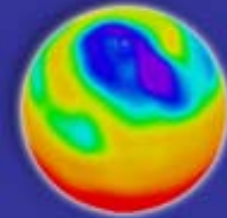# BRITISH ATMOSPHERIC DATA CENTRE
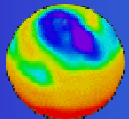
# **Presenting a multi-terabyte dataset via the web**

**Ag Stephens**

BADC Data Scientist
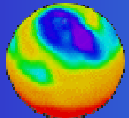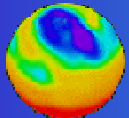
11 November 2003

**http://badc.nerc.ac.uk**

# Presentation outline

- An introduction to the BADC.

- The project stages for delivering a Live Access Server:

  1. Project design.

  2. Tools to convert the data.

  3. Designing a caching architecture.

  4. Aggregation of data files.

  5. Setting up the Live Access Server.

- Further demonstration and conclusions.

# What is the BADC?

• The NERC-designated data centre for atmospheric science.

• Over 20 TB of data.

• Serving around 5,000 users.

• Agreement with Met Office and ECMWF to distribute data.

http://badc.nerc.ac.uk

# How people use the BADC



http://badc.nerc.ac.uk

Int

1.

2.

3.

4.

5.

6.

7.

8.

/cdat

**http://badc.nerc.ac.uk**

# LAS Project Stage 3:  Caching

- Cache copy of directory structure.

- Cache algorithms written in Python.

- Control data volumes.

- Analyse and process request sizes.

- Cache of about 1 TB initially.

**http://badc.nerc.ac.uk**

Climate Data Markup Language (CDML) files are created by the cdscan utility.

CDML contains the following sections:

<dataset>  - general information at the dataset level.

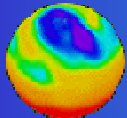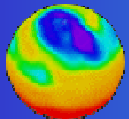<axis> - axis dimension information.
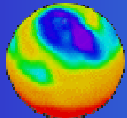
<variable> - relating to individual variables.
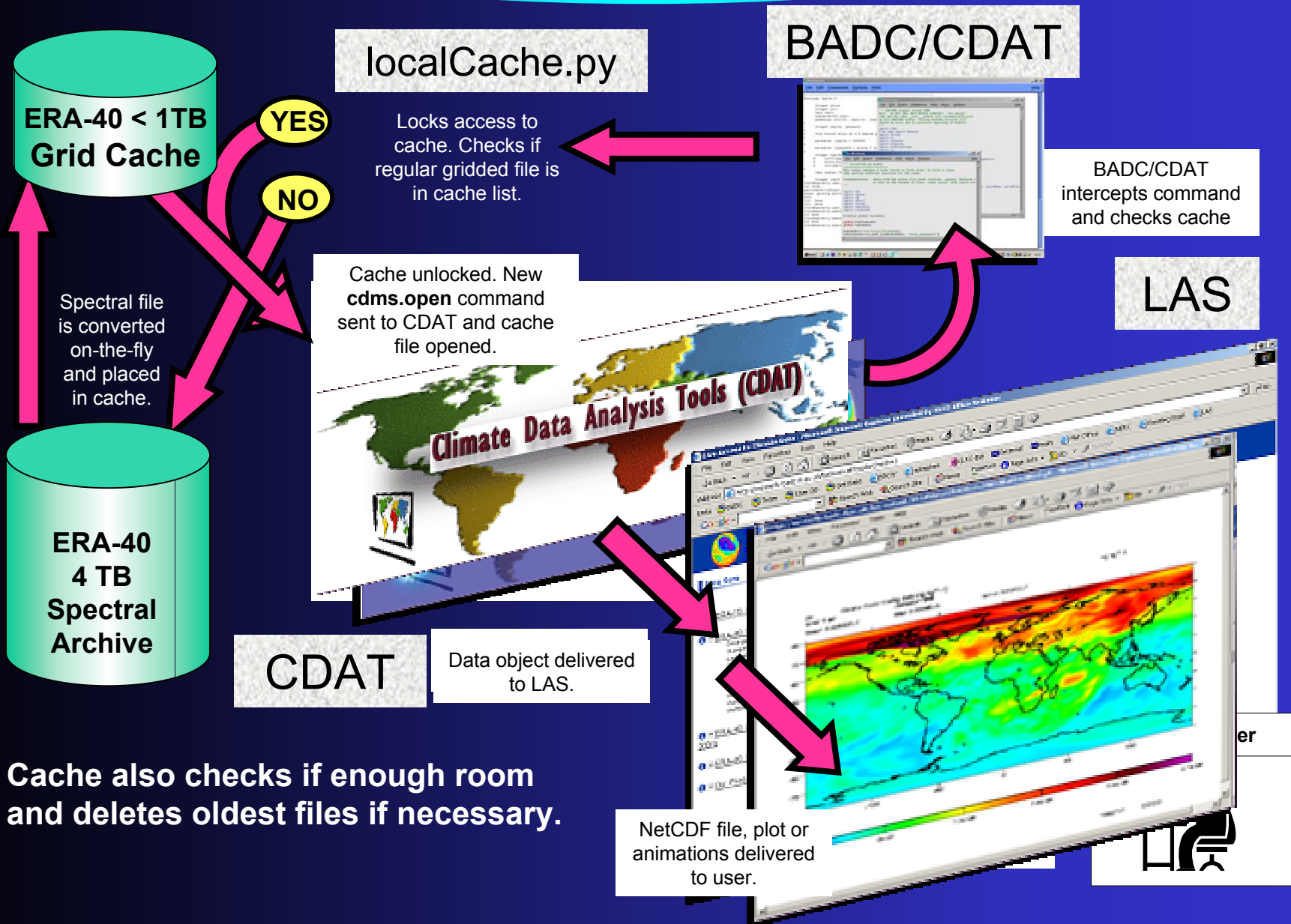
*3,000,000 files from one 21KB XML file!*

