# Science Drivers: Why JASMIN?

**Bryan Lawrence**

Rising demand     The Data Commons     Looking Forward     Summary
●○○○○○○○○     ○○○○○     ○○○○○     ○
The Big Picture

## New ways of thinking: Data Intensive Science

The four paradigms of science:

1. **E**xperimental Science - Repeatable experiments observing and/or interfering with and observing nature.

2. **T**heoretical Science - Mathematics and Abstraction and analytical solutions compared with experiment.

3. **C**omputational Science - Numerical solutions of equations and simulations of structure and evolution.

4. **D**ata Intensive Science - All the above but mediated by so much data that the science starts ab-inito from (often) someone else's data.
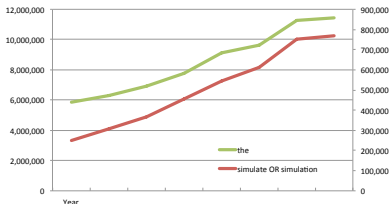


Jim Gray
(Tony Hey et al ...)

We are still working through the consequences of the advent of data intensive science, but JASMIN is one!

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand | The Data Commons | Looking Forward | Summary

Where's this coming from? Scientific Pull underpinned by Technology Push
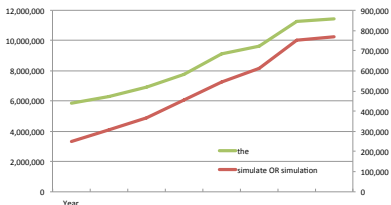
## The growth of simulation

A proxy for the use of simulation in science: the growth in the numbers of (scientific) papers which mention *simulate* or *simulation* in their abstracts (via ISI).



In the last fifteen years (2001-2015), the number of papers published have doubled, but the number of papers with simulation in the abstract have tripled.

Rising demand | The Data Commons | Looking Forward | Summary

Where's this coming from? Scientific Pull underpinned by Technology Push

## The growth of simulation

A proxy for the use of simulation in science: the growth in the numbers of (scientific) papers which mention *simulate* or *simulation* in their abstracts (via ISI).
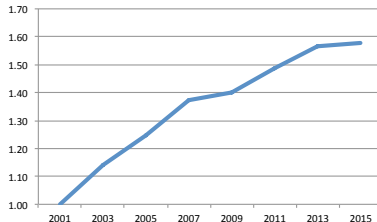


In the last fifteen years (2001-2015), the number of papers published have doubled, but the number of papers with simulation in the abstract have tripled.



Even if we conclude that a doubling of papers doesn't reflect a doubling in science of the community, we can still conclude that the proportion of the community doing, or exploiting, simulation has grown by 50%!

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand ○○●○○○○○○○
The Data Commons ○○○○○
Looking Forward ○○○○○
Summary ○

Where's this coming from? Scientific Pull underpinned by Technology Push

## Core Science Requirements

Schematic for Global Atmospheric Model

Big International Drivers:

WWRP — WORLD WEATHER RESEARCH PROGRAMME

GAW

Copernicus — Europe's eyes on Earth

WCRP — World Climate Research Programme

| Today: | Observations | Models |
|--------|--------------|--------|
| Volume | 20 million = $2 \times 10^7$ | 5 million grid points<br>100 levels<br>10 prognostic variables = $5 \times 10^9$ |
| Type | 98% from 60 different satellite instruments | physical parameters of atmosphere, waves, ocean |

| Soon: | Observations | Models |
|-------|--------------|--------|
| Volume | 200 million = $2 \times 10^8$ | 500 million grid points<br>200 levels<br>100 prognostic variables = $1 \times 10^{13}$ |
| Type | 98% from 80 different satellite instruments | physical and chemical parameters of atmosphere, waves, ocean, ice, vegetation |

