# Cyberinfrastructure Challenges
# (from a climate science repository perspective)

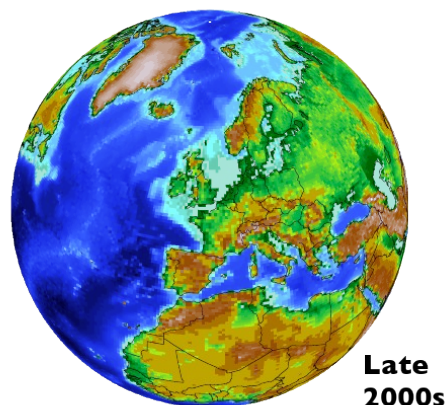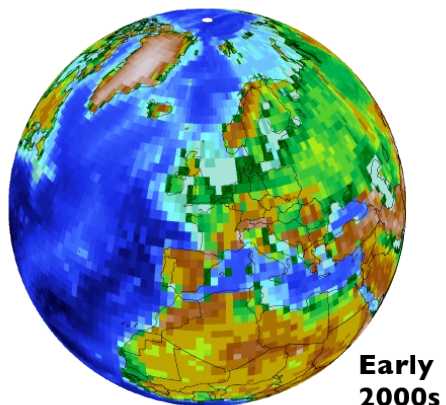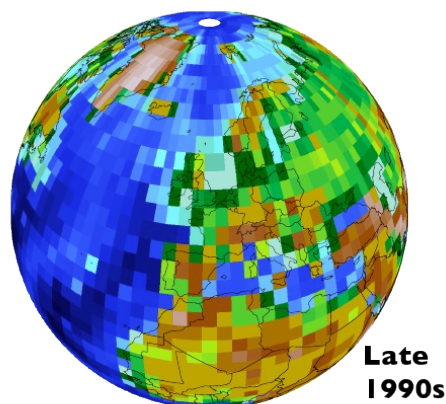Bryan Lawrence

CEDA

Rutherford Appleton Laboratory

# Simulation Data Deluge

Fifth coupled model intercomparison project (CMIP5) (running now)

- Petabytes of output
- Globally synchronised petascale cache(s)
- Millions of Datasets aimed at different user communities!
- Comprehensive Metadata Structures
- Comprehensive Services

CMIP5 is a GLOBAL problem (the simulations are generated globally and consumed globally)!

Solutions need to be global!

Early 1990s

Late 1990s

Early 2000s

Late 2000s

Globes courtesy of Gary Strand (NCAR)

National Problem Too!

