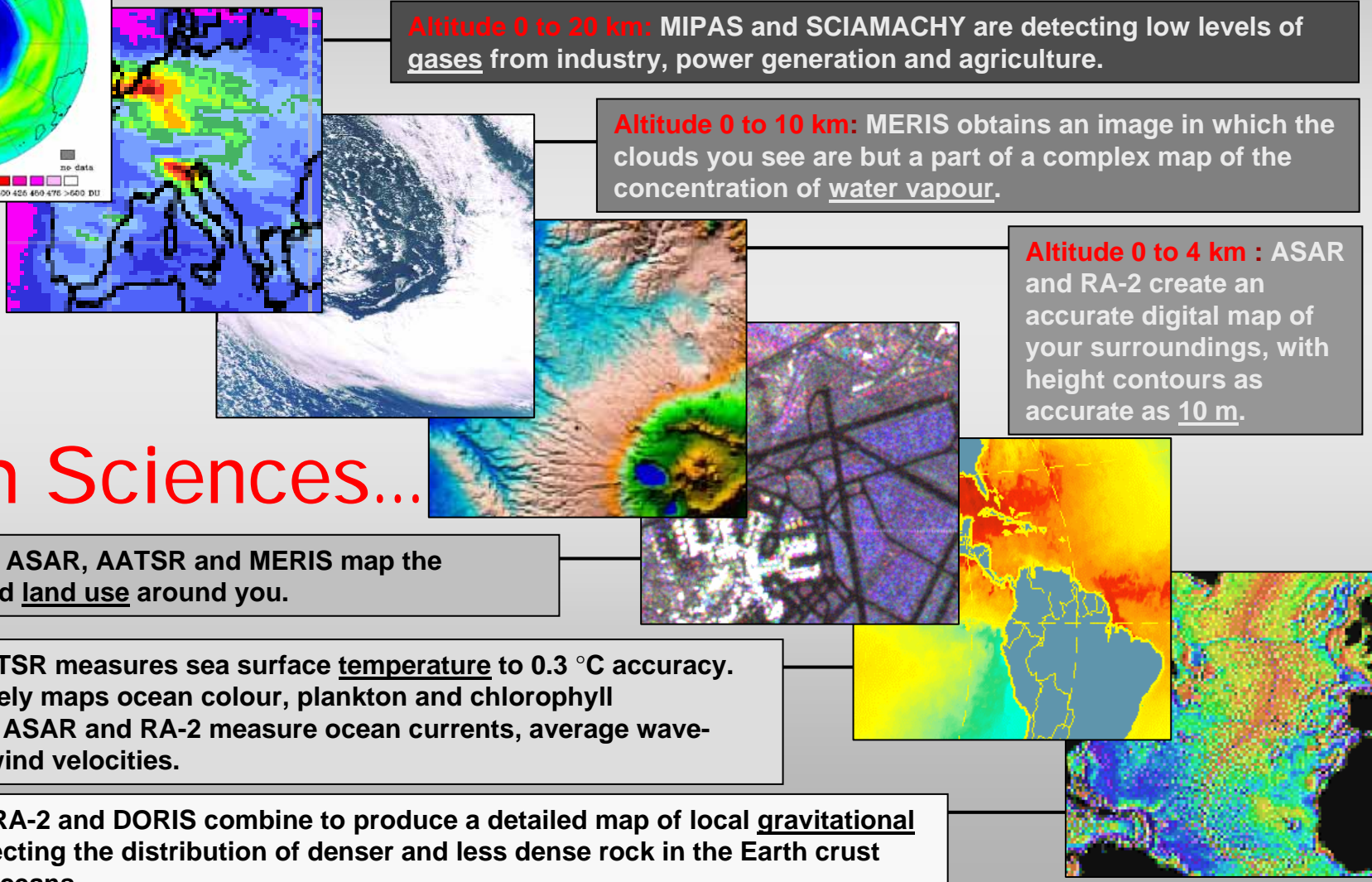
Altitude 0 to 4 km : ASAR and RA-2 create an accurate digital map of your surroundings, with height contours as accurate as 10 m.

...for Earth Sciences...