→ Factor 10 per day   → Factor 2000 per time step
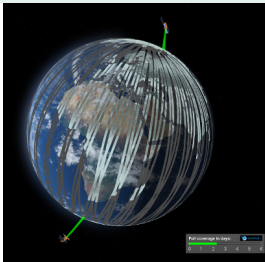
→ but many more time steps needed

aerosol cci
cloud cci
fire cci
ghg cci
glaciers cci
antarctic ice sheet cci
ice sheets greenland cci
land cover cci
ocean colour cci
ozone cci
sea ice cci
sea level cci
sst cci
soil moisture cci
cmug cci

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence – RAL, June 2016

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

# The Sentinels: Big EO data crucial to NERC science!



## Sentinels

Sentinel 1A (2014), 1B (2016)
Sentinel 2A (2015) *2B (2017?)*
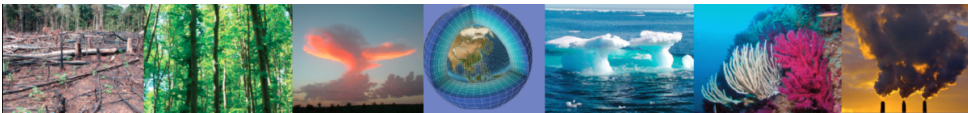Sentinel 3A (2016) *3B (2018?)*
Data rate: o(6) PB/year

NERC SCIENCE OF THE ENVIRONMENT

COMET: Centre for Observation and Modelling of Earthquakes, Volcanoes, and Tectonics



Inteferogram measuring deformation
Napa Valley 8/2014
Napa

San Pablo Bay

COMET
norut
PPOLabs

(Picture credits: ESA, Arianespace.com, PPO.labs-Norut–COMET-SEOM Insarap study, ewf.nerc.ac.uk/2014/09/02/new-satellite-maps-out-napa-valley-earthquake/ )

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○●○○○○
From a modelling perspective

The Data Commons
○○○○○

Looking Forward
○○○○○

Summary
○

## Infrastructure Requirements



### Infrastructure Strategy for the European Earth System Modelling Community 2012-2022

**Key science questions**

1. How predictable is climate on a range of timescales ?
2. What is the sensitivity of climate and how can we reduce uncertainties ?
3. What is needed to provide reliable predictions of regional climate changes ?
4. Can we model and understand glacial-interglacial cycles ?
5. Can we attribute observed signals to understand processes ?

**Recommendations:**

1. Provide a blend of high-performance computing facilities . . .
2. Accelerate the preparation for exascale computing . . .
3. Ensure data from climate simulations are easily available and well documented, especially for the climate impacts community.
4. Build a physical network connecting national archives with transfer capacities exceeding Tbits/sec.
5. Strengthen the European expertise in climate science and computing . . .

enes
EUROPEAN NETWORK
FOR EARTH SYSTEM MODELLING

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

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

Rising demand
○○○○○●○○○
Interdisciplinary Science!

The Data Commons
○○○○○

Looking Forward
○○○○○

Summary
○

# The Propagation of Direct Numerical Simulation



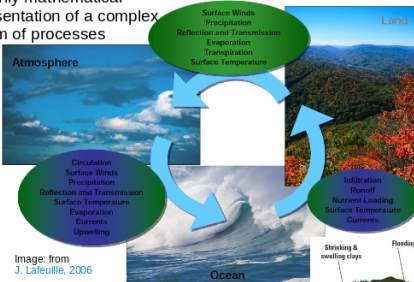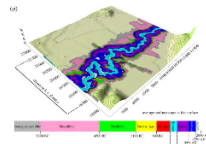Primarily mathematical representation of a complex system of processes

Image: from
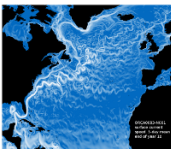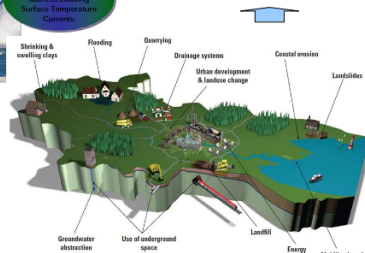J. Lafeuille, 2006

Coulthard and Van De Wiel IDoi: 10.1098/rsta.2011.0597

http://www.bgs.ac.uk/research/environmentalModelling/home.html

More communities want to observe and simulate the world at ever higher resolution!

