# Using open source tools to reduce large data sets in a distributed environment
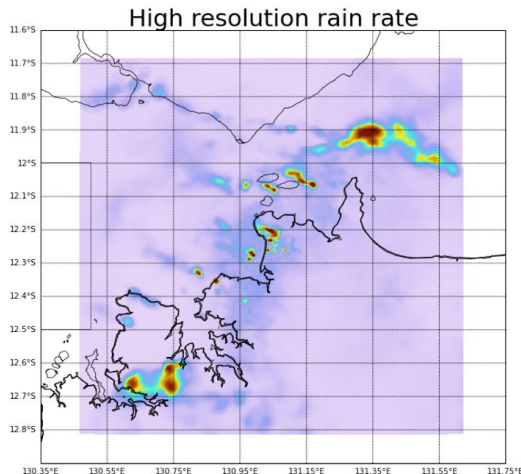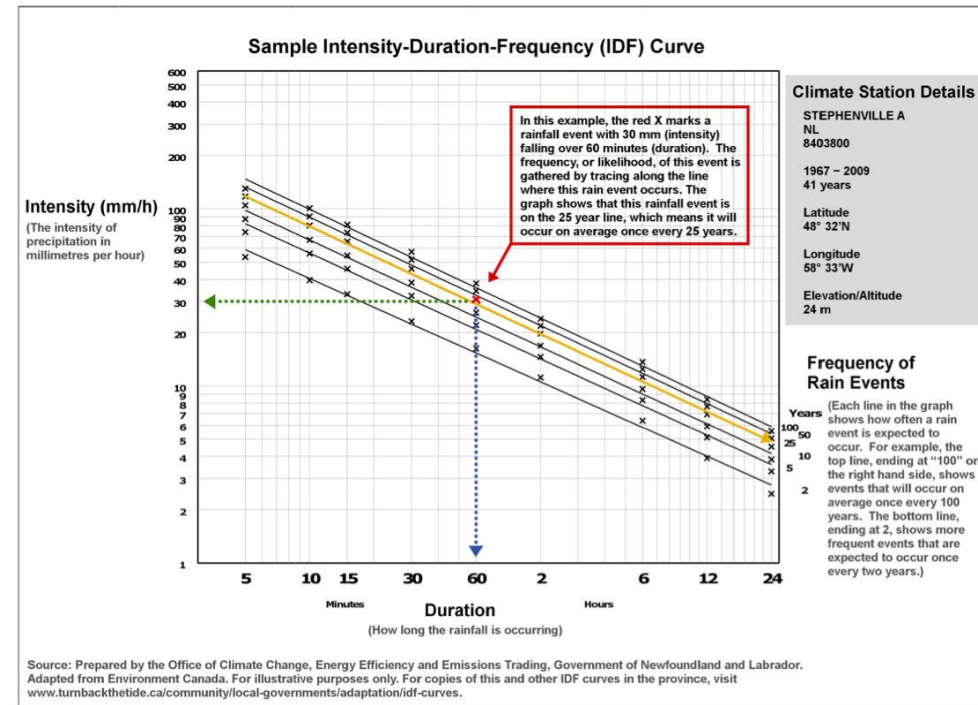
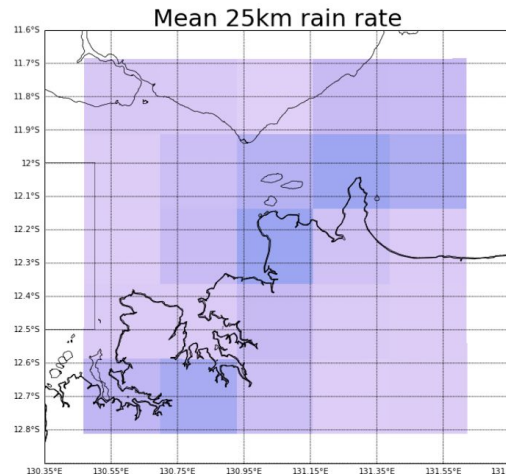*A view from a reluctant and very poor data guy*



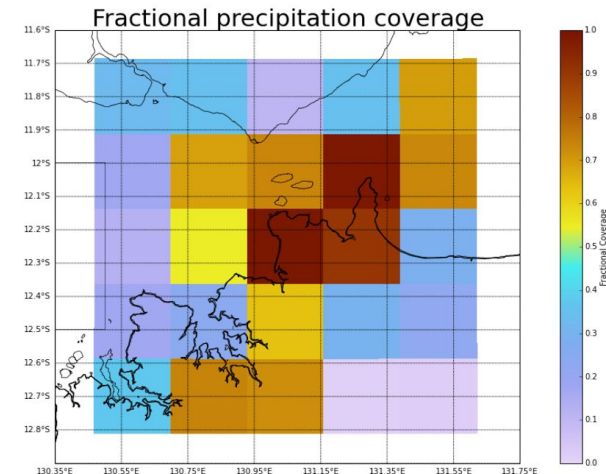*Scott Collis and Jonathan Helmus and many others!*

# The Problem:

- Currently many centers are downscaling GCM runs using WRF to 12.5km to look at local rainfall.
- This does not match to high density hydrological systems (cities). Or natural scales of precipitation.
- Most urban centers have a WSR-88D radar nearby capable of resolutions less than 4km.
- We are investigating the impact of resolution on modelled streamflow.



