

# BESS Summer School

## Argo Quality Control Analysis

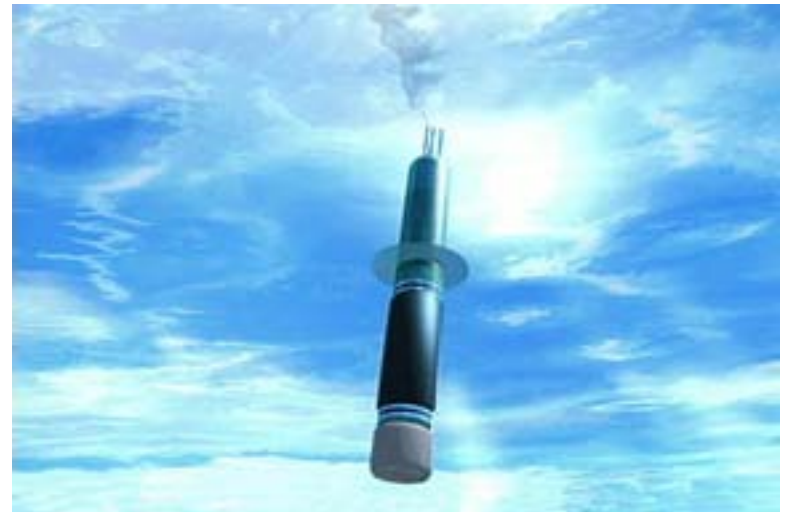
**Supervisors:** Robert Gurney, Francois Taiani

### **Analysis sub-group**

- Alastair Gemmell
- Jennine Jonczyk
- Dusan Kostic
- Ulirich Weber

### **Visualisation sub-group**

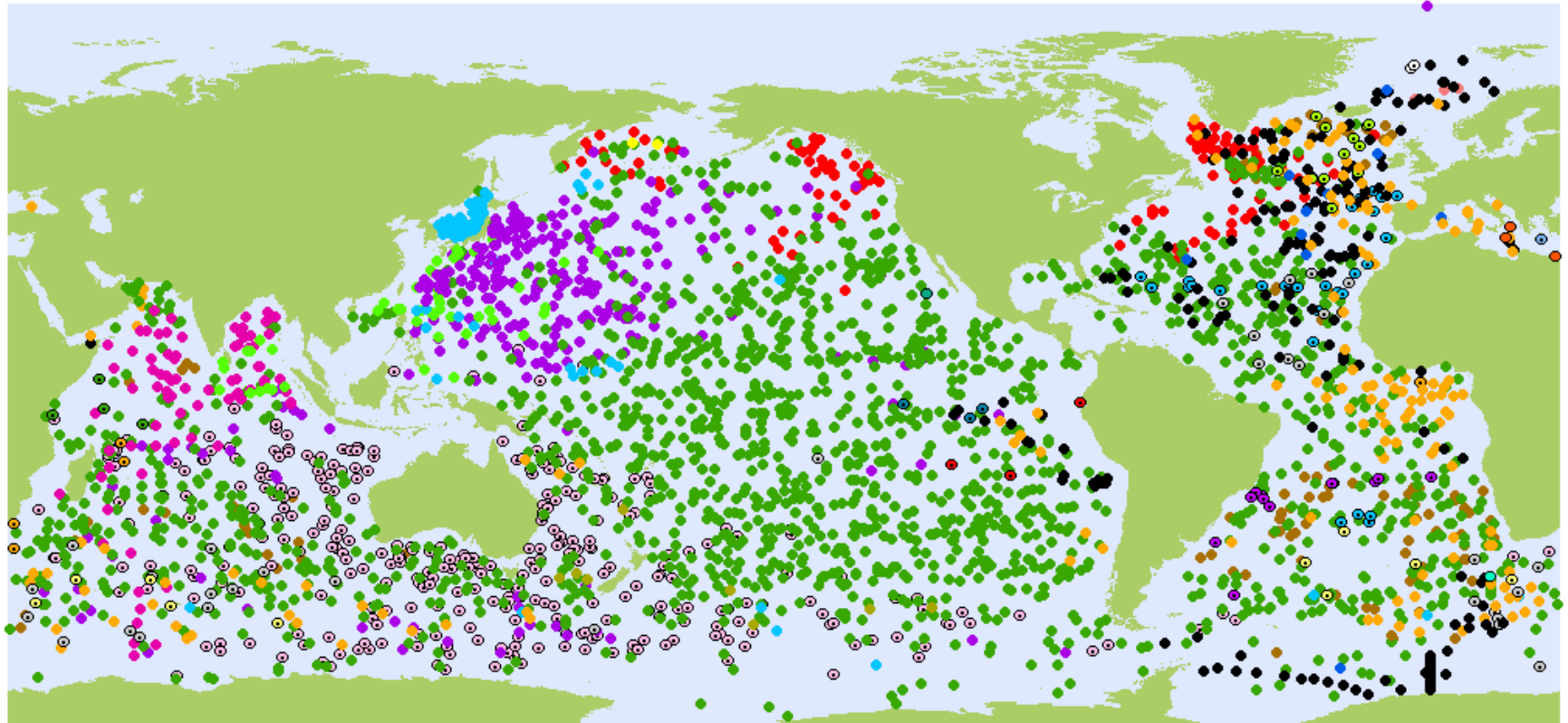
- Radoslaw Guzinski
- Fei Ma
- Eleanor Mackay
- Mark Wilkinson



# What are Argo floats?

- Autonomous profiling floats. Phased in since ~2000
- Goal of 3000 floats was achieved by 2007
- 10 day cycle
  - Mostly float at (up to) 2000m depth
  - Every 10 days ascend, collecting data on Temperature and Salinity at frequent depth intervals
  - At surface transmit latest profile data via satellite to operational centres and Argo data centres
- Paradigm shift in sampling of T and S in the ocean
- Undergo real time QC 'quick and dirty'
  - This is updated months later by the results of detailed DELAYED MODE QC. This is what we are looking at

# Where are Argo floats?



3301 Argo Floats

○ ARGENTINA (1 0)	● CHINA (46)	● GABON (1)	● ITALY (3)	● MEXICO (1)	● RUSSIAN FEDERATION (2)	● UNITED STATES (1 750)
○ AUSTRALIA (331)	● ECUADOR (3)	● GERMANY (185)	● JAPAN (287)	○ NETHERLANDS (32)	● SAUDI ARABIA (0)	
● BRAZIL (1 2)	● EUROPEAN UNION (1 0)	● GREECE (1)	● KENYA (4)	● NEW ZEALAND (8)	● SOUTH AFRICA (1)	
● CANADA (1 29)	○ FINLAND (2)	● INDIA (82)	● SOUTH KOREA (85)	● NORWAY (3)	● SPAIN (2 7)	
● CHILE (3)	● FRANCE (1 65)	● IRELAND (1 2)	● MAURITIUS (4)	○ POLAND (0)	● UNITED KINGDOM (1 02)	

May 2011

# The data for this project

## ■ Argo Data

- Detailed delayed mode Quality Control.
- Profiles are coded as containing percentages of **Good Data**:
  - Code **0**: GD = **0**
  - Code **1**: **0** < GD <= **25**
  - Code **2**: **25** < GD <= **50**
  - Code **3**: **50** < GD <= **75**
  - Code **4**: **75** < GD < **100**
  - Code **5**: GD = **100**

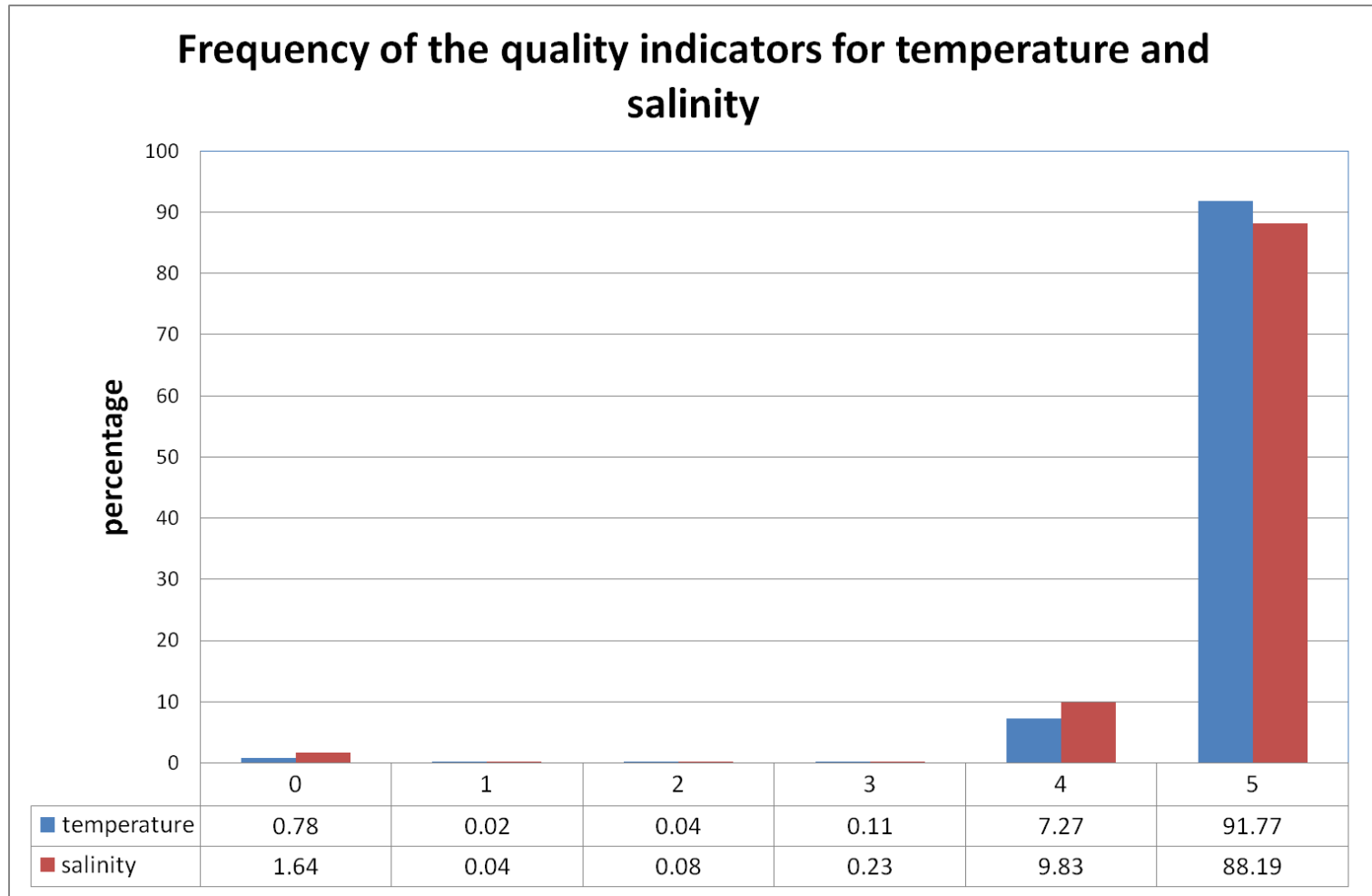
## ■ Accept/Reject decisions from operational centres

- Courtesy of Jim Cummings at US Navy via the GODAE project
- Available for:
  - BMRC: Australia's **B**ureau of **M**eteorology **R**esearch **C**entre
  - FNMOC: US Navy's **F**leet **N**umerical **M**eteorology and **O**ceanography **C**enter
  - MEDS: Canada's **M**arine **E**nvironment **D**ata **S**ervice
  - UKMO: **UK**'s **M**et **O**ffice

# Motivation

- Data assimilation extremely important for accuracy of forecasts
- Argo floats important source of these data
- Operational centres must QUICKLY decide to accept or reject recent Argo profiles
- They don't have time to wait for detailed QC done by Argo project
- So...how accurately can they detect profiles which (upon detailed analysis) are later shown by the Argo project to contain bad data?