More complexity!

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence – RAL, June 2016

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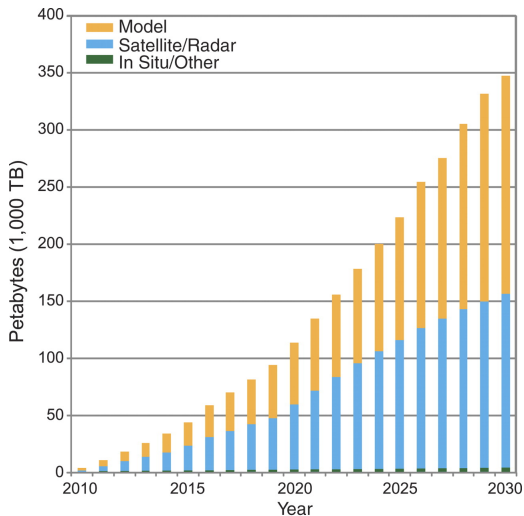NATURAL ENVIRONMENT RESEARCH COUNCIL

## Communities



Many interacting communities, each with their own software, compute environments, observations etc.

Figure adapted from Moss et al, 2010

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

Centre for Environmental Data Analysis
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
NATURAL ENVIRONMENT RESEARCH COUNCIL
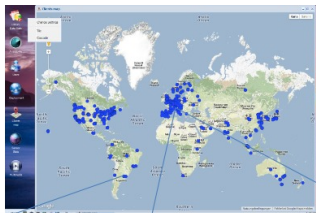
## More Data

Fig. 2 The volume of worldwide climate data is expanding rapidly, creating challenges for both physical archiving and sharing, as well as for ease of access and finding what's needed, particularly if you're not a climate scientist.

(BNL: Even if you are?)



J T Overpeck et al. Science 2011;331:700-702

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○●
Consequences

The Data Commons
○○○○○

Looking Forward
○○○○○

Summary
○

## The trend



**Slide courtesy of Stefan Kindermann, DKRZ and IS-ENES2**

**Individual End Users**
- Limited resources (bandwidth, storage,..)

**Organized User Groups**
- Organize a local cache of required files
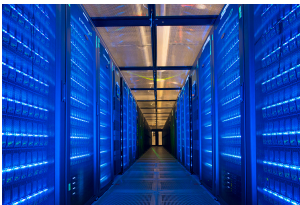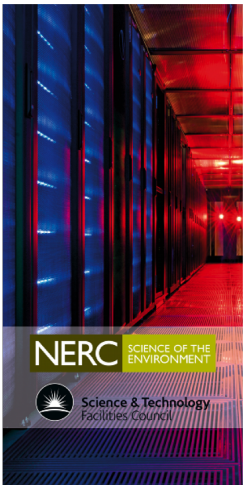- Most of group don't access ESGF, use cache instead!

**Data Centre Service Group**
- Provides access to ESGF replica cache
- May also provide access to data near compute resources
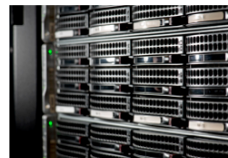- (BADC, DKRZ, IPSL, KNMI, UC)

Trend

Needed: Replacement for „*Download and Process at Home*" Approach

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○○
Hardware

The Data Commons
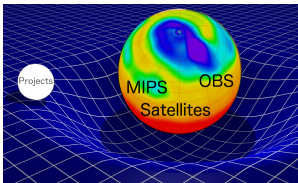●○○○○

Looking Forward
○○○○○

Summary
○

# So we have built a Data Intensive Computing System: JASMIN







- ▶ 16 PB Fast Storage
  (Panasas, many Tbit/s bandwidth)

- ▶ 1 PB Bulk Storage

- ▶ Elastic Tape

- ▶ 4000 cores: half deployed as hypervisors, half as the "Lotus" batch cluster.