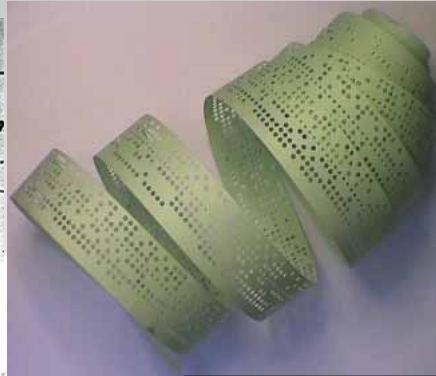
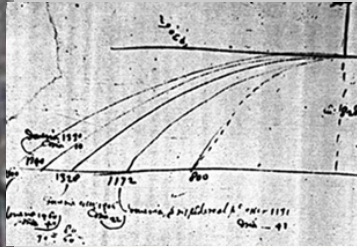
Ground level: ASAR, AATSR and MERIS map the vegetation and land use around you.

Sea level: AATSR measures sea surface temperature to 0.3 °C accuracy. MERIS precisely maps ocean colour, plankton and chlorophyll distributions. ASAR and RA-2 measure ocean currents, average wave-heights and wind velocities.

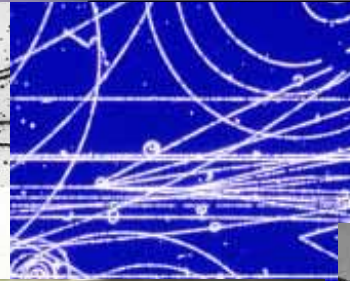
Underwater: RA-2 and DORIS combine to produce a detailed map of local gravitational strength, detecting the distribution of denser and less dense rock in the Earth crust beneath the oceans.



Technik – Lebensdauer Technologie Zyklen - Metadaten



The
Maintenance Service
Department
TABULATING MACHINES
DIVISION
of
Remington Rand



METHODS THAT REVEAL THE FUTURE OF BUSINESS

Remington Rand
PUNCHED-CARD ACCOUNTING

and the Opportunities
it holds for You . . .



Open Access

Exzellente Wissenschaft braucht
langfristig gesicherten,
ungehinderten und umfassenden
Zugang zu erstklassigen
Informationsquellen.

„ ... Science is above all,
communication.“

Josiah Williard Gibbs, physicist

NSF ‘Atkins/Blue Ribbon’ Report on Cyberinfrastructure

- *‘the primary access to the latest findings in a growing number of fields is through the Web, then through classic preprints and conferences, and lastly through refereed archival papers’.*
- *‘archives containing hundreds or thousands of terabytes of data will be affordable and necessary for archiving scientific and engineering information’.*



Hin zu einem Web der Kultur und Wissenschaft

- Kultur- und Wissenschaftskrise im Informationzeitalter
- Vision eines (offenen) Webs der Kultur und Wissenschaft
- → Berlin Erklärung und Verwirklichung der Vision



CERN, “e-science” und offener Zugang

CERN Convention 1954:

... shall provide ... research of pure scientific and fundamental character...

... shall have no concern with work for military requirements and the results of its experimental and theoretical work shall be published or otherwise made generally available.

60% der CERN
Publikationen sind
öffentlich im Web
zugänglich.
CDS, der CERN
Document Server
ist eine recht
grosse e-Bibliothek.



Zusammenfassung

Wir brauchen eine leicht zugängliche, nachhaltige Infrastruktur mit:

- hohen Rechnerleistungen,
- grossen Daten-Speicherkapazitäten,
- schnellen Datenleitungen,
- interoperablen Datenarchiven,
- einfachem Zugriff zu innovativen Datenbearbeitungswerkzeugen,
- Registraturen von Service-Einrichtungen, Datenbeständen (inklusive Dokumentationen, Publikationen, etc.)
- schnellen Suchalgorithmen,
- einer Plattform, wo Forschungsergebnisse et al. öffentlich verfügbar gemacht und diskutiert werden
- einer Open-Access Philosophie
- einer Langzeit-Archivierungs- und –Daten-Zugangs-Strategie

Offene Fragen bzgl. fachbezogener, nationaler, internationaler
Lösungen - **Absprache, Zusammenarbeit, Standards notwendig**

e-Science: Modell für die Wissenschaft des 21-ten Jahrhunderts → die ganze e-Science massiv fördern

Offene e-Bibliotheken: notwendige Bedingung → auch massiv fördern

Dauerhafte Aufbewahrung von Daten und Wissen: eine grosse Herausforderung an die Wissenschaft und Technologie

Wir brauchen eine neue, virtuelle, digitale, verteilte und von der Wissenschaft “zertifizierte” Bibliothek von Alexandria mit dem Wissen und der Kultur der Welt

Besten Dank