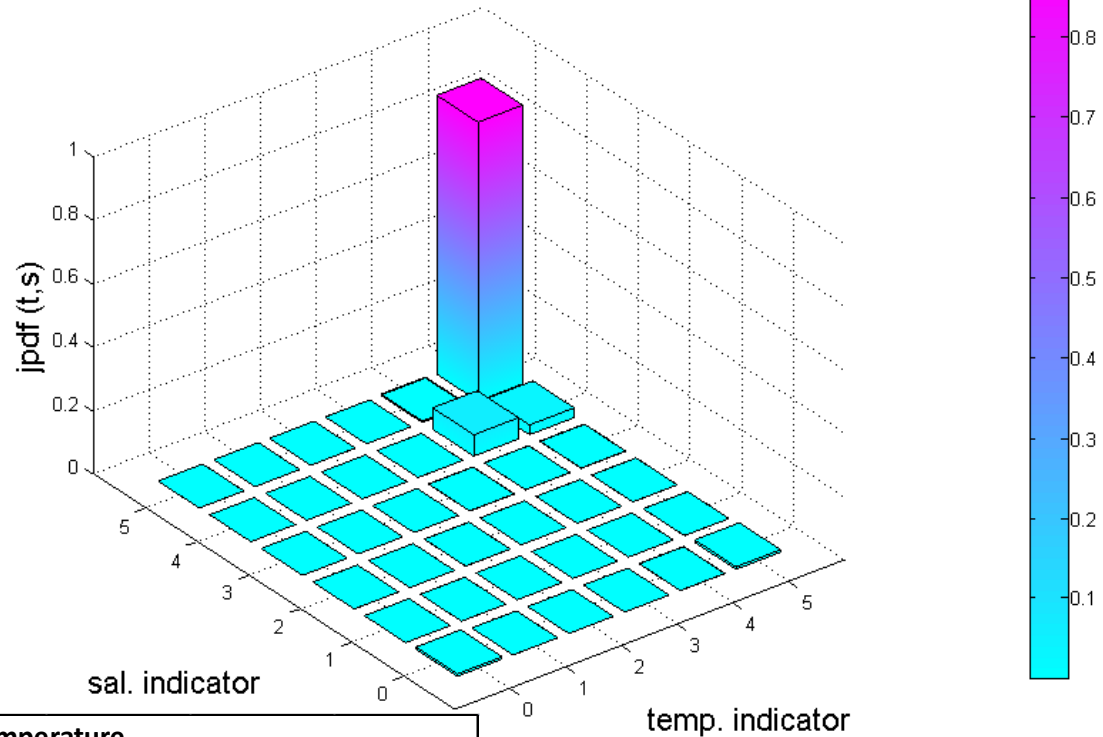
# Distribution of quality indicators



# Indicators as dependent variables

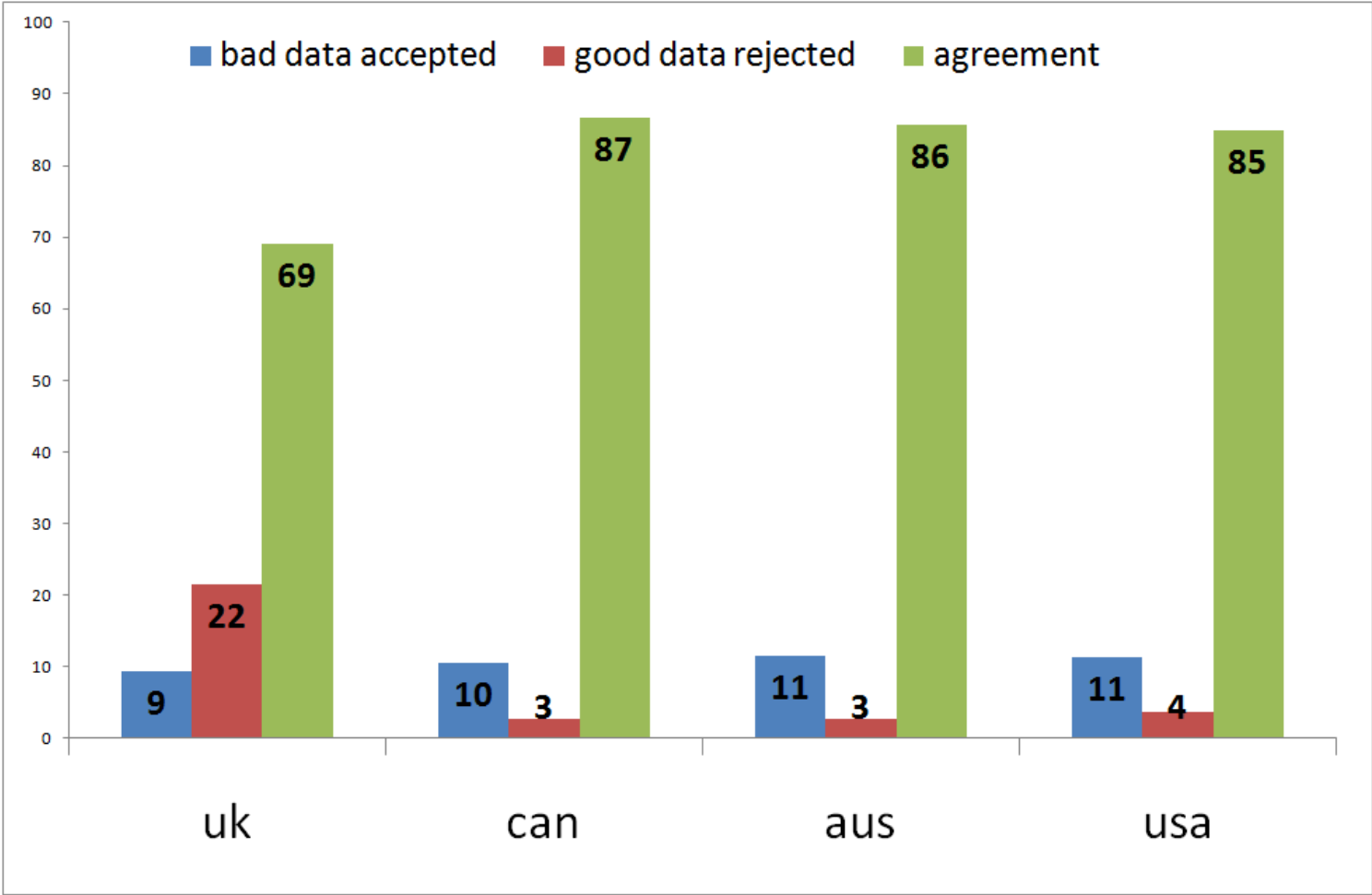
Empirical Joint Distribution for T&S indicators

		0	1	
salinity	0	1155	4	
	1	1	22	
	2	8	0	
	3	2	0	
	4	9	0	
	5	29	0	
		No	153350	

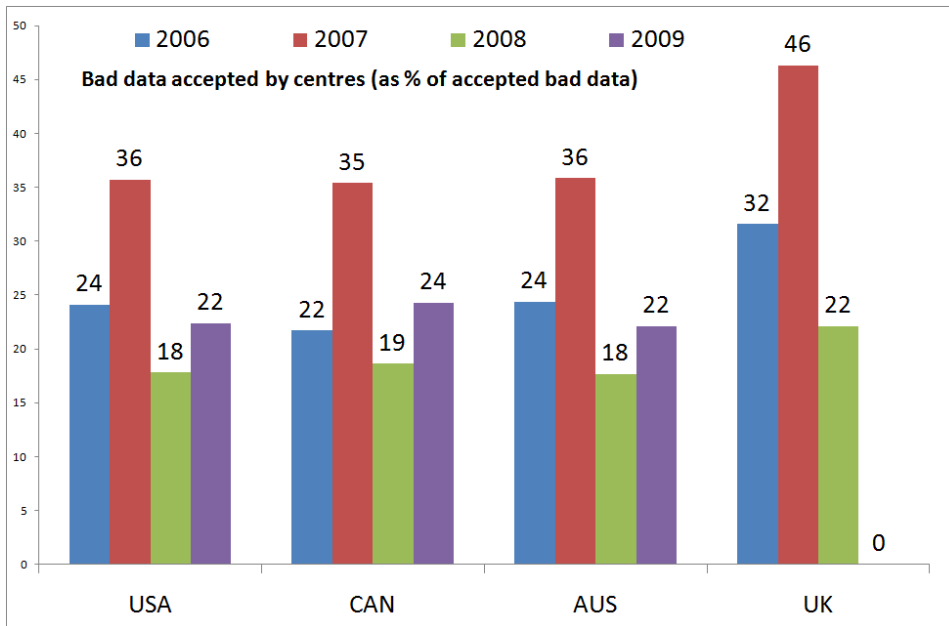


		temperature					
		0	1	2	3	4	5
salinity	0	7.53E-03	2.61E-05	2.61E-05	3.91E-05	8.35E-04	7.95E-03
	1	6.52E-06	1.43E-04	6.52E-06	0.00E+00	5.22E-05	1.76E-04
	2	5.22E-05	0.00E+00	3.65E-04	1.96E-05	5.22E-05	3.59E-04
	3	1.30E-05	0.00E+00	1.30E-05	9.72E-04	1.96E-04	1.06E-03
	4	5.87E-05	0.00E+00	6.52E-06	3.91E-05	6.71E-02	3.08E-02
	5	1.89E-04	0.00E+00	3.26E-05	3.91E-05	4.40E-03	8.77E-01
sum			1.0				

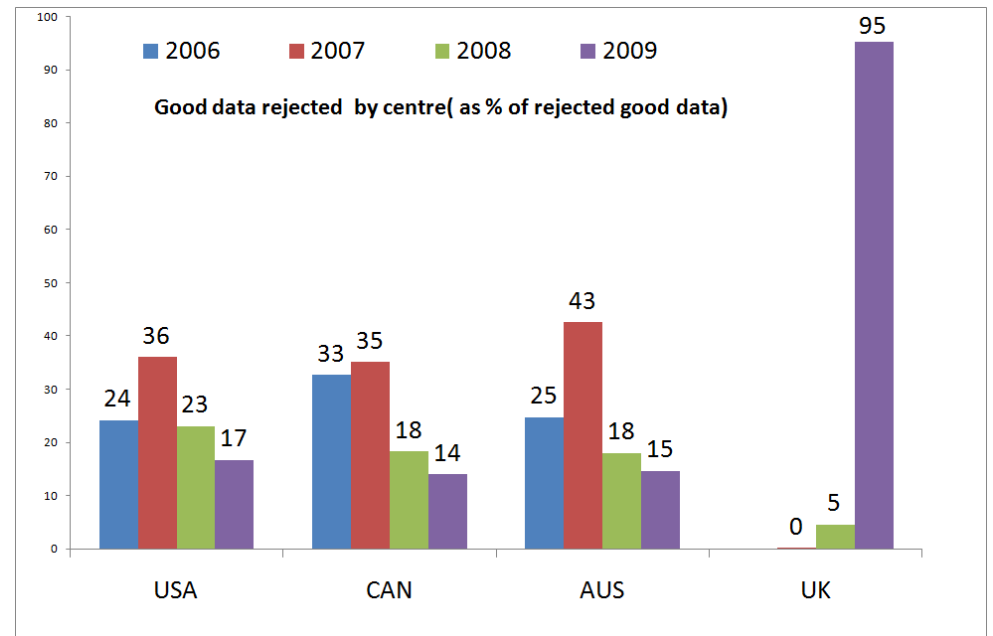
# First Look.....