- ▶ Some high memory nodes, a range, bottom heavy.

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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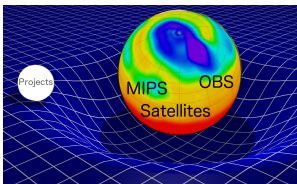NATURAL ENVIRONMENT RESEARCH COUNCIL

# JASMIN — The Data Commons



- ▶ Provide a state-of-the art storage and computational environment
- ▶ Provide and populate a managed data environment with key datasets (the "archive").
- ▶ Encourage and facilitate the bringing of data and/or computation alongside/to the archive!
- ▶ Provide flexible methods of exploiting the computational environment.

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○○
Data Gravity in the Commons

The Data Commons
○●○○○○

Looking Forward
○○○○○

Summary
○

# JASMIN — The Data Commons



Projects

MIPS
Satellites
OBS

- Provide a state-of-the art storage and computational environment
- Provide and populate a managed data environment with key datasets (the "archive").
- Encourage and facilitate the bringing of data and/or computation alongside/to the archive!
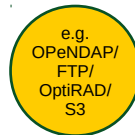- Provide flexible methods of exploiting the computational environment.

e.g.
CEMS

**Platform as a Service**
-----
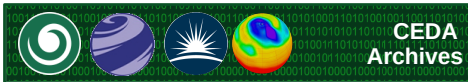We provide you the "Platform"; you can LOGIN and exploit the batch cluster.

e.g.
BIOLINUX

**Infrastructure as a Service**
-----
We provide you with a cloud on which you INSTALL your own computing.

e.g.
OPeNDAP/
FTP/
OptiRAD/
S3

**Software as a Service**
-----
We provide you with REMOTE access to data VIA web and other interfaces.
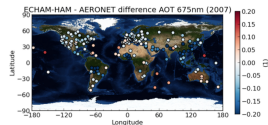
**CEDA Archives**

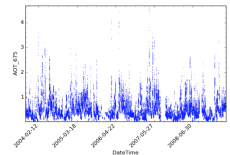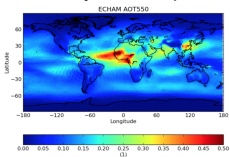**JASMIN – Data Intensive Computer**
Storage, Compute and Network Fabric
Batch Compute, Private Cloud, Disk, Tape

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
OOOOOOOOOO
Software in the Commons

The Data Commons
OO○●OO

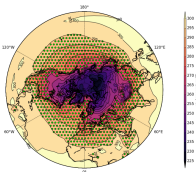Looking Forward
OOOOO

Summary
O

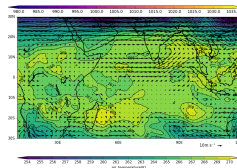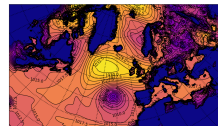## Tools

# JASMIN Analysis Platform

Community Intercomparison Suite:



CF-python, CF-plot, cfview:



```
import cf, cfplot as cfp
f=cf.read('/opt/graphics/cfplot_data/tas_A1.nc')
g=f.subspace(time=15)
cfp.gopen()
cfp.cscale('magma')
cfp.mapset(proj='npstere')
cfp.con(g)
cfp.stipple(f=g, min=265, max=295, size=100, color='#00ff00')
cfp.gclose()
```

http://www.cistools.net/
http://cfpython.bitbucket.org
http://ajheaps.github.io/cf-plot
...and many more ... all shared and (hopefully) kept up to date on the JAP.