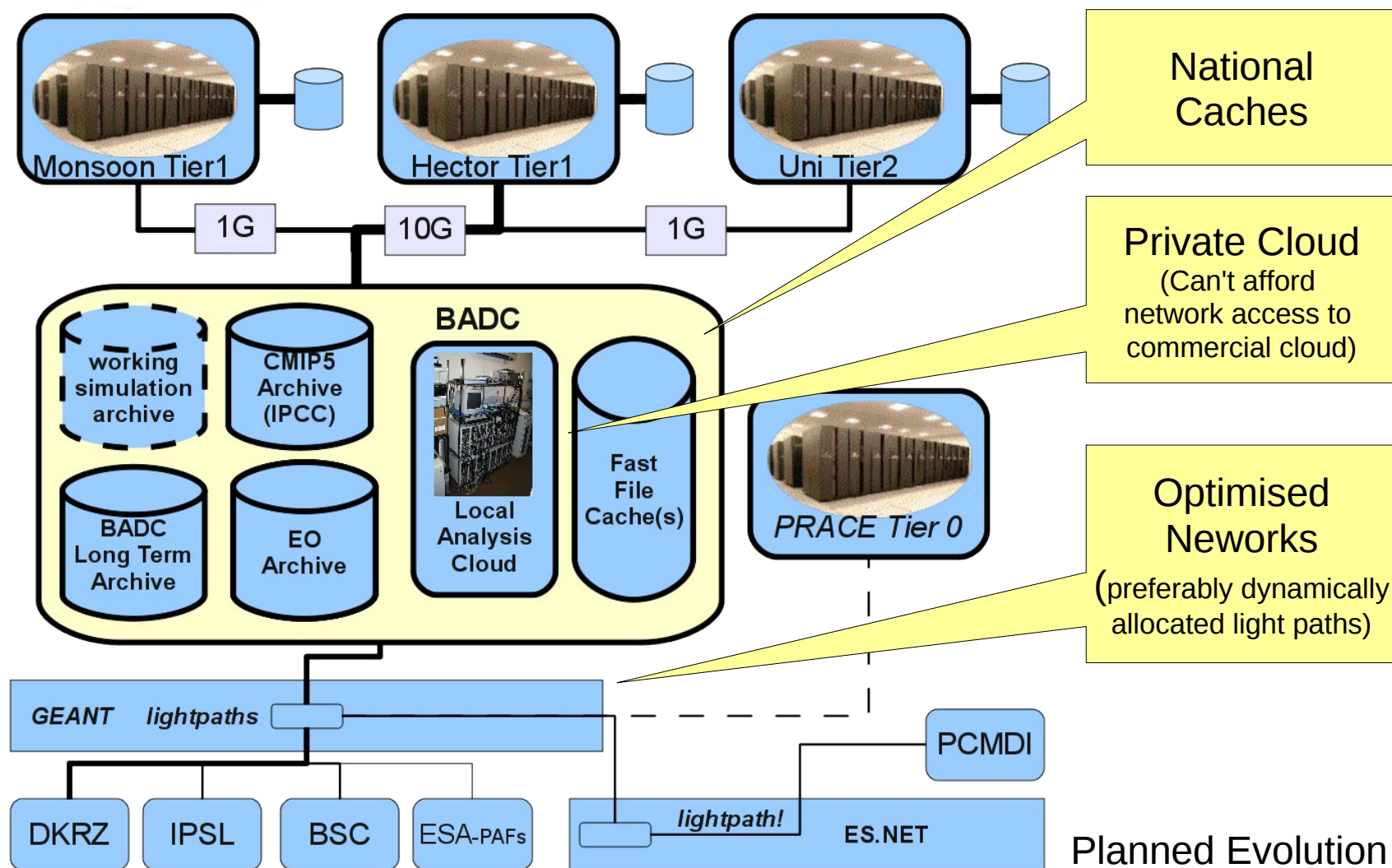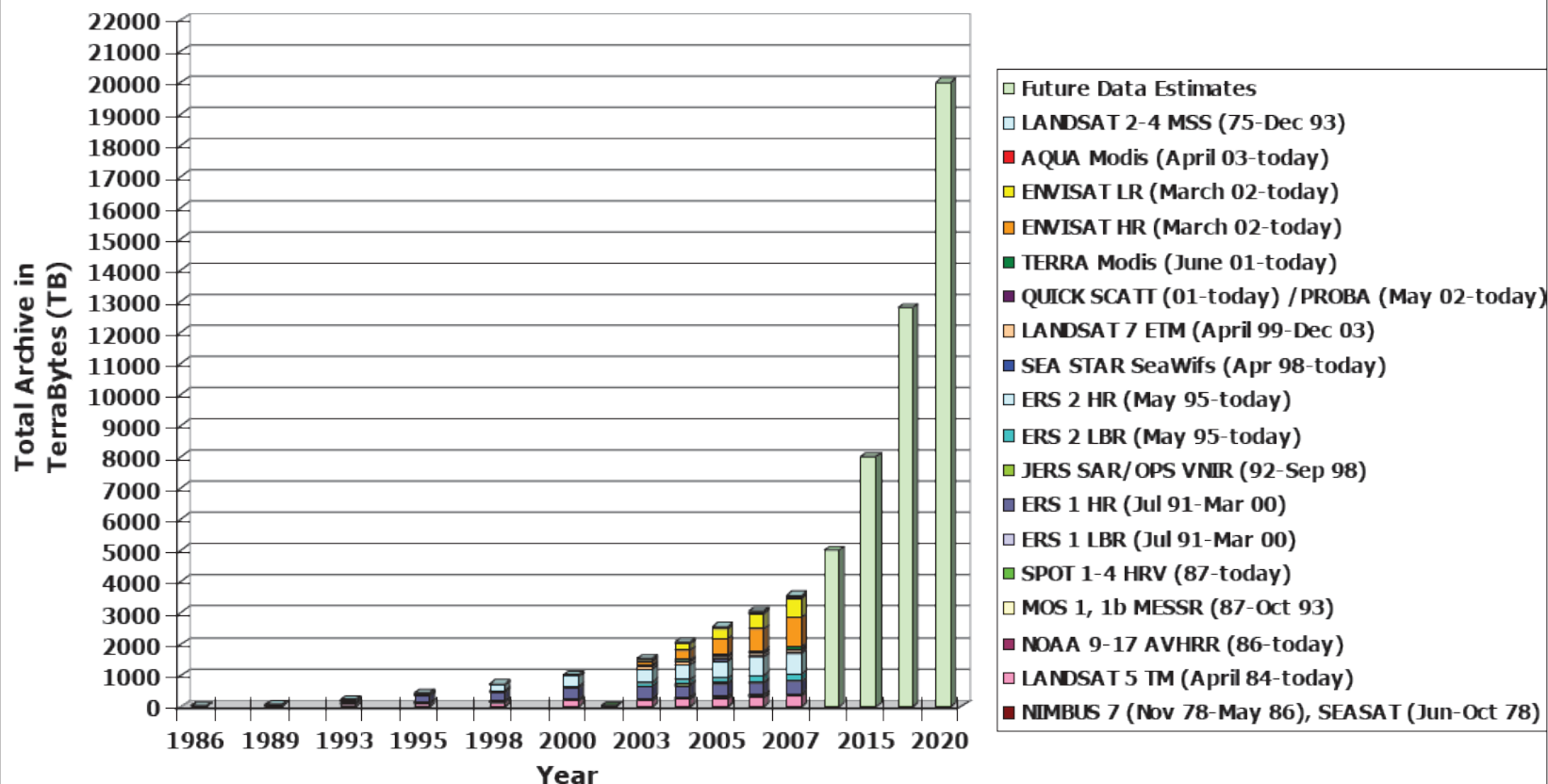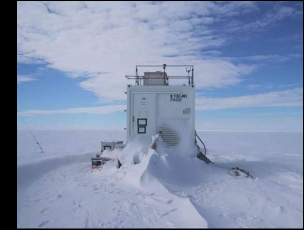Evolution of ESA's EO Data Archives between 1986-2007 and future estimates (up to 2020)