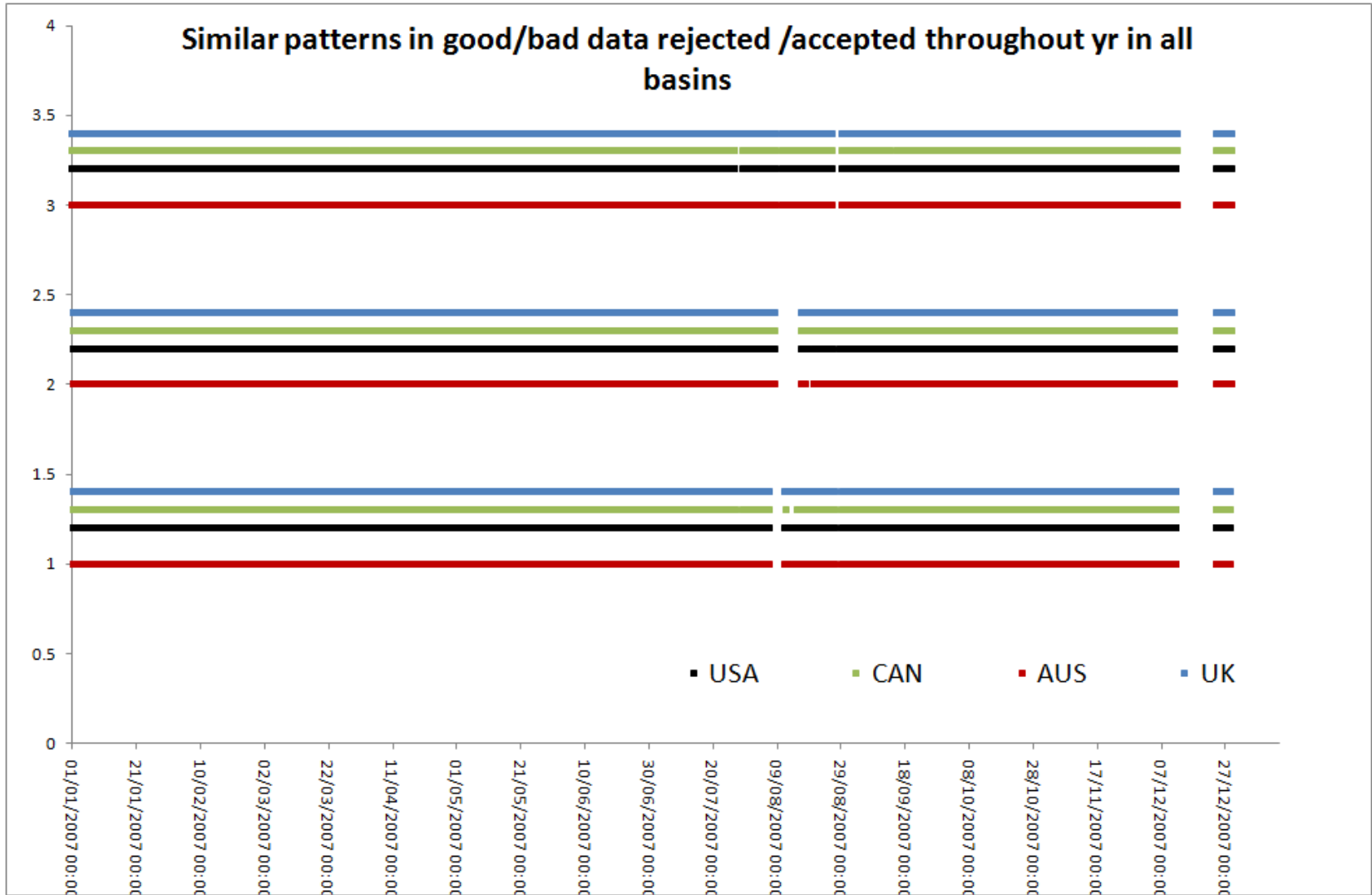


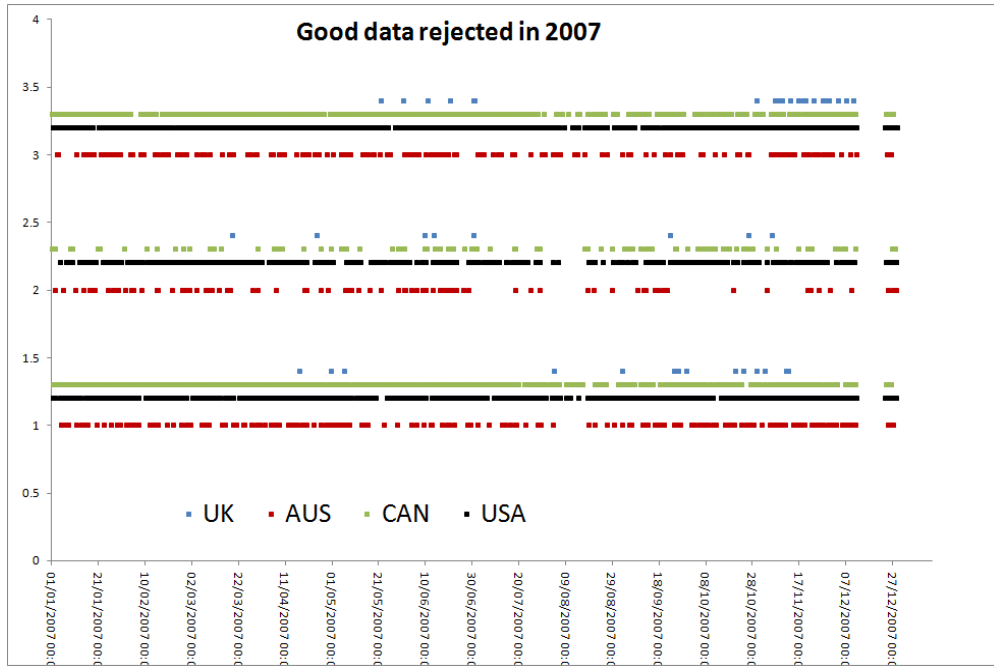


■ Q: is it worse to accept bad data or reject good data?



# Any particular times?

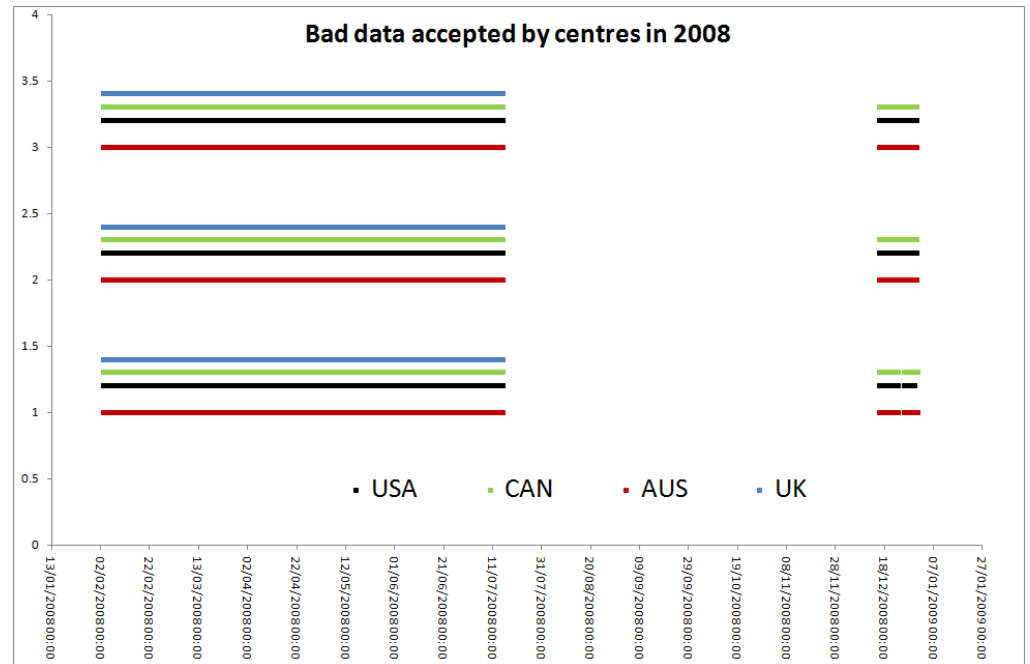




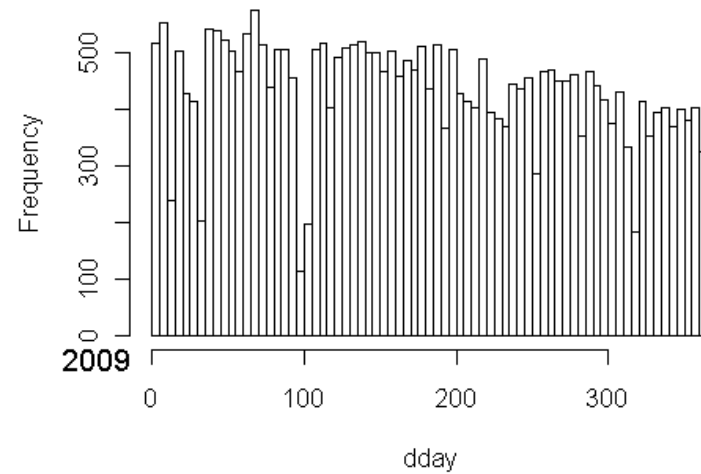
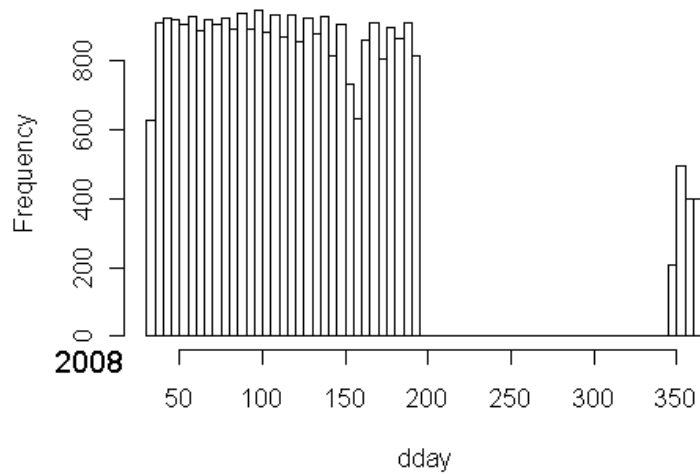
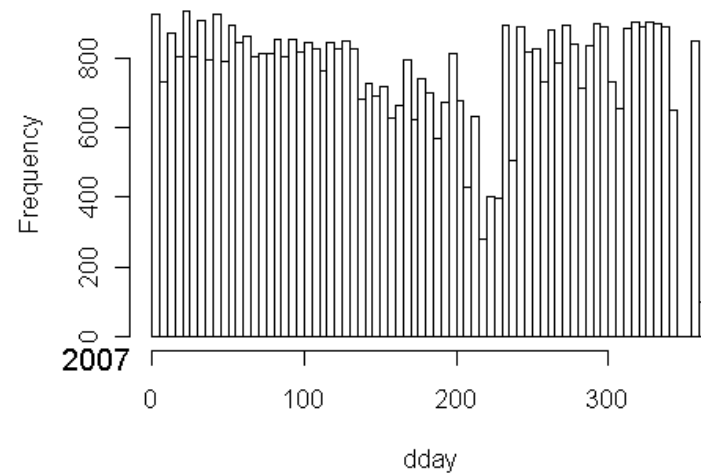
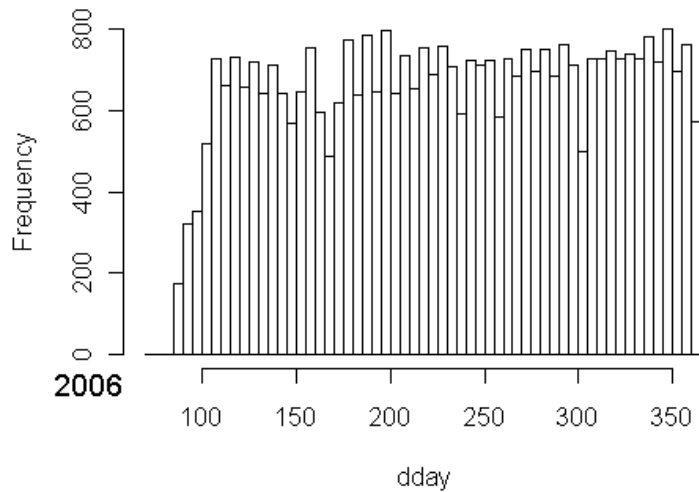
EXCEPT:

Good data rejected in 2007

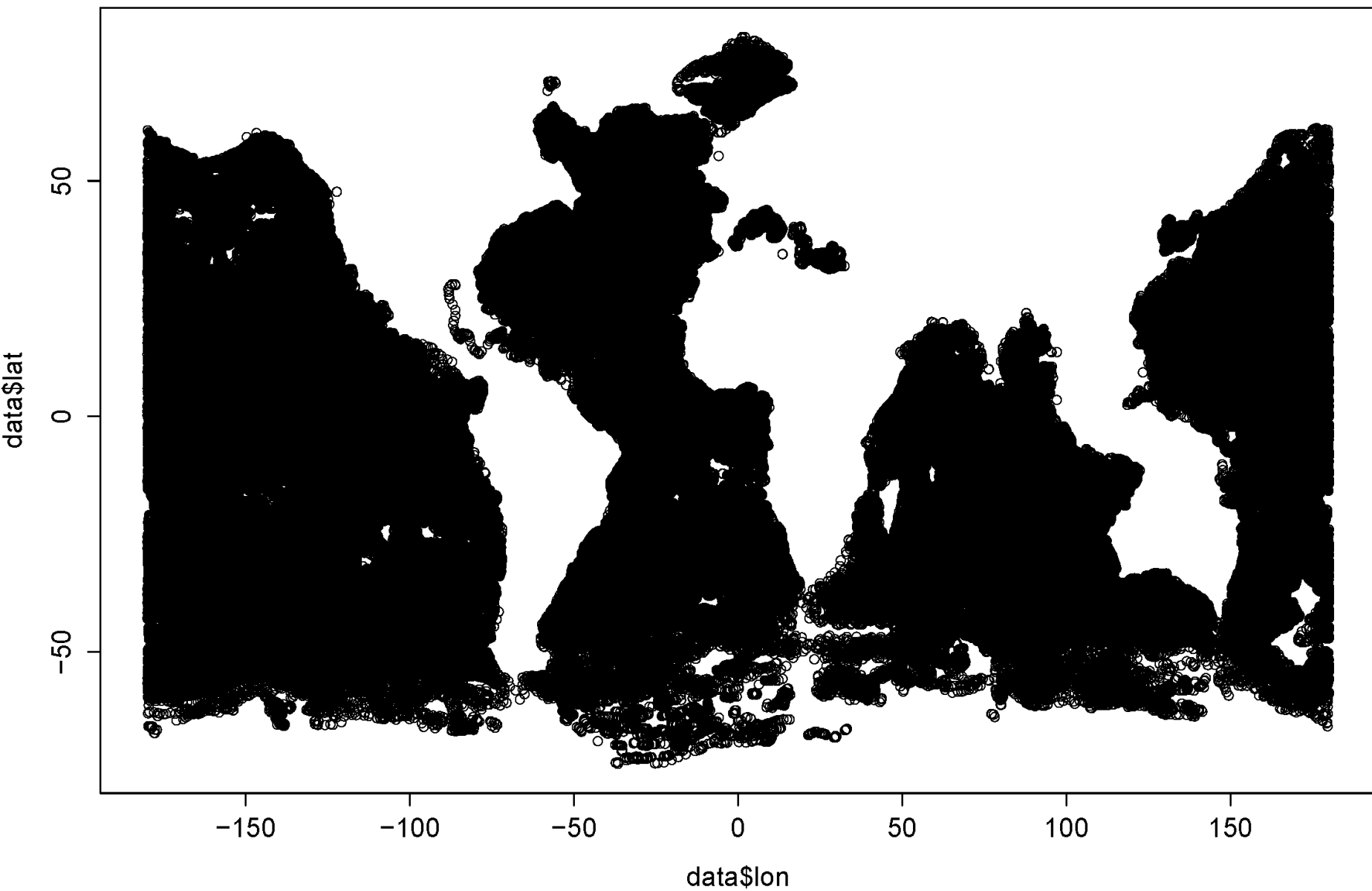
Bad data accepted in 2008



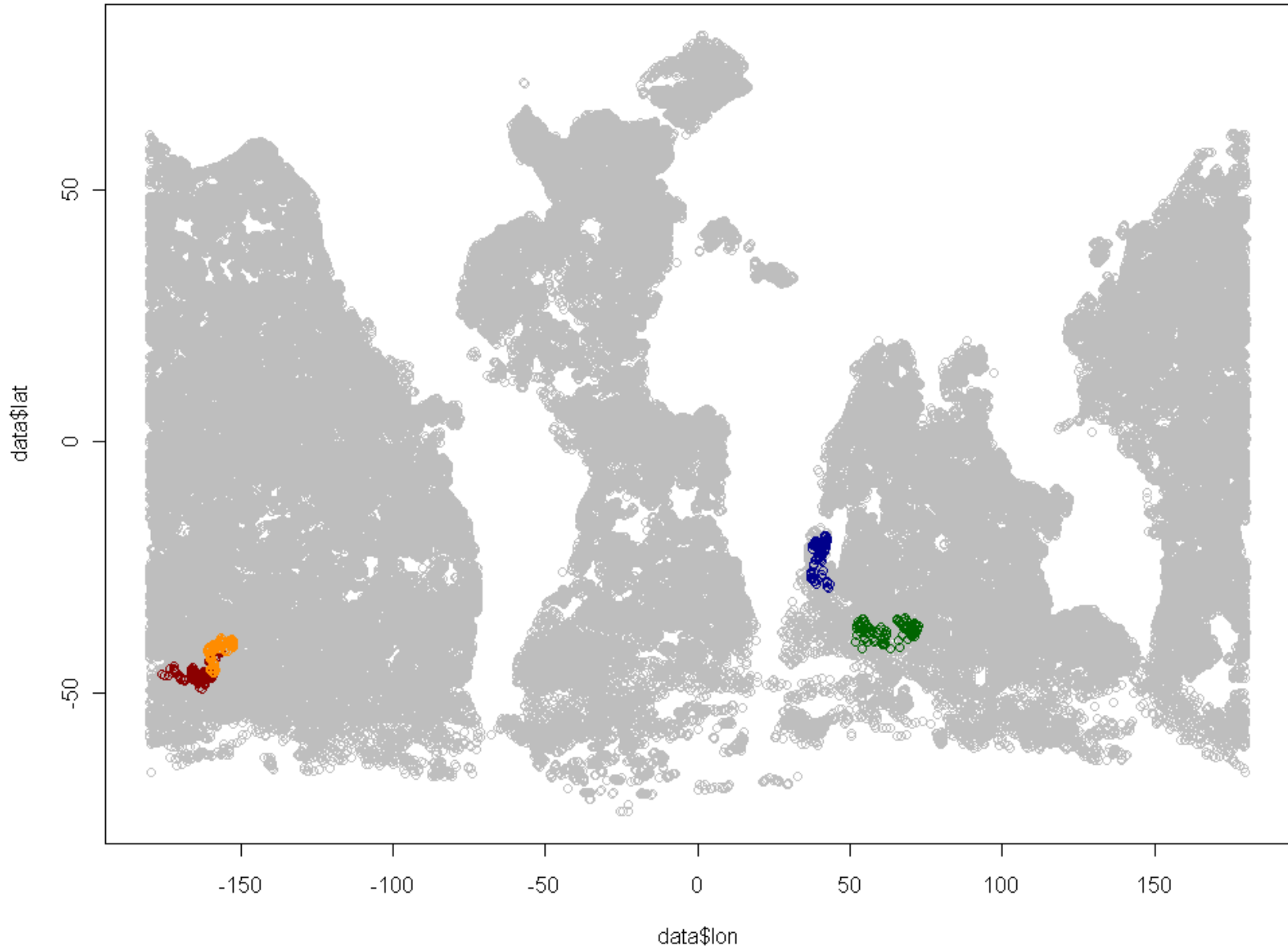
## Data availability in days/year (bin=3)



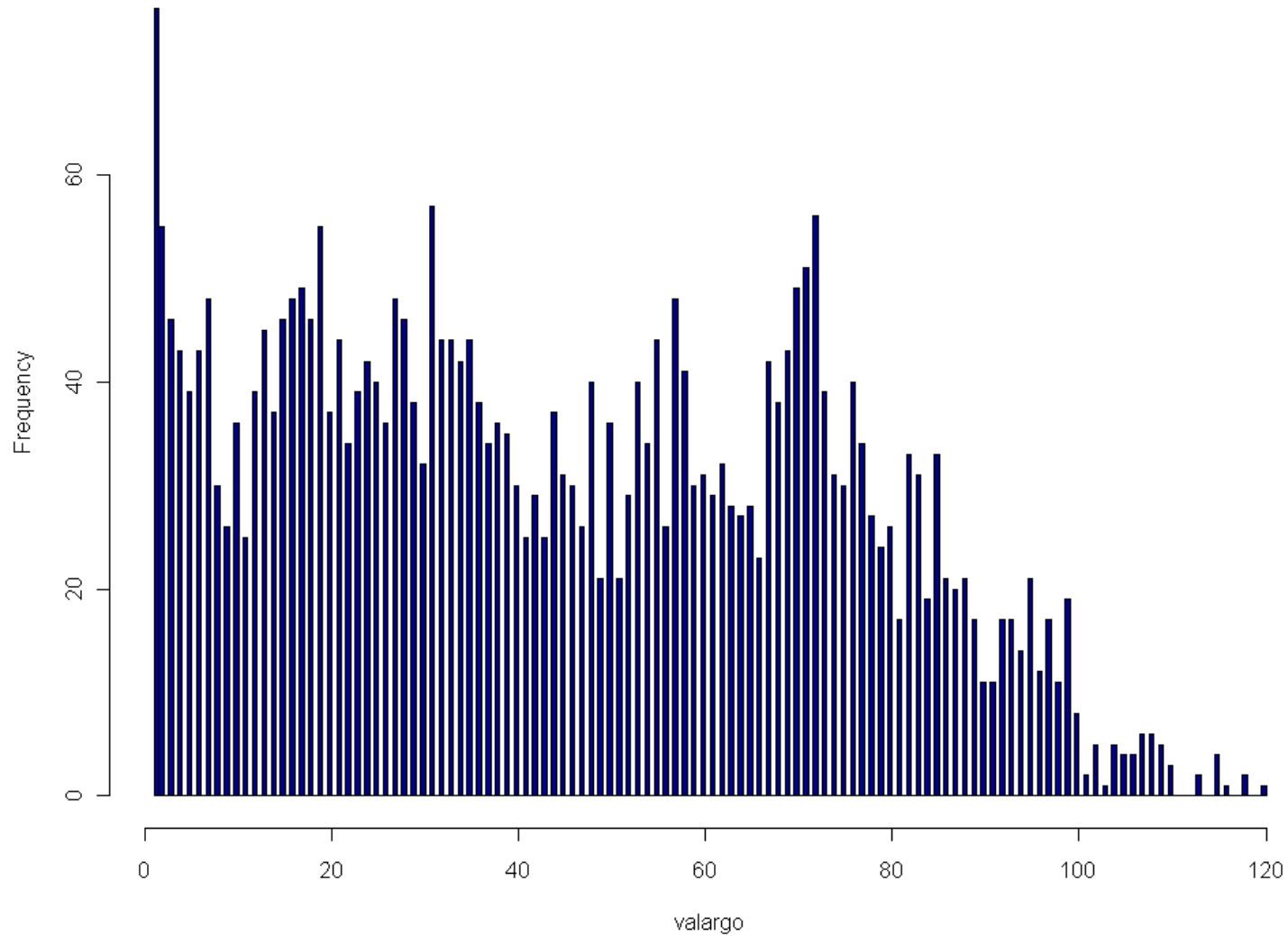
## *Spatial extend of QC dataset*



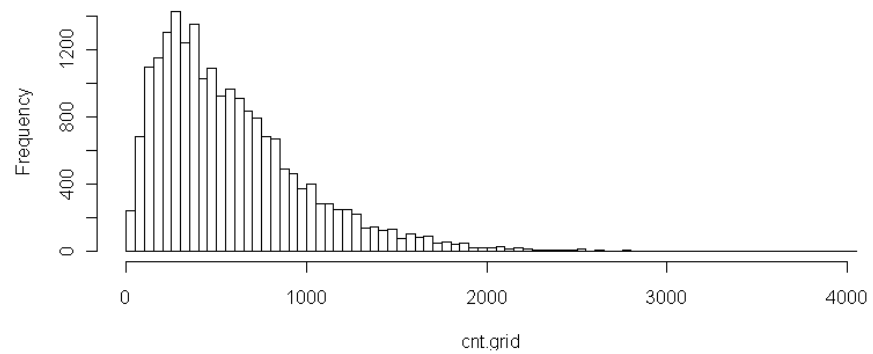
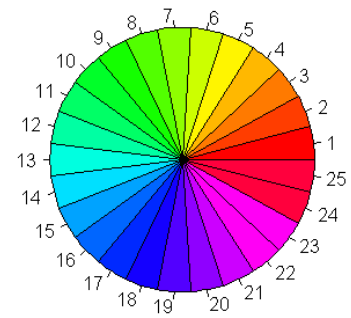
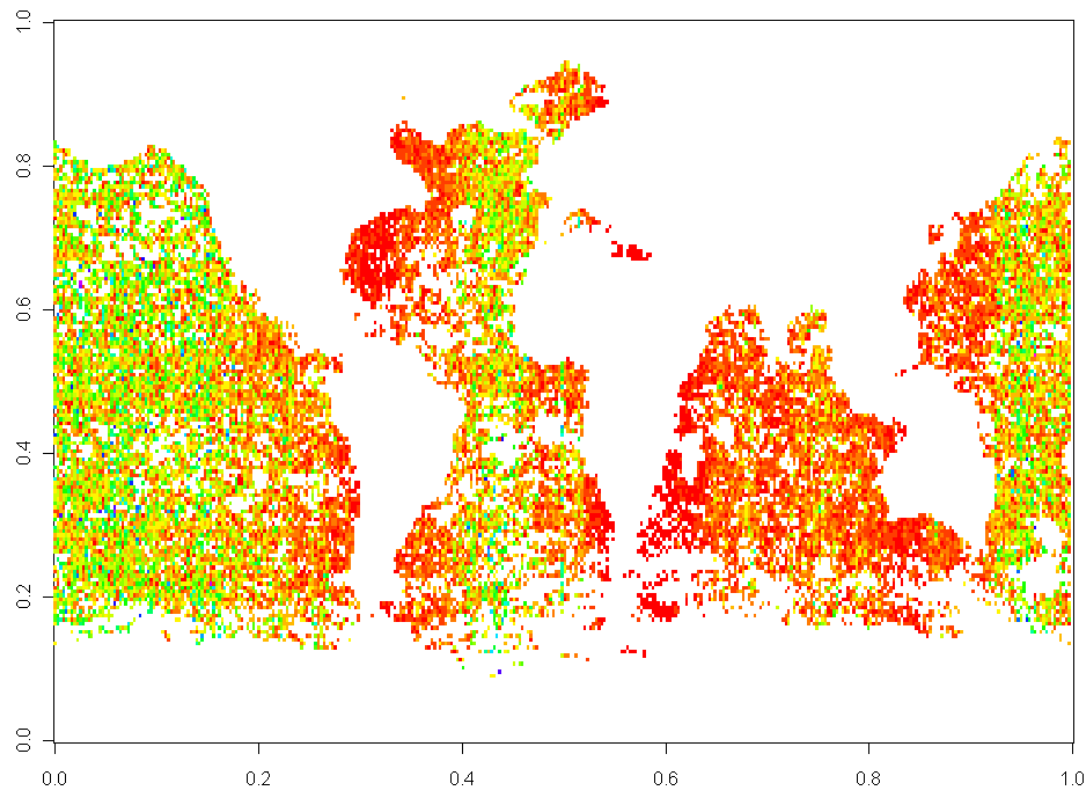
***Spatial extend of all.valid.argo.float (-999 excluded) + coloured paths of sel.argo.float***