See Duncan Watson-Parris et al, 2016 (doi:10.5194/gmd-2016-27)

National Centre for
Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○○
Algorithms

The Data Commons
○○○●○

Looking Forward
○○○○○

Summary
○

## Common Software/Algorithm Patterns

Supporting a wide variety of
algorithms and workflows:
(but much to do to exploit
parallelism)



"Big Data Ogres"
by analogy with the Berkeley
Dwarves for computational
patterns.

**Different Problem Architectures, e.g:**

1. Pleasingly Parallel (e.g. retrievals over images)
2. Filtered pleasingly parallel (e.g. cyclone tracking)
3. Fusion (e.g. data assimilation)
4. (Space-)Time Series Analysis (FFT/MEM etc)
5. Machine Learning (clustering, EOFs etc)

**Important Data Sources, e.g:**

1. Table driven (eg. RDBMS + SQL)
2. Document driven (e.g XMLDB + XQUERY)
3. Image driven (e.g. GeoTIFF + your code)
4. (Binary) File driven (e.g. NetCDF + your code)

**Sub-Ogres: Kernels & Applications, e.g:**

1. Simple Stencils (Averaging, Finite Differencing etc)
2. 4D-Variational Assimilation/ Kalman Filters
3. Data Mining Algorithms (classification/clustering) etc
4. Neural Networks

*Modified from Jha et al 2014 arXiv:1403.1528[cs]*

National Centre for
Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
OOOOOOOOOO
Algorithms

The Data Commons
OOOO⬤O

Looking Forward
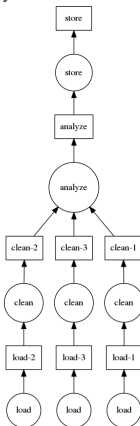OOOOO

Summary
O

# Uncommon software solutions

Plethora of parallel architectures and tools out there



Contrast between two very different parts of our workflow:

▶ Many of our analysis tasks are build once, use once, throwaway. No room for optimisation (or MPI), yet

▶ Much of our workflow is repeatable: "build", "run", "move", "reduce/reformat", "analyse". Much room for automation.

Whatever tools we use, we'll need to get use to generating, understanding, and exploiting concurrency in more complicated ways:



Much to do, as infrastructure providers, and users, to harness these tools to accelerate our workflows!

(These two examples: dsk, and cylc, representing analysis and scheduling, reduction and proliferation.)

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

## Local growth thus far?

JASMIN phase 1 and phase 2: total disk storage:



Within that, the archive growth rate:



Note growth rates:

- from early 2013 to early 2015: 1.25 PB/year
- since early 2015: 2.5 PB/year (ARCHER upgrade?)

▶ Note the steep rise in 2012/2013 associated with CMIP5 (we expect CMIP6 to be ten times as big and arriving in the same sort of time duration).

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

# Capacity for future growth?



J1 Storage (June 2016) (6,262 RAW TB)
- Allocated (1,270)
- Used (4,677)
- Unallocated (312)



J2 Storage (June 2016) (11,533 RAW TB)
- Allocated (3,743)
- Used (5,309)
- Unallocated (2,479)

▶ J1 bladesets are more or less full, with a little headroom for user data analysis.

▶ 2 PB (Raw) of room for new allocations (≈ one year until JASMIN is full!)

It's worth noting that while we (have) need(ed) parallel file systems for performance and ease of management at the petascale, they don't come without problems in terms of managing and using them efficiently.

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

## Sentinel Data Rates

### If we just consider Sentinels 1 to 3



Sentinels (Incoming, TB/day)

Legend: Best Case, Worst Case



Sentinels (Volume in archive, PB)

Legend: stored_best, stored_worst

▶ Uncertainty in data rates grows with time. Not surprisingly we consider lower data growth rates "best" case scenarios.

▶ Uncertainty plays out in the archive with differences of petabytes in the volume stored in the not too distant future.

It's clear that we won't be storing all Sentinel data on disk, and while we have the tape (and library) capacity to store the Sentinel data, we don't yet have the software systems in place to make it easy for the user community to exploit a rolling cache for whole mission processing.