http://badc.nerc.ac.uk

Work required to configure LAS:

1. Configuring Apache webserver (RedHat Linux).

2. Configuring Tomcat Java Servlet Engine.

3. Interfacing to MySQL database.

4. Ingesting CDML files into LAS.
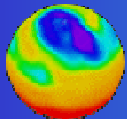
5. Security layer (pending).

**NCAS**
NATURAL
ENVIRONMENT
RESEARCH COUNCIL

**http://badc.nerc.ac.uk**

**CLRC**

**NATURAL ENVIRONMENT RESEARCH COUNCIL**

# How it all fits together

**ERA-40 < 1TB Grid Cache**

**YES**

**NO**

localCache.py

Locks access to cache. Checks if regular gridded file is in cache list.

BADC/CDAT

BADC/CDAT intercepts command and checks cache

Spectral file is converted on-the-fly and placed in cache.

Cache unlocked. New **cdms.open** command sent to CDAT and cache file opened.

LAS

**ERA-40 4 TB Spectral Archive**

Climate Data Analysis Tools (CDAT)

**CDAT**

Data object delivered to LAS.

**Cache also checks if enough room and deletes oldest files if necessary.**

NetCDF file, plot or animations delivered to user.

# BADC LAS Demo 1: 1 month to NetCDF

http://badc.nerc.ac.uk

**ERA-40 Re-analysis Data**

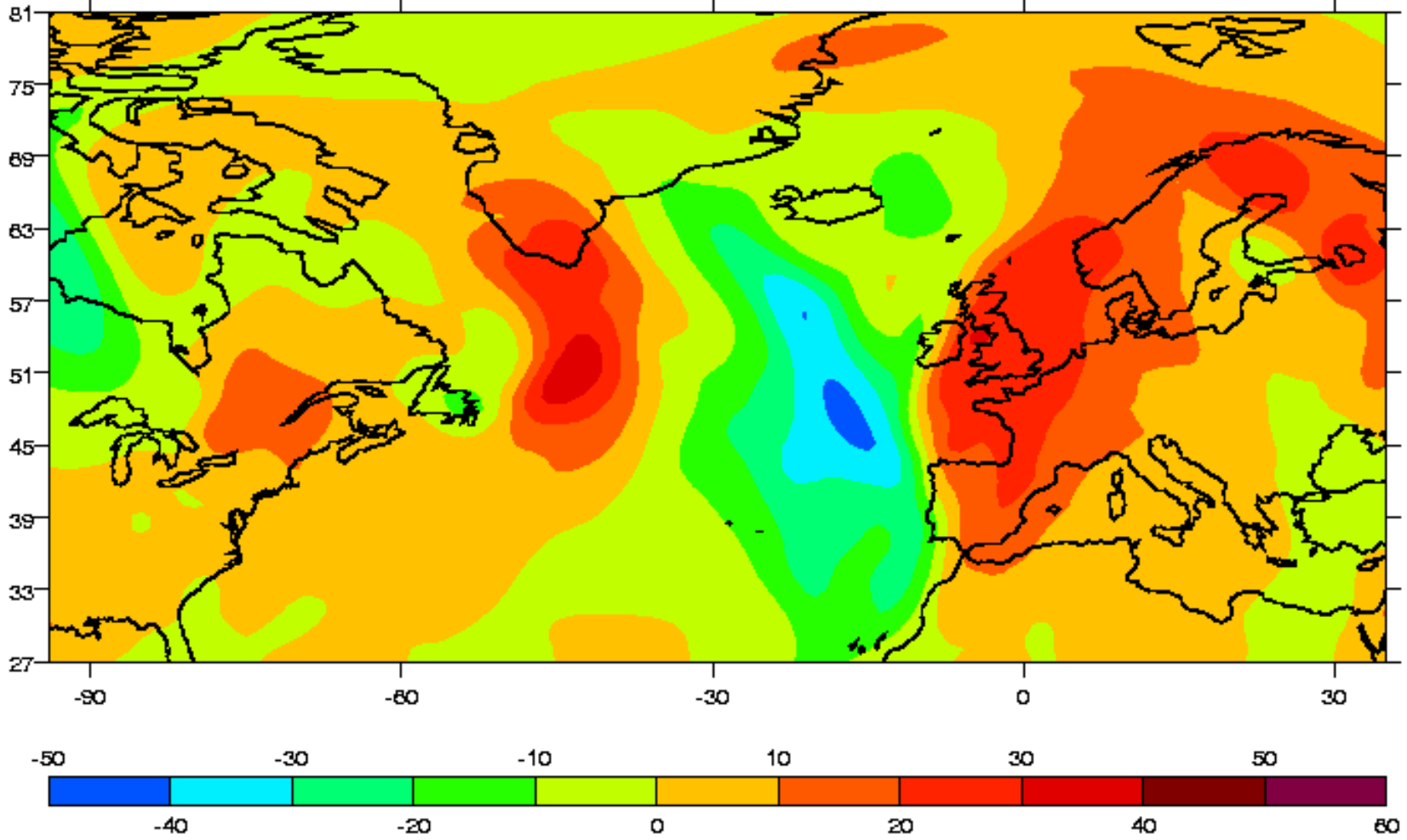v        V-velocity (ms**-1)                                      1987/10/10      0:0:0.0
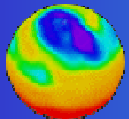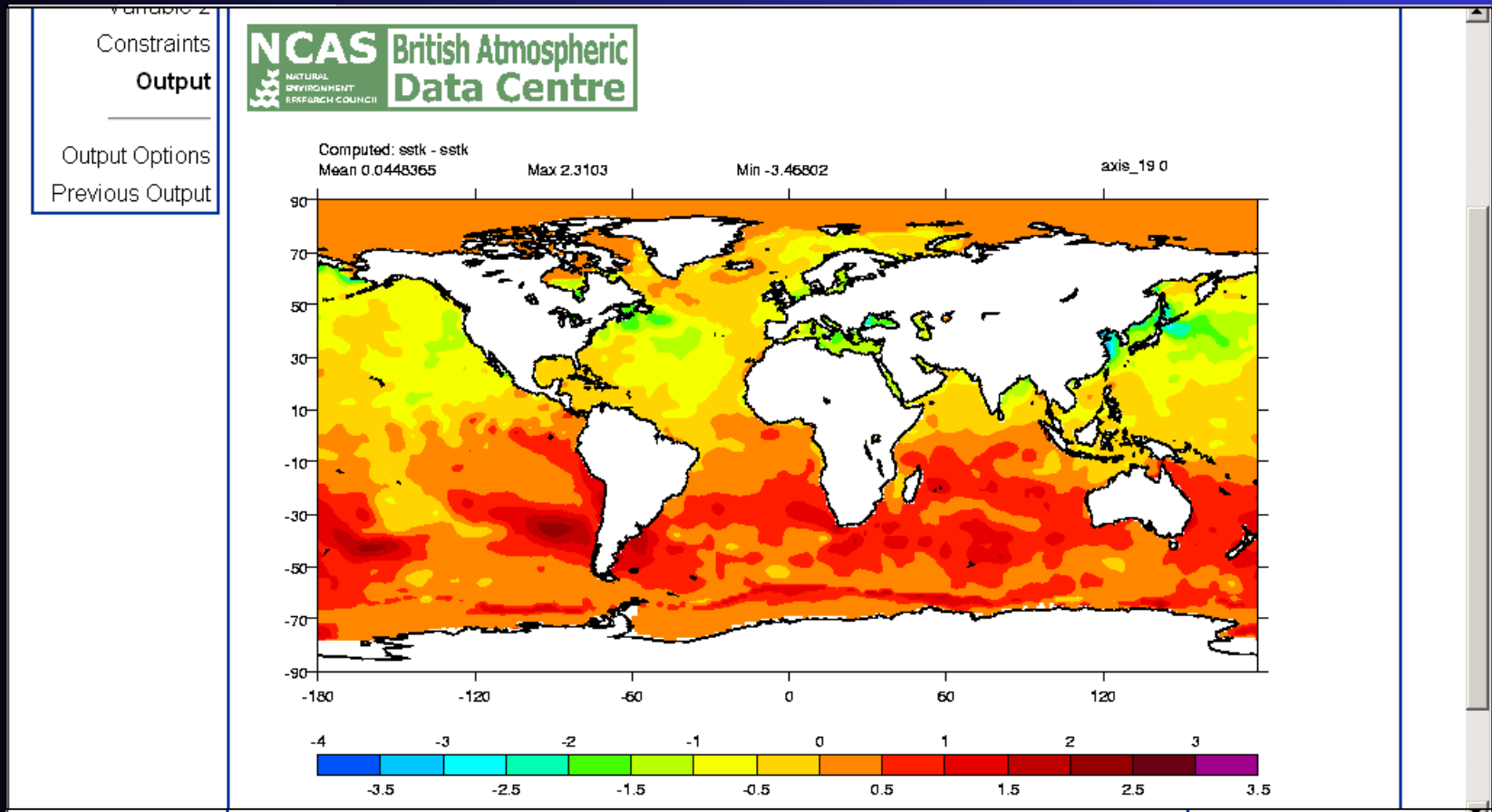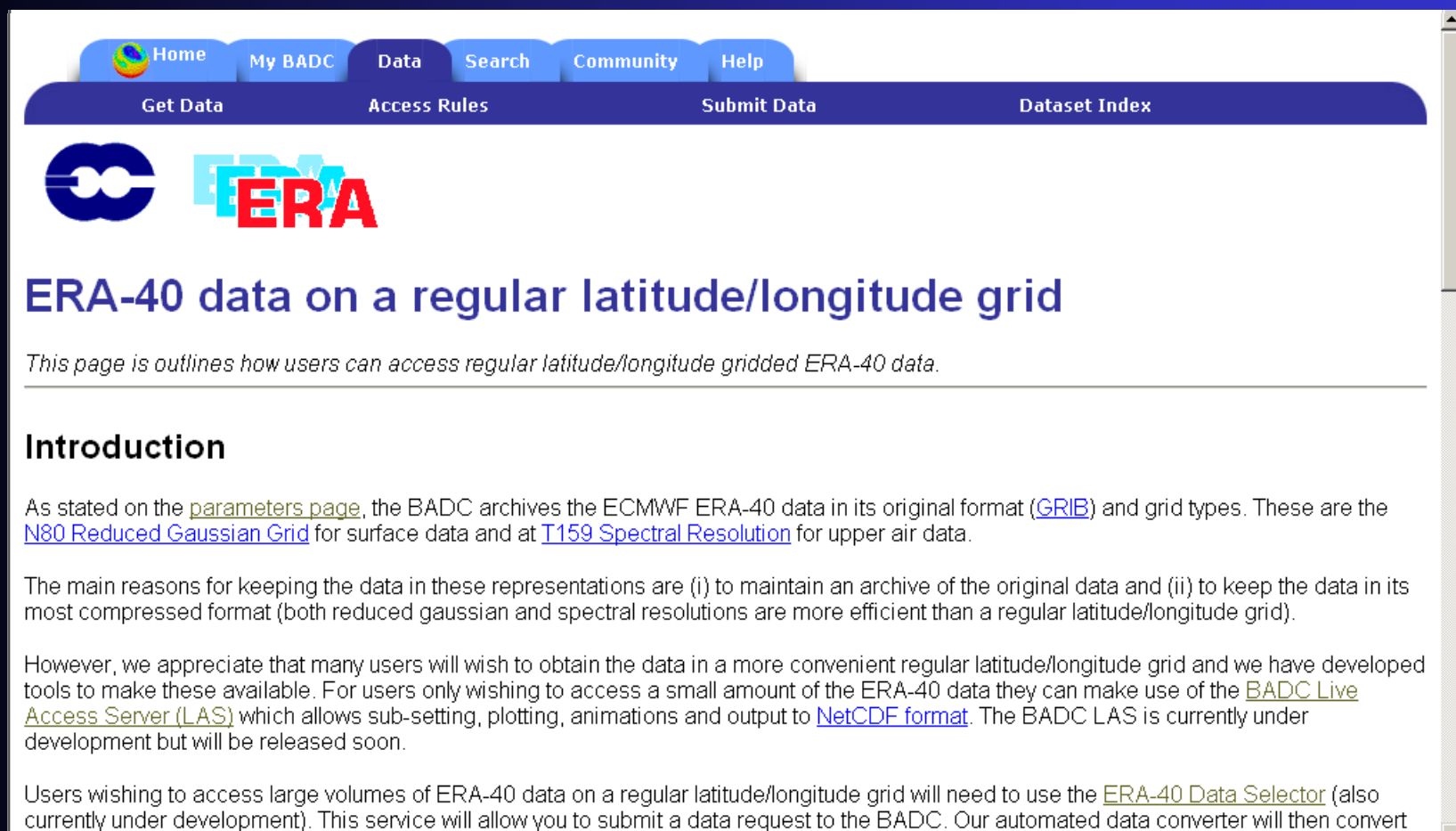
Mean 1.12676             Max 36.7089                 Min -42.7208