Sample Intensity-Duration-Frequency (IDF) Curve

In this example, the red X marks a rainfall event with 30 mm (intensity) falling over 60 minutes (duration). The frequency, or likelihood, of this event is gathered by tracing along the line where this rain event occurs. The graph shows that this rainfall event is on the 25 year line, which means it will occur on average once every 25 years.

**Climate Station Details**

STEPHENVILLE A
NL
8403800

1967 – 2009
41 years

Latitude
48° 32'N

Longitude
58° 33'W

Elevation/Altitude
24 m

**Frequency of Rain Events**

(Each line in the graph shows how often a rain event is expected to occur. For example, the top line, ending at "100" on the right hand side, shows events that will occur on average once every 100 years. The bottom line, ending at 2, shows more frequent events that are expected to occur once every two years.)

Intensity (mm/h) (The intensity of precipitation in millimetres per hour)

Duration (How long the rainfall is occurring)

Source: Prepared by the Office of Climate Change, Energy Efficiency and Emissions Trading, Government of Newfoundland and Labrador. Adapted from Environment Canada. For illustrative purposes only. For copies of this and other IDF curves in the province, visit www.turnbackthetide.ca/community/local-governments/adaptation/idf-curves.



High resolution rain rate



Mean 25km rain rate



Fractional precipitation coverage

**(a)**                **(b)**                **(c)**

# Solution: Radar. New Problem: Lots of data.

- A radar system images precipitating systems.
- It measures reflectivity factor, or the 6th moment of DSD as well as a number of other measurements depending on the sophistication of the system.
- So many formats. Complex geometry.
- Produces many thousands of files per year.
- Fortunately the problem is pleasantly parallel. The data is highly granular.
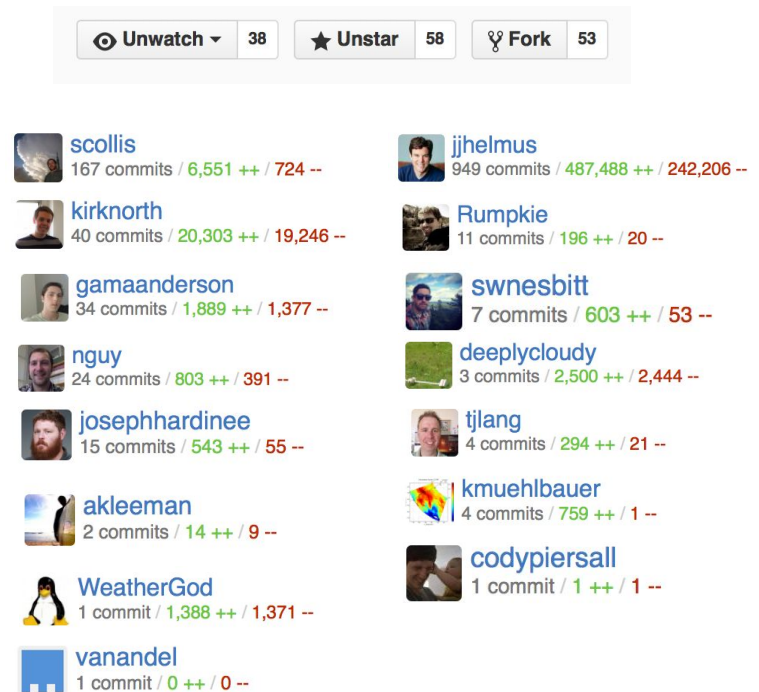


By Original Photograph By: Andrew J Oldaker - http://en.wikipedia.org/wiki/Image:LabNexrad.jpg, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=3124873