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○

The Data Commons
○○○○○

Looking Forward
○○○●○

Summary
○

The carrying capacity of the Commons

## Model Data Rates



PRIMAVERA data flows (courtesy Matthew Mizielinski)

CMIP6



- ▶ Typical modelling project: o(1 year) running primary simulations, but more data arriving up to two years.

- ▶ Data will be used for many years thereafter by both the original modelling community, and many others.

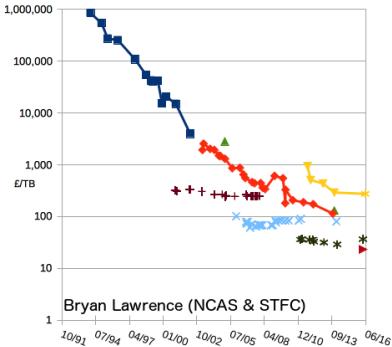- ▶ Data rates and volumes still unknown, but at least 10 PB over the 2017/2019 period, and possibly much much more.

- ▶ . . . as much of which needs to be on disk as we can manage!

National Centre for
Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○

The Data Commons
○○○○○

Looking Forward
○○○○●

Summary
○

The carrying capacity of the Commons

## Storage Costs

### The cost of storage at STFC over 25 years



Bryan Lawrence (NCAS & STFC)

(With thanks to Peter Chiu, Jonathan Churchill, and Tim Folkes)

### Kryder's Law is definitely slowing!

► Disk (Blue and red lines)
► Parallell Disk (Yellow lines)
► Tape Generations show as unconnected points (often same tapes, different drives!)
► Tape is likely to be cheaper for the foreseeable future (disk technology advances slowing rapidly)
► (The worry is that market forces may drive tape and even disk into extinction!)

Note that as the volume increases, the cost of software to manage the large volume storage needs to be added to the raw cost of disk.

We might be able to move to object stores, which will bring our raw cost per TB down, but we're all going to have to learn to live without file systems!

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

Rising demand
○○○○○○○○○○
From demand to Futures

The Data Commons
○○○○○

Looking Forward
○○○○○

Summary
●

## Drivers and Choices

Massive increase in data
production driven by Moore's Law
from

- ▶ Growth in the use of simulation,
- ▶ Growth in resolution and complexity of simulation,
- ▶ Growth in the number and frequency of observations,

Compounded by

- ▶ Complex interactions between communities

Resulting in massive increase in

- ▶ Number and volume of datasets in the CEDA archive,
- ▶ Number of communities, and volumes of data they need to handle, and
- ▶ Interactions and the need for quality documentation of data and computational provenance.

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

Why JASMIN?
Bryan Lawrence - RAL, June 2016

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

## Drivers and Choices

Massive increase in data production driven by Moore's Law from

- ▶ Growth in the use of simulation,
- ▶ Growth in resolution and complexity of simulation,
- ▶ Growth in the number and frequency of observations,

Compounded by

- ▶ Complex interactions between communities

Resulting in massive increase in

- ▶ Number and volume of datasets in the CEDA archive,
- ▶ Number of communities, and volumes of data they need to handle, and
- ▶ Interactions and the need for quality documentation of data and computational provenance.

All this accompanied by a complex and varied computational requirements:

- ▶ People need access to data, and to share data, and to provide their own services on their data.
- ▶ Data volumes are getting so large that parallelised workflows are necessary, but it's hard to build "throw-away parallel codes".
- ▶ "Bringing compute to the data" requires support for complex computational environments.

Leading to

- ▶ Complex infrastructure requirements, the death of old software friends (e.g filesystems), the rise of new paradigms (containerisation), and
- ▶ Much new learning for us, the community of scientists who want exploit data to address society's pressing problems!

**National Centre for Atmospheric Science**
NATURAL ENVIRONMENT RESEARCH COUNCIL

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