# BADC LAS Demo 3: Departure plot

# BADC LAS Demo: Metadata



**http://badc.nerc.ac.uk**

Advantages of our approach:

- Multiple TB via one interface – some virtual!

- Users saved from files and formats.

- New options for sub-setting and plotting.

- Automatic monitoring of data usage.

- Caching system available for other purposes.

- Knowledge of CDAT and LAS for other projects.

**NCAS**
NATURAL
ENVIRONMENT
RESEARCH COUNCIL

**http://badc.nerc.ac.uk**

CLRC

**NATURAL
ENVIRONMENT
RESEARCH COUNCIL**

# What have we learnt?

Disadvantages:

- No automatic response to massive requests.

- Limits to configurations of plots and animations.

- Caching database is slow.

- Only one dataset presented so far.

We plan to:

- Implement parallel LASes (ECMWF, UM, COAPEC).

- Implement a time algorithm to keep users informed.

- Generate user-defined LASes on-the-fly.

- Allow comparison of different datasets.

- Re-think the caching database interaction for speed.

- Look to parallelise the background file conversions.

**http://badc.nerc.ac.uk**

NCAS
NATURAL
ENVIRONMENT
RESEARCH COUNCIL

CLRC

NATURAL
ENVIRONMENT
RESEARCH COUNCIL

# Useful Links

**BADC:**    http://badc.nerc.ac.uk

**CDAT:**    http://esg.llnl.gov/cdat

**CDML:**    http://esg.llnl.gov/cdat/cdms_html/cdms-6.htm

**LAS:**    http://ferret.pmel.noaa.gov/Ferret/LAS/ferret_LAS.html

**Pyfort:**    http://pyfortran.sourceforge.net

NCAS
NATURAL ENVIRONMENT RESEARCH COUNCIL
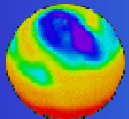
CLRC

NATURAL ENVIRONMENT RESEARCH COUNCIL

# LAS Project Overview

**ARCHIVE**

**CACHE**

**DELIVERY**

Conversion / Caching

4 TB

8 TB

Aggregation / User Interface

18 TB

**Python**
**EMOS**
**Pyfort**
**grib2ctl.pl**
**gribmap**

**CDAT**
**LAS**

**Spectral & Gaussian**
**Permanent**
**GRIB**

**1 degree grid**

**Temporary**
**GRIB**

**1 degree grid**

**Short-term**
**NetCDF/plots**