# Why Pangeo?

## What is it and why we need it

Theo McCaie

**Met Office** Informatics Lab

Grains of sand

A mug full
(25th square)

A teaspoon full
(16th square)

The volume of a
person
(33rd square)

In a cargo
container
(39th square)

5th largest container
ship in the world
(53rd square)

**2x**
the worldwide
container
ship capacity
(64th square)

| 1 | 2 | 4 | 8 | 16 | 32 | 64 | 128 |
| 256 | 512 | 1024 | 2048 | 4096 | 8192 | 16384 | 32768 |
| 65536 | 131K | 262K | 524K | 1M | 2M | 4M | 8M |
| 16M | 33M | 67M | 134M | 268M | 536M | 1G | 2G |
| 4G | 8G | 17G | 34G | 68G | 137G | 274G | 549G |
| 1T | 2T | 4T | 8T | 17T | 35T | 70T | 140T |
| 281T | 562T | 1P | 2P | 4P | 9P | 18P | 36P |
| 72P | 144P | 288P | 576P | 1E | 2E | 4E | 9E |

Generating more data in a week than in total 7 years ago

# Total market capitalisation: Cars Vs Computers

hello@informaticslab.co.uk    @informatics_lab    informaticslab.co.uk

# We are looking for a way to:

- Help people make better decisions based on complex data?

- Allow bespoke solutions to domain specific problems?

- Empower our ever more data-literate consumers?

- Unlock all the value in hundreds of PB of data?

- Provide improved accuracy without ever more flops?

- Empower analysts to do what they want not what they can?

# Iris

- Open Source (BSD license)

- Encapsulates Dask for distributed calculations

- Re-gridding/projection/interpolation

- Unit conversion

- Reads/converts various formats (Grib, NetCDF, fieldsfiles...)

- Automatic plotting via matplotlib and holoviews
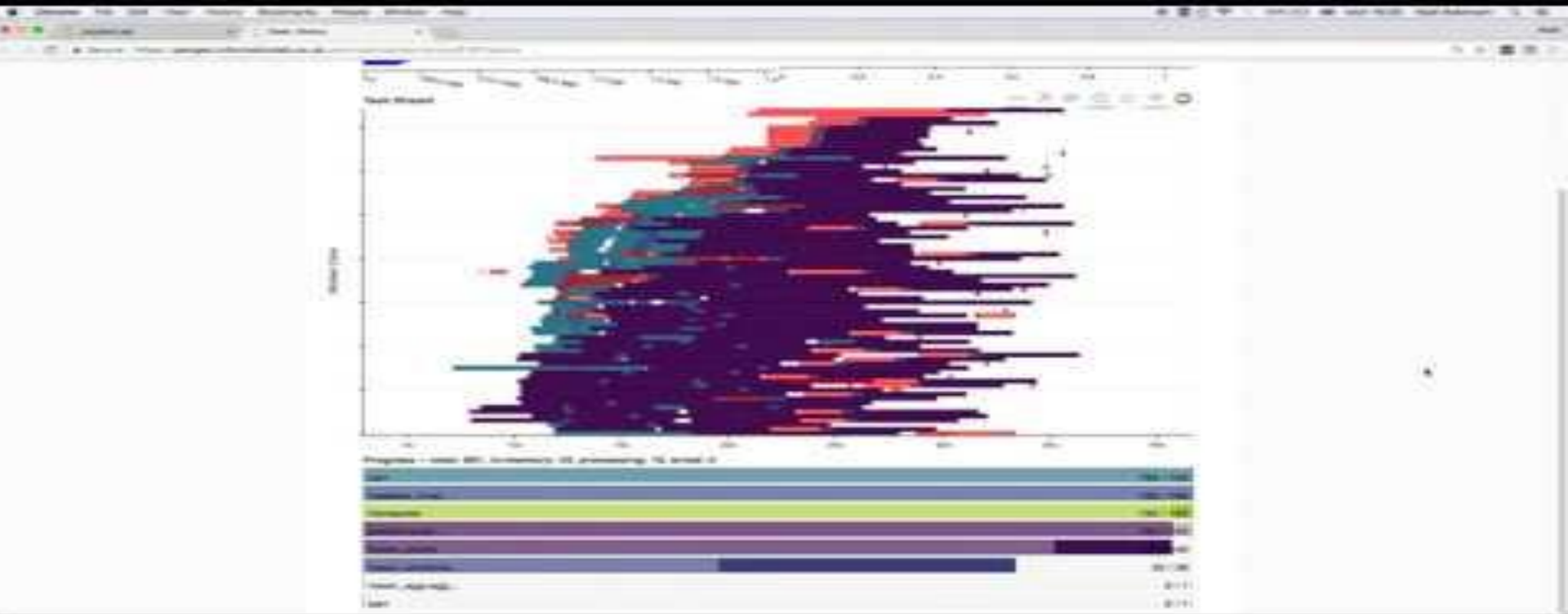
- 8 FTEs at the Met Office working on it!