## *Histogram of Valid Argo.floats per measurement*



# 1 deg grid Cumulative Argo Counts

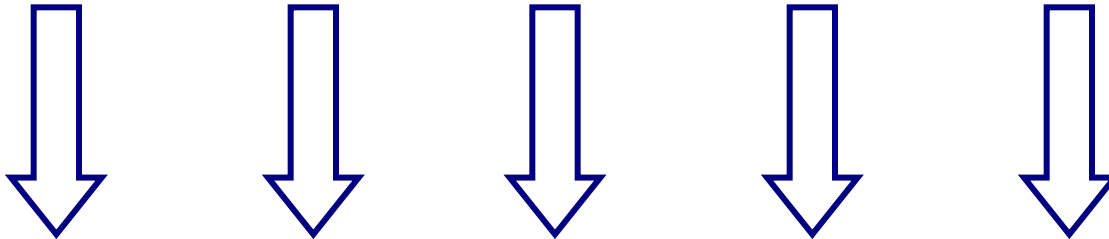




# Clustering

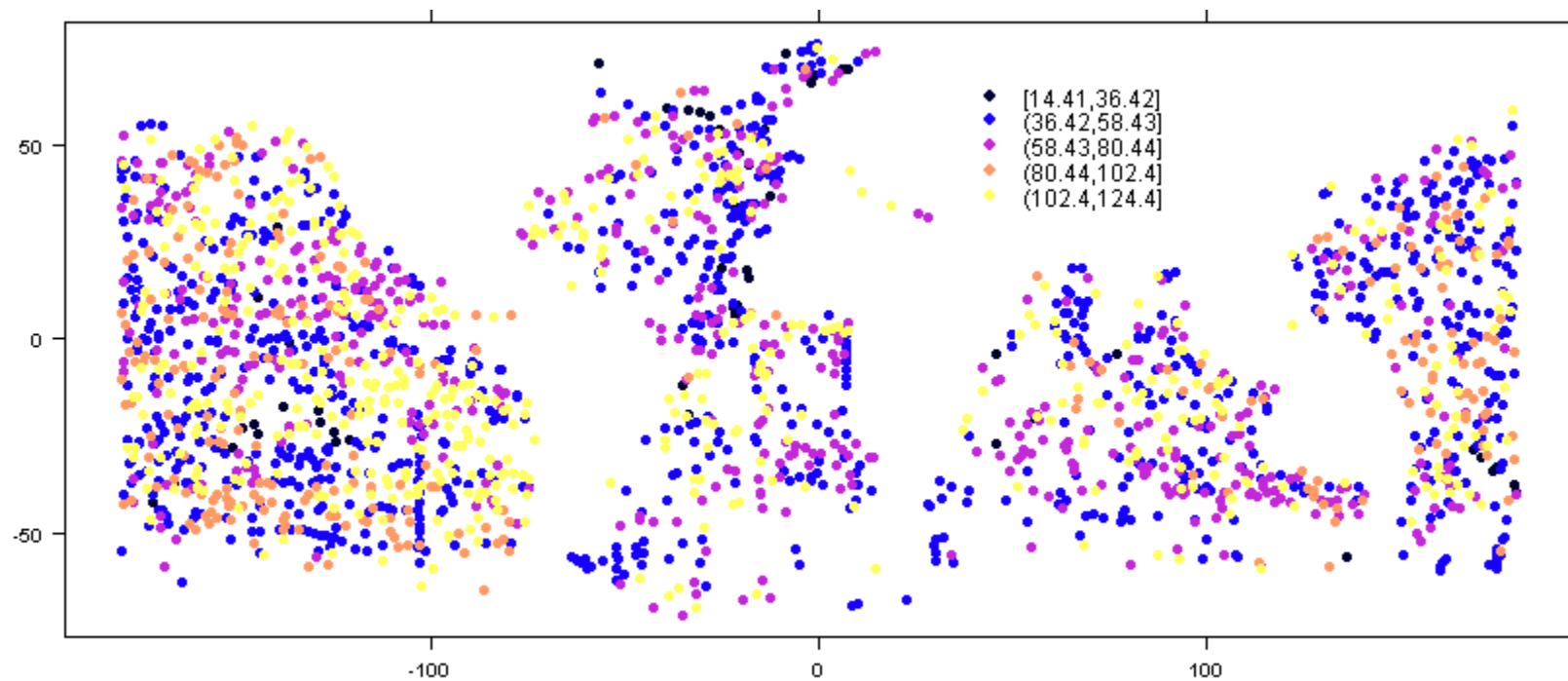
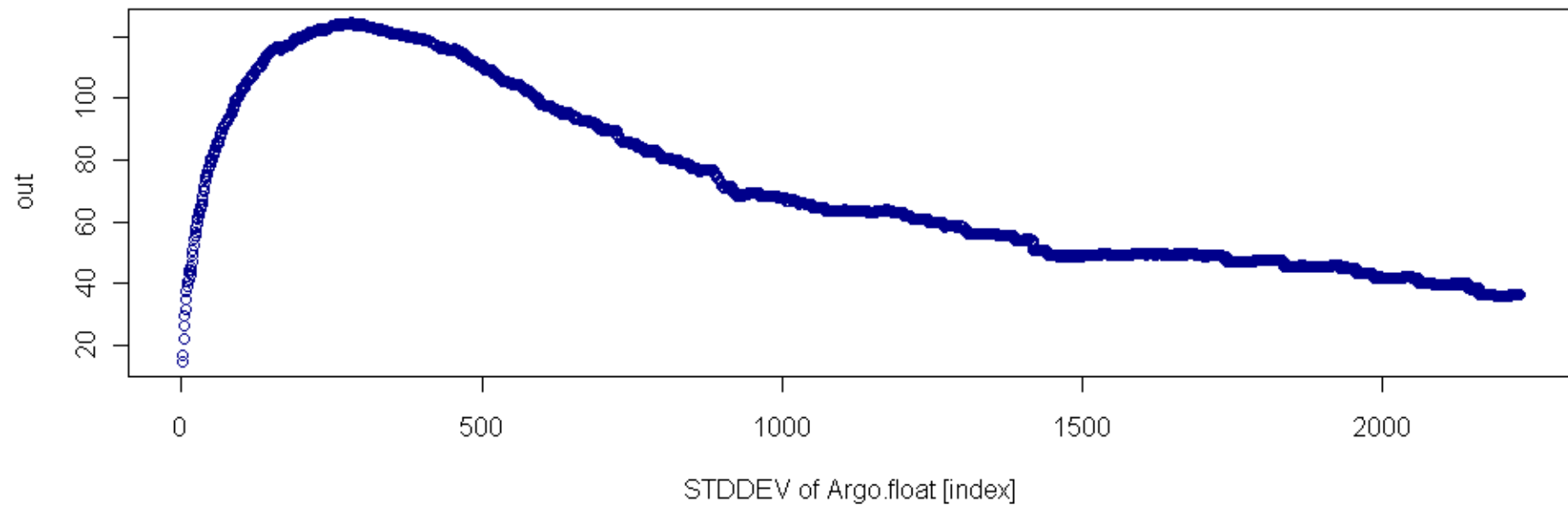
	Day 1	Day 2	...	...	Day 365
<i>d_mode_qc</i>	0	0	0	0	0
<i>bmrc-qc</i>	55	0	0	55	55
<i>fnmoc_qc</i>	0	0	0	0	0
<i>meds_qc</i>	55	50	0	55	0
<i>ukmo_qc</i>	0	0	0	0	0

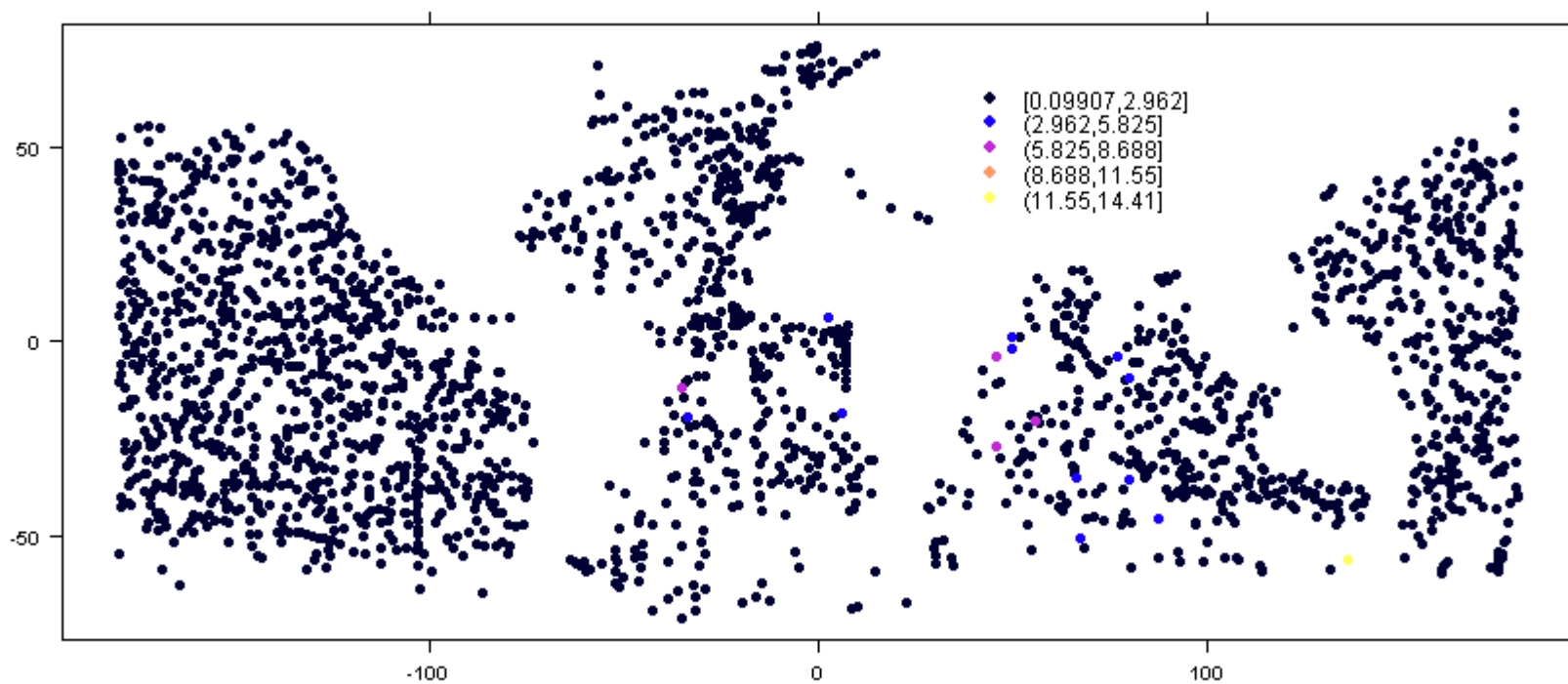
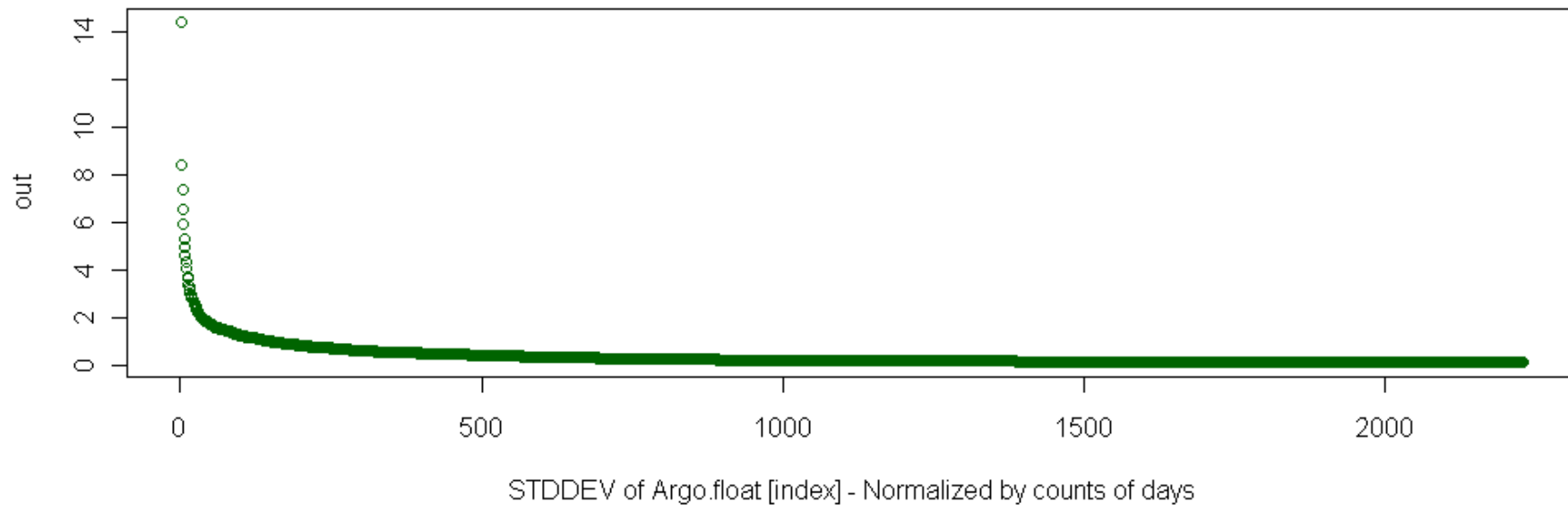
**TOTAL, 2**