## Open source software + Cluster tools make the problem tractable

# Py-ART the Python ARM Radar Toolkit

## https://github.com/ARM-DOE/pyart

- Py-ART uses a moderately complex data model which closely mirrors the CF-Radial community file format.
- It makes heavy use of Python dictionaries to form a self –describing radar object. Rich IO layer allowing a very large number file formats to be read into the data model. Primary output format is CF-Radial. We are investigating adding ODIM as an output format.
- Community codebase on GitHub, main fork is DoE maintained and moderated. Set of core dependences (all in Anaconda Python) with many optional dependencies which, when present, increase functionality.
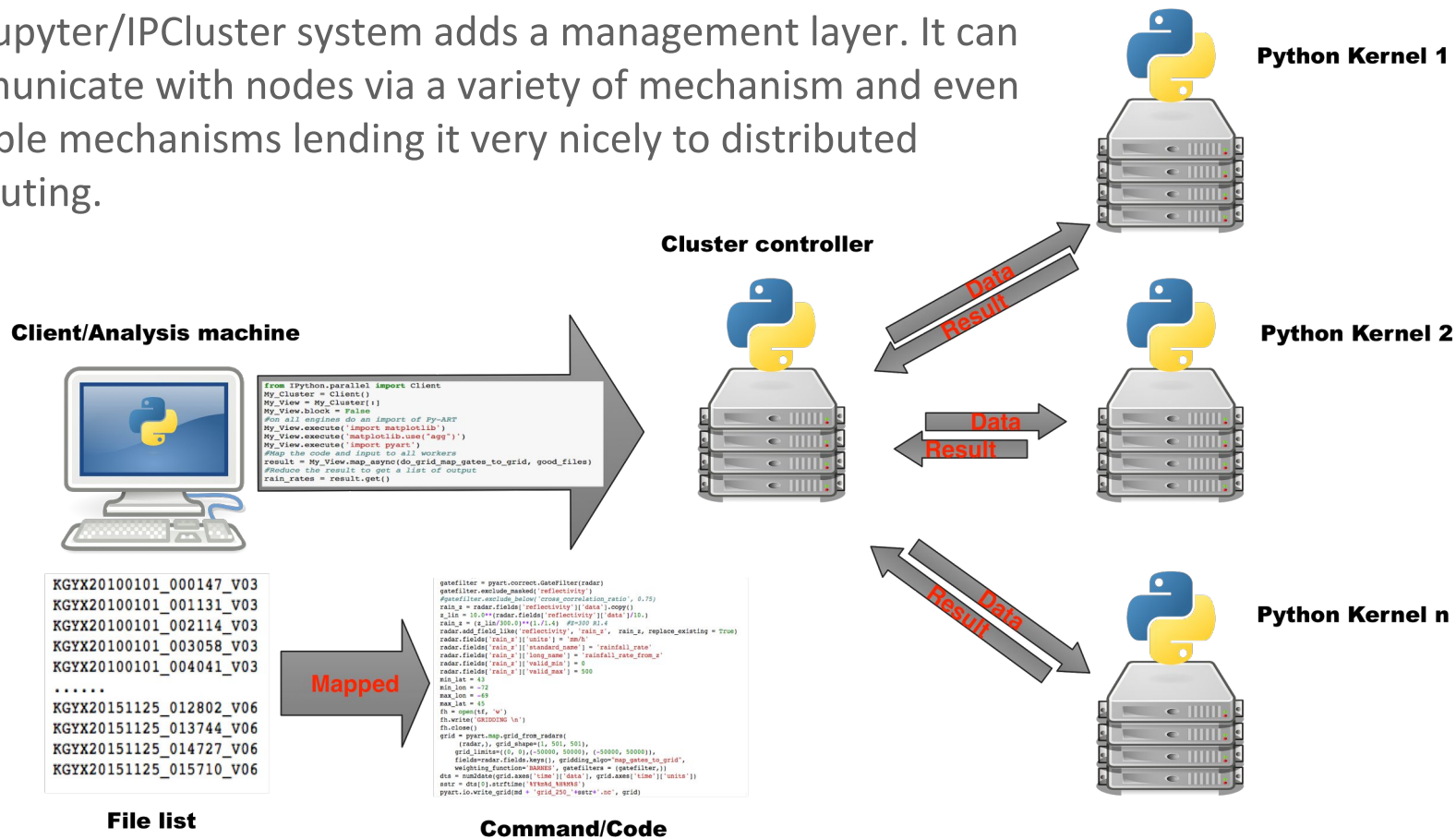


⊙ Unwatch ▾ | 38    ★ Unstar | 58    ⑂ Fork | 53

**scollis**
167 commits / 6,551 ++ / 724 --

**jjhelmus**
949 commits / 487,488 ++ / 242,206 --

**kirknorth**
40 commits / 20,303 ++ / 19,246 --

**Rumpkie**
11 commits / 196 ++ / 20 --

**gamaanderson**
34 commits / 1,889 ++ / 1,377 --

**swnesbitt**
7 commits / 603 ++ / 53 --

**nguy**
24 commits / 803 ++ / 391 --

**deeplycloudy**
3 commits / 2,500 ++ / 2,444 --

**josephhardinee**
15 commits / 543 ++ / 55 --

**tjlang**
4 commits / 294 ++ / 21 --

**akleeman**
2 commits / 14 ++ / 9 --

**kmuehlbauer**
4 commits / 759 ++ / 1 --

**WeatherGod**
1 commit / 1,