# Pangeo



**=**

- Responding to demand elastically

- Interactive analysis to encourage "flow", not fire-and-forget batch jobs

- Laziness/just-in-time

- Thin web client views to interact with data

- Agile, bespoke, product creation

https://bit.ly/2O9qJr3

# On the cloud?

- Doesn't have to be but...

# Workloads are volatile by nature

On Prem this gives two options:

- Very big cluster
- quick results
- Inefficient utilisation

vs

- Smaller cluster
- suppress volatility with queues
- higher utilisation

# Volatile workloads in the cloud

1. Scheduler creates many jobs for an individual user

2. Many schedulers submit to the same orchestrator, smoothing volatility somewhat

3. Orchestrator asks for more cloud resources in response to spikes in demand

4. Many Many users use the cloud smoothing demand a lot

5. The Cloud installs more racks in response to demand

Fast

Slow

# What's next

- Data discovery - "Hey pangeo, get me data on UK storms last year"

- Science to service - one click APIs and dashboards.

- Environmental data analysis for all - deployments, tools, APIs and tutorials to make environmental data accessible to SMEs

- On demand hardware configured clusters

- Improved data storage and handling

# What's next: our data challenge

- Fast metadata access

- Parallel and elastic over fast and immediately consistent

- Universally addressable

- Reluctant to embrace one size fits all

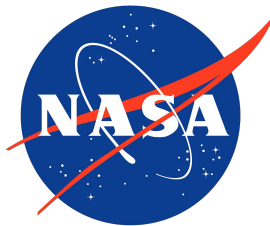- Still need to deal with the high volume and velocity data as is currently generated

**Zarr**

# Why pangeo?

# Innovation is hard... and always has been

*introduction of innovation...is a struggle against stupidity and envy, apathy and evil, secret opposition and open conflict of interests, a horrible period of struggle with man, a martyrdom even if success ensues.*
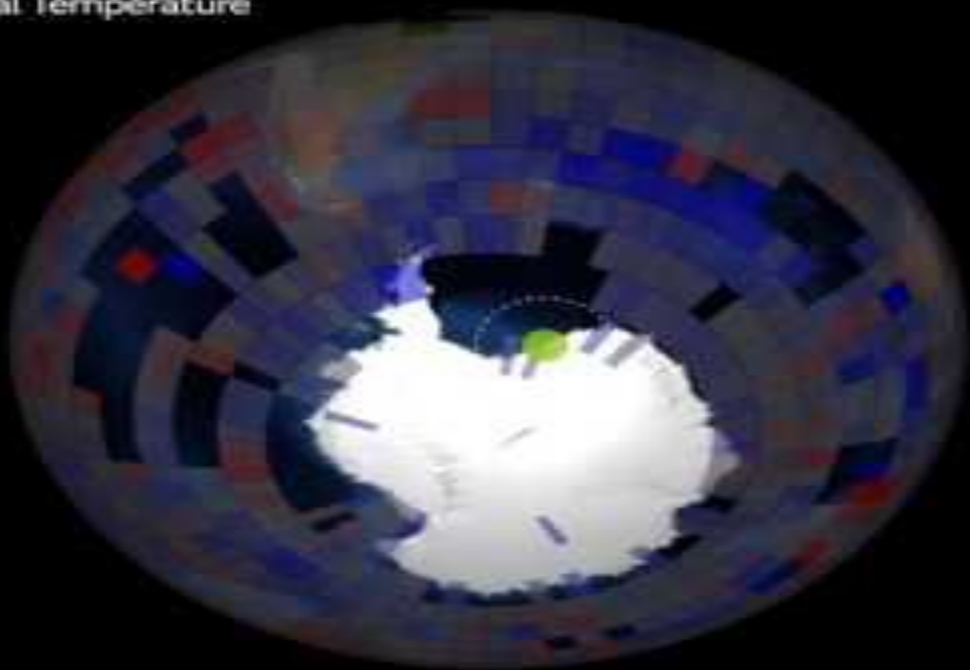
Diesel, 1858-1913

**Met Office**   **INFORMATICS LAB**

# Thank you!

✉ hello@informaticslab.co.uk     🐦 @informatics_lab     🌐 informaticslab.co.uk