	110	50	0	110	55
--	-----	----	---	-----	----

**STDDEV(TOTAL, 2)**

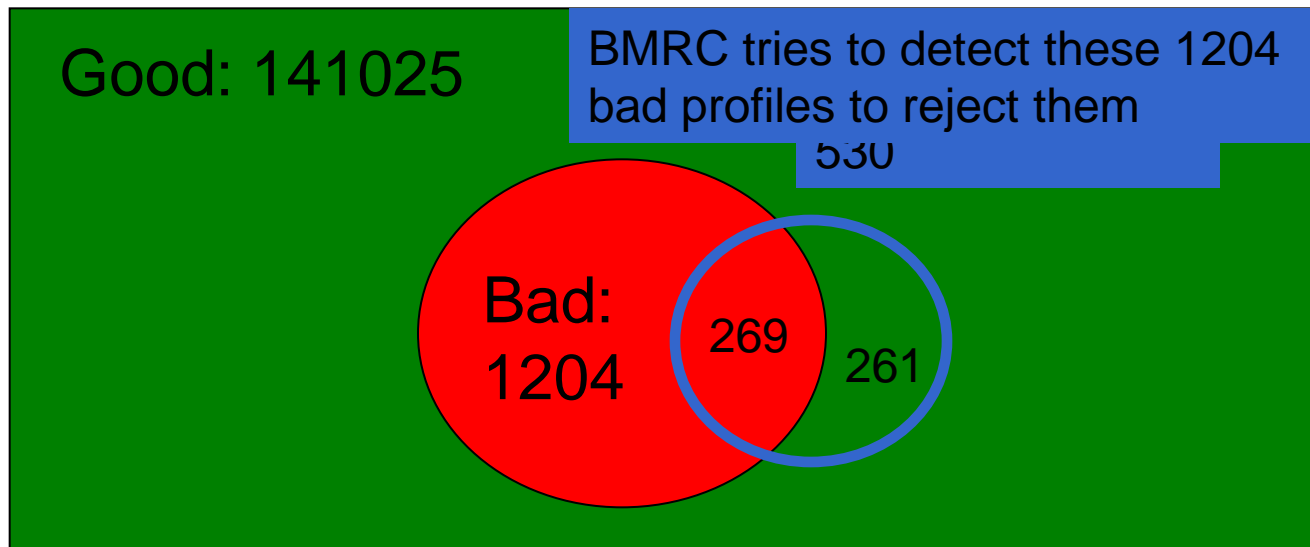




# Recall and Precision

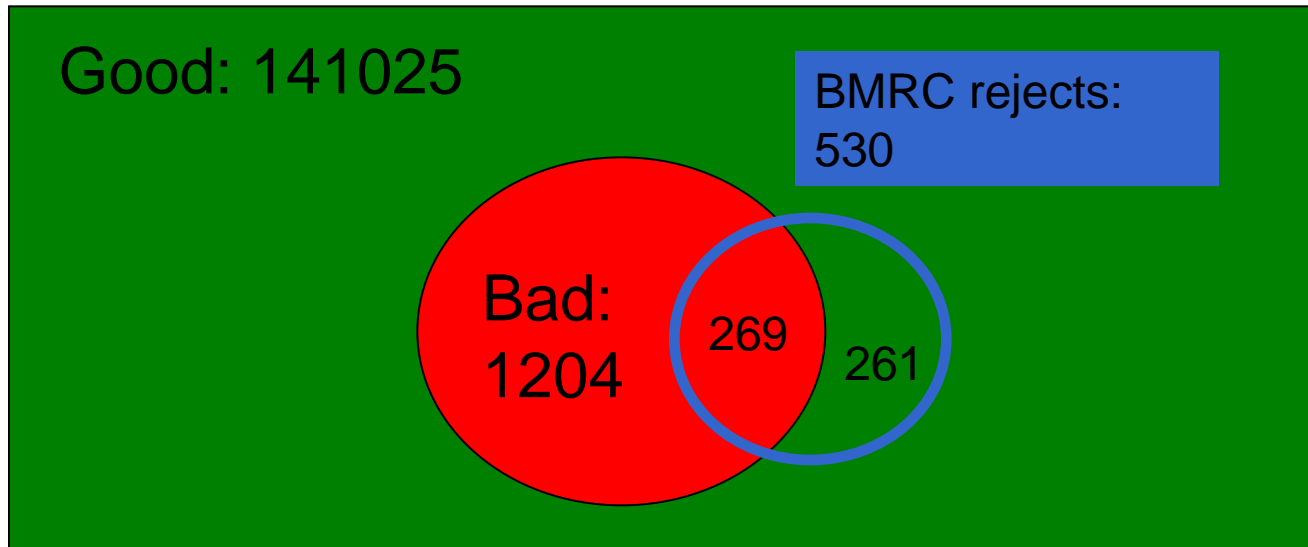
- Metrics to measure degree of success of the operational centres in detecting the bad data

Total: 142229



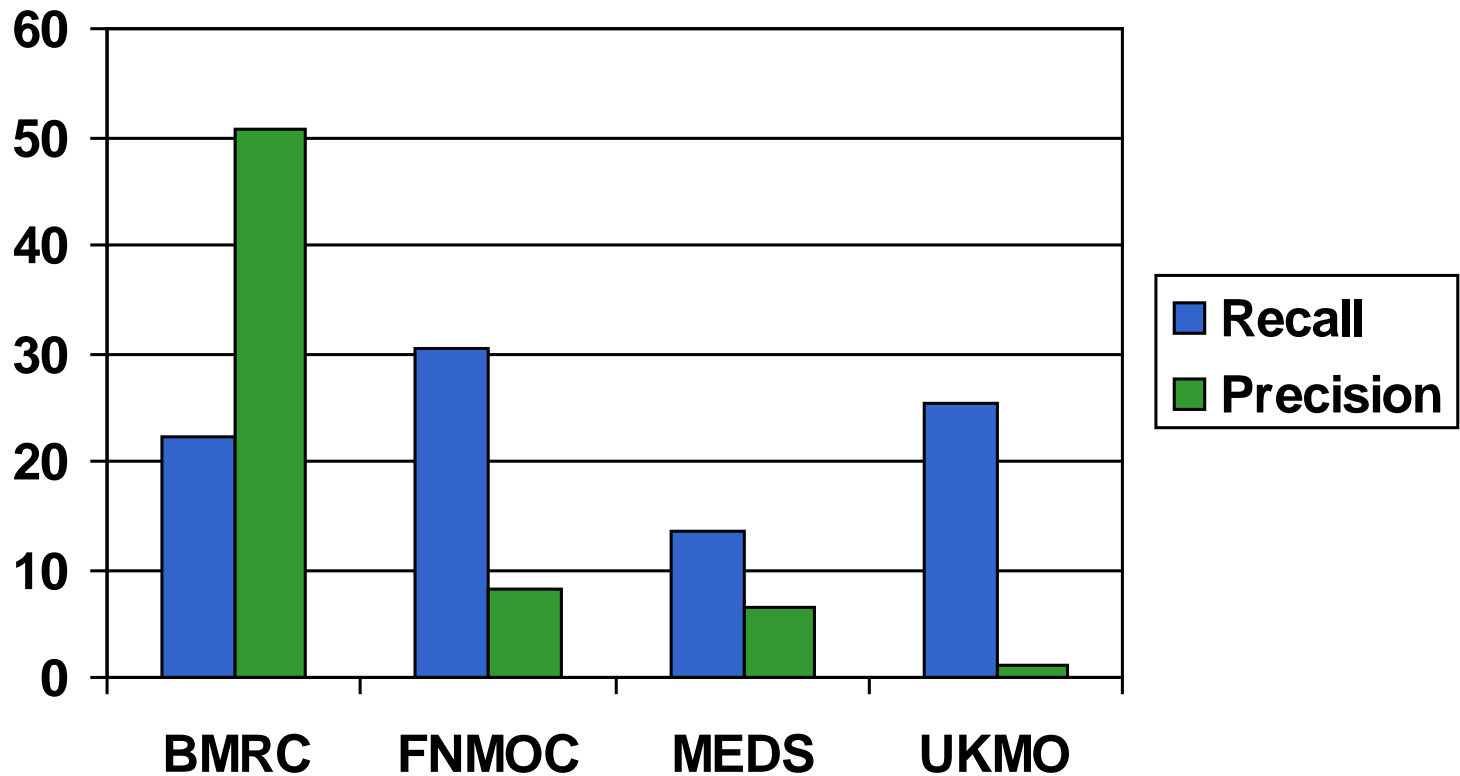
# Recall and Precision

Total: 142229

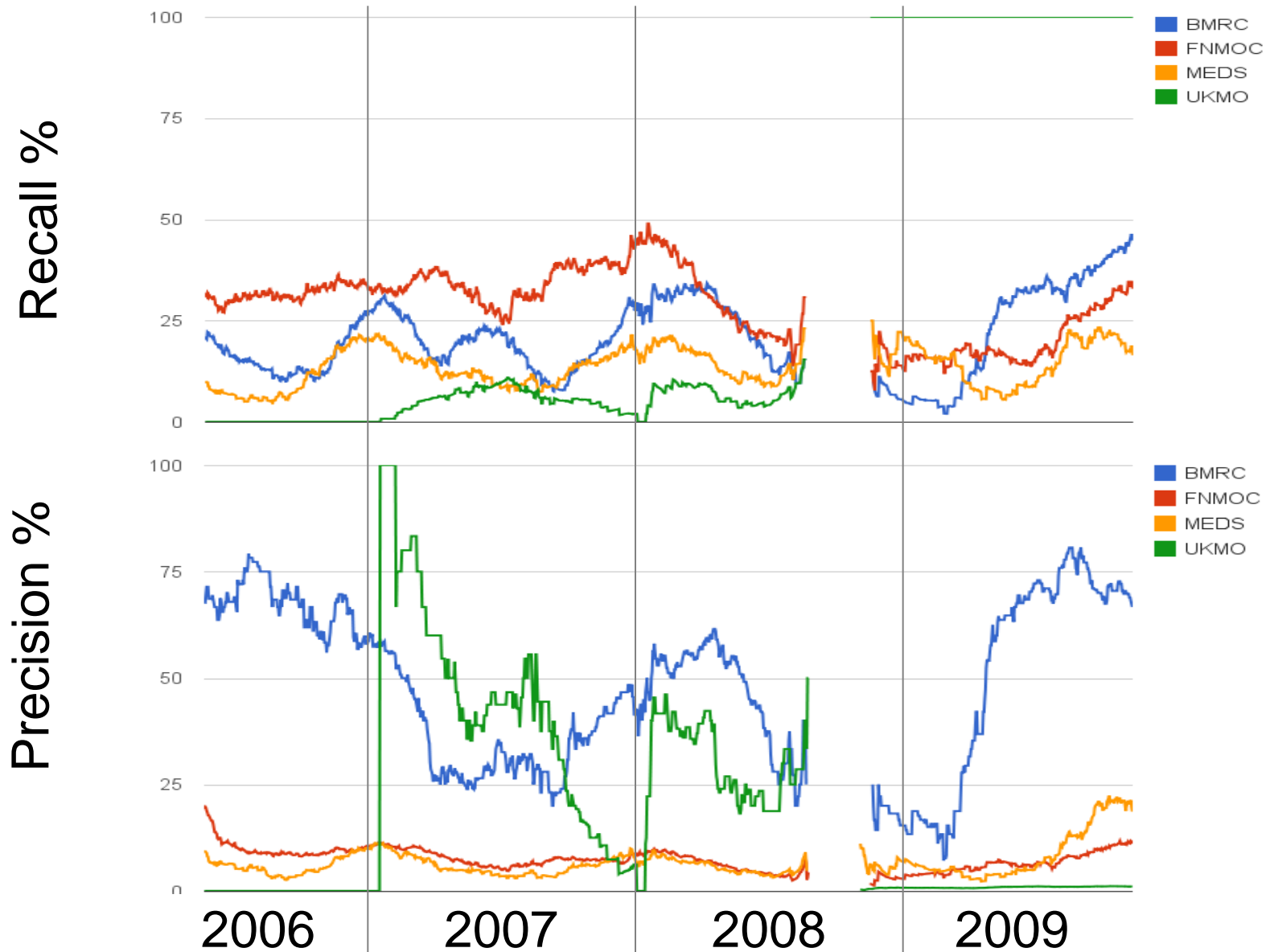


- Recall = what proportion of the **bad data** did they detect?  
→  $269 / 1204 = 22.3\%$
- Precision = what proportion of their **rejects** were actually bad?  
→  $269 / 530 = 50.8\%$