Source: ESA GSCB Workshop June 2009

# Observatories and Sensor (networks)

# Storage can't keep up!

Figure 2

Information Versus Available Storage

(All data, not just scientific data)

Regardless of how good we are at data systems, science will not escape the general trend: more data being produced than can be stored, which means we need to work smarter:

- Better a priori discrimination of what we should keep
  - Don't even bother writing it to any storage.
- Better documentation of what we have produced, to inform initial decisions about what to keep.
  - Decide quickly about whether to move it to working storage.
- Appraisal of what we have kept (if it's big – don't bother if it's small)
  - Avoid holding data which is irrelevant.

British Atmospheric Data Centre
NATIONAL CENTRE FOR ATMOSPHERIC SCIENCE
NATURAL ENVIRONMENT RESEARCH COUNCIL

Centre for Environmental Data Archival
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Earth Observation
NATURAL ENVIRONMENT RESEARCH COUNCIL

# And the Challenge?

Simulation + Earth Observation
+ Sensor Networks
( +looking into the past )

=     Information about the environment

(all individually increasing their output and proliferating in a heterogeneous and geographically distributed manner)

(which needs integration into a coherent view and interpretation)

Cyberinfrastructure Challenges: from the global large scale data transport and storage, national caches, to automatic/manual metadata creation/entry (*reliable tools to get the metadata to drive it all*) and the systems (including ontology systems) to interpret it all.

**British Atmospheric Data Centre**
NATIONAL CENTRE FOR ATMOSPHERIC SCIENCE
NATURAL ENVIRONMENT RESEARCH COUNCIL

Centre for Environmental Data Archival
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL

**National Centre for Earth Observation**
NATURAL ENVIRONMENT RESEARCH COUNCIL

**Science Driven**

Not really up to High volume, Long distance

Point and Click Doesn't Scale

Open Access ≠ insecure and/or overloaded

Internal, External, Annotation (A,B,C,D,E & more)

Calculation & "Sophisticated" visualisation: need more standard APIs

Getting metadata is hard: need much better tools. NOTHING SCALES without metdata!

The solution to HETEROGENEITY is STANDARDISATION (with FLEXIBILITY) + MODEL DRIVEN ARCHITECTURES!

Semantics Matter! Need to get beyond serialisation and simple unstructured Relationships (linked data, I'm looking at you!)

... have structure: need "Data Models" (independent of Storage schema)

Usage Conventions

Energy Cost versus Availability

Import role for metamodels

Existing Cyberinfrastructure too FLAKEY

**Technology Driven**

Central stack:
- Portals
- Applications (inc Scripts)
- Service Clients
- Transport
- Security: AAA + Policy
- Services
- External Metadata
- Data Items
- Internal Metadata
- File Formats/ Database Type
- Distribution Management
- Transport
- File System
- Physical Servers

**British Atmospheric Data Centre**
NATIONAL CENTRE FOR ATMOSPHERIC SCIENCE
NATURAL ENVIRONMENT RESEARCH COUNCIL

**Centre for Environmental Data Archival**
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL

**National Centre for Earth Observation**
NATURAL ENVIRONMENT RESEARCH COUNCIL

# Rewards

# Curation

# Citation

# Licenses & IPR

# Trust

# Reliance

# Plans