# Recall and Precision by Centre

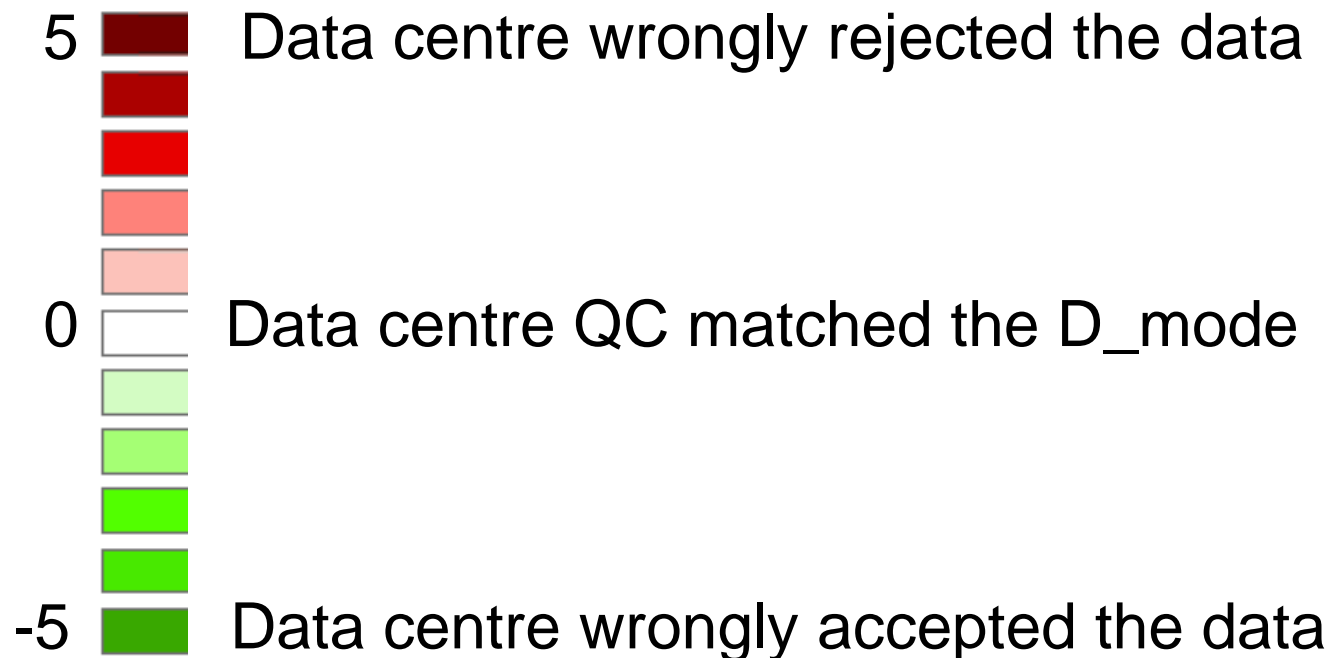


# Recall and Precision over time



# Google Earth – Salinity quality visualisation

- Salinity Quality code (QC) data was used
- Q.C. is ranked 0 to 5
- D\_Mode QC was subtracted from data centre QC to produce a ranking system



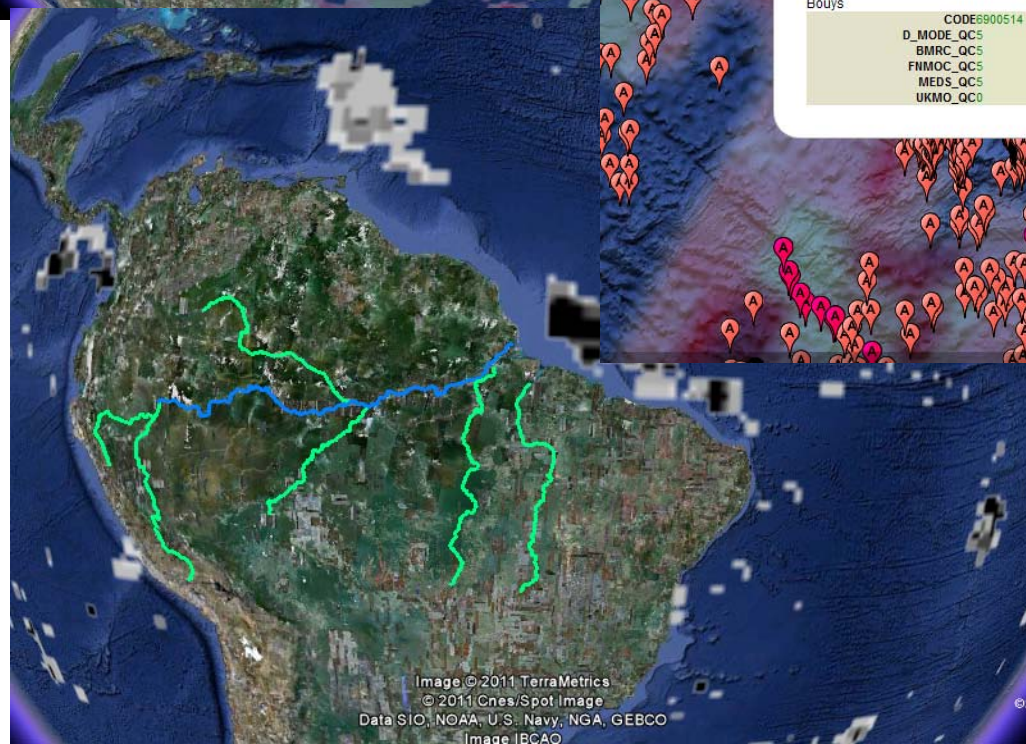
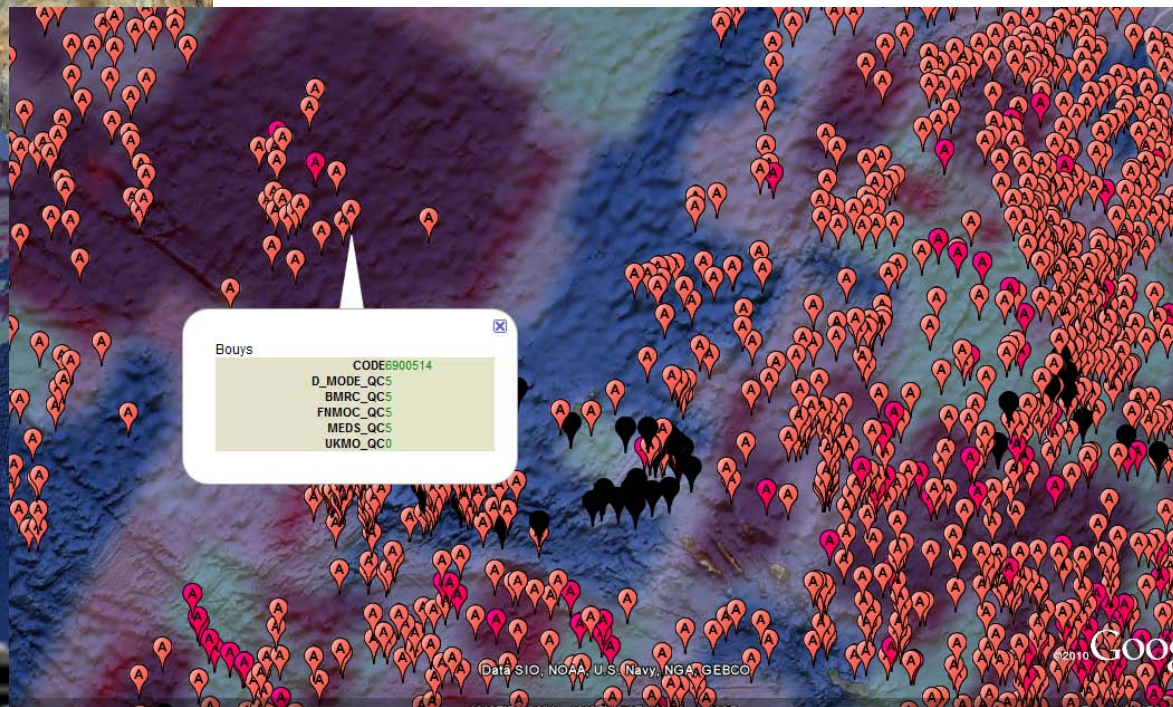
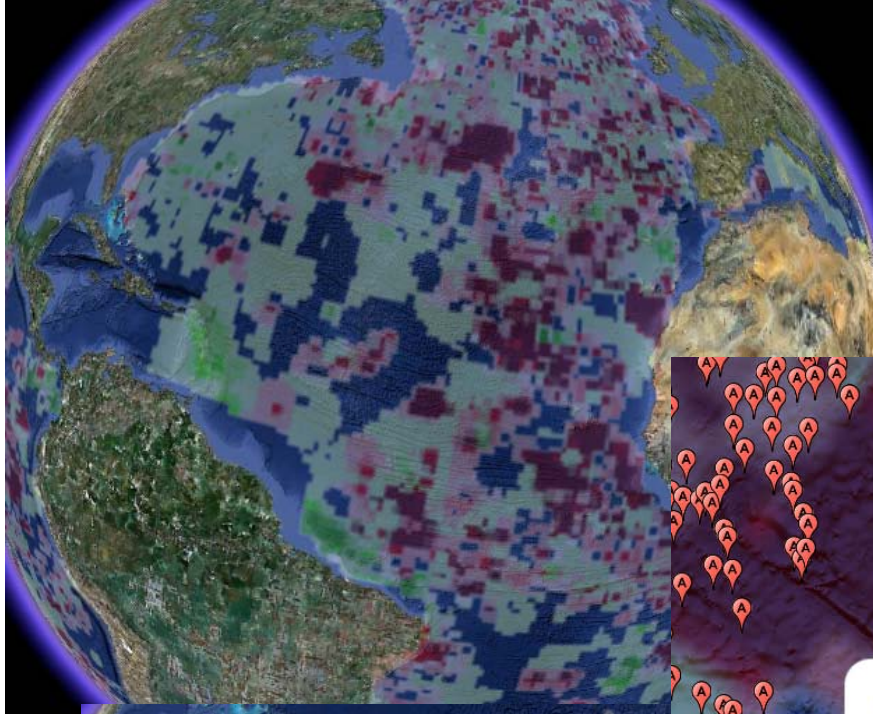


# Google Earth – Salinity quality visualisation

- Raster maps were produced in ArcGIS and then exported into Google Earth. This file can easily be downloaded from the web:

[http://dl.dropbox.com/u/29469014/Centres\\_sal.kmz](http://dl.dropbox.com/u/29469014/Centres_sal.kmz)

- **Inverse distance weighting (IDW)** is a method for [multivariate interpolation](#)



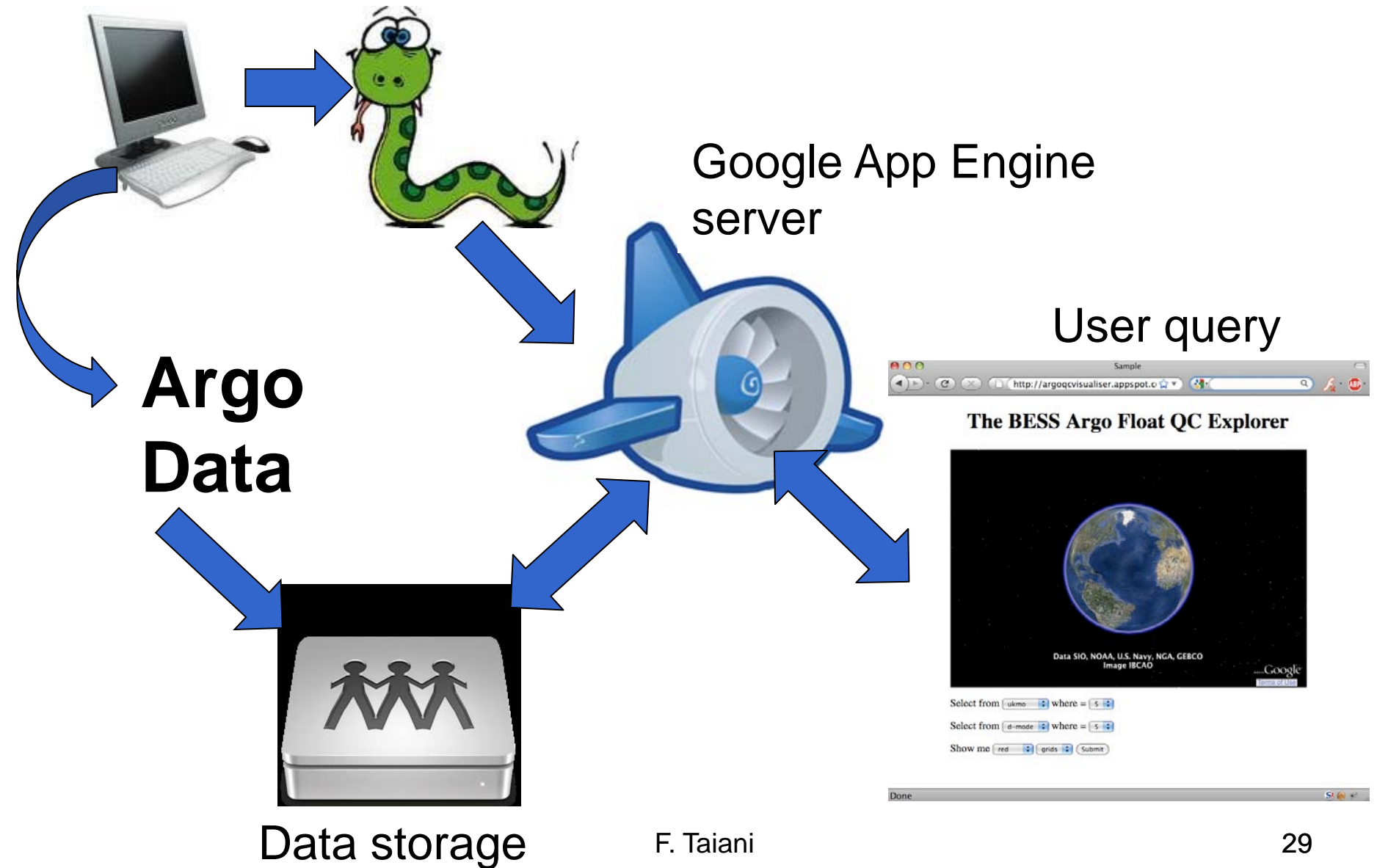
# Algorithm Input / Output

- Two data sets - two values e.g:  
select where d\_mode\_qc = good AND bmrc\_qc = bad
- Other inputs:
  - Colour, presentation, temperature or salinity, comparison operator
- Output: KML file showing
  - Location of the selected floats on the globe
  - Percentage of floats that met the condition in an area

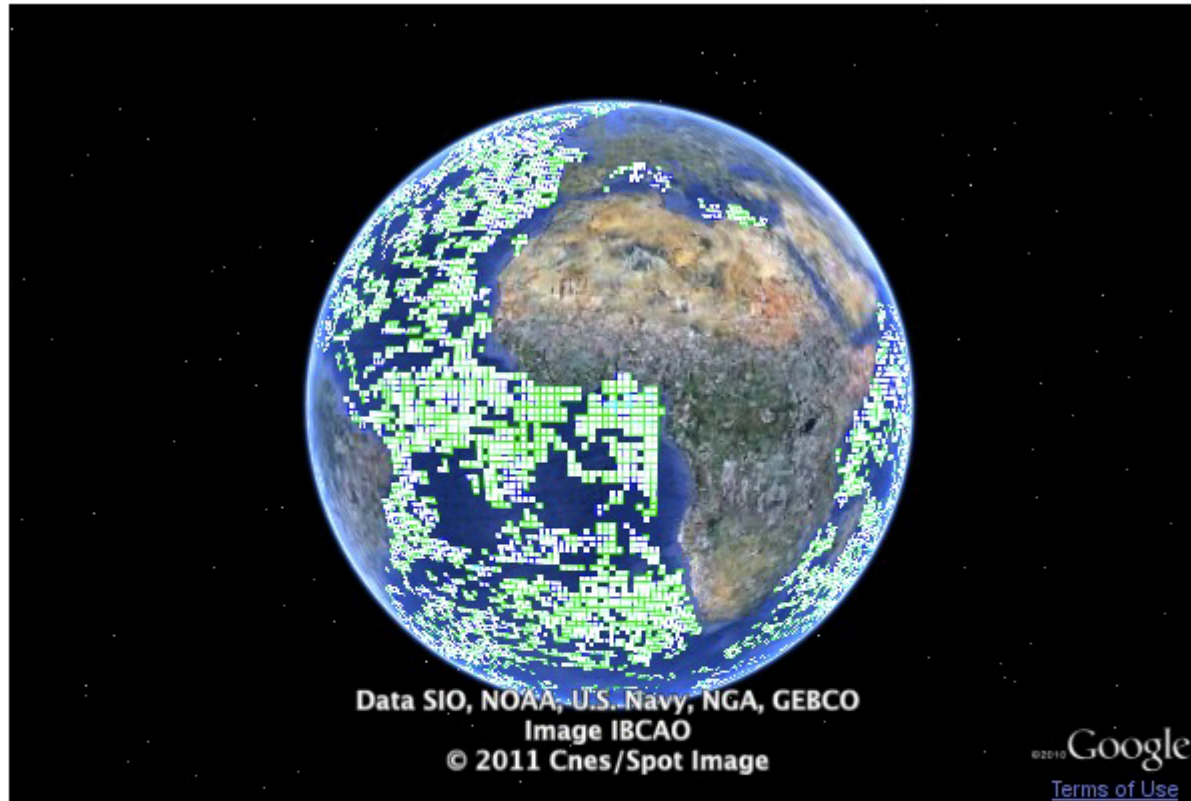
# Internal Work

- Data kept in text file
- Open the file and:
  - Scan the file line by line
  - Evaluate the condition
  - IF pins:
    - Create KML pin
  - IF grid:
    - Bin the result
    - Compute statistic
    - Create Grid
  - Output KML file

# Online visualisation



# Query options



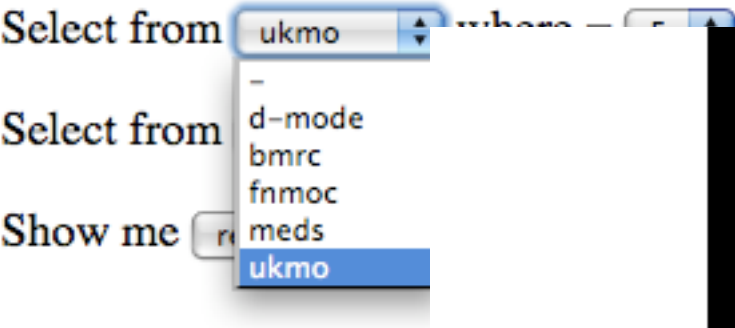
Please select the year:

2006  2007  2008  2009

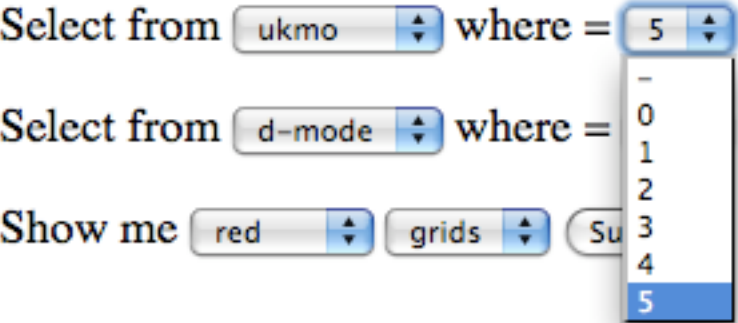
# User Queries



1. Select data centre

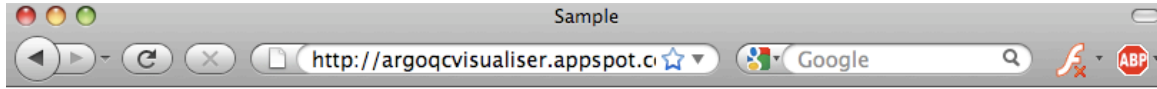


2. Select data qc value



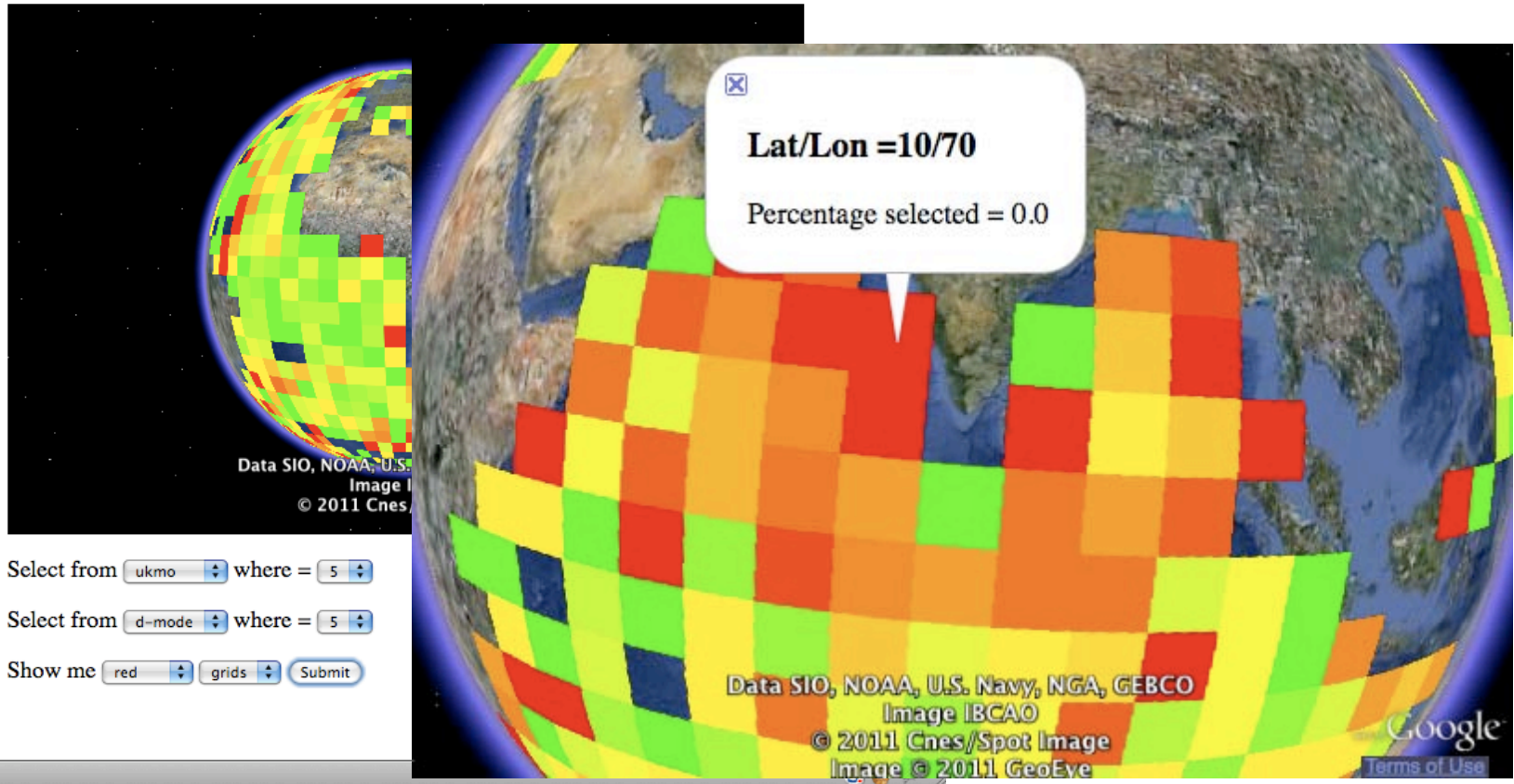
3. Submit!

# And the result...



## The BESS Argo Float QC Explorer

ZOOM...



<http://argoqcvisualiser.appspot.com/>



# Summary 1

- All centres have relatively low recall rates for bad data
- Precision rates are much lower for MEDS and FNMOC suggesting fundamental differences in qc methodology compared to BMRC and UKMO
- Most accepted data is of good quality
- There is a high consistency between temperature and salinity qc values
- No particular time of year that causes more good or bad data

# Summary 2

- 'Spatial-temporal clustering' may be a good tool for evaluating the consistency of argo floats in terms of data quality
- Visualisation tools such as Google Earth have the potential to provide a powerful way to visualise large spatial datasets
- Online application of visualisation tools can provide the ability for rapid and simple queries for a large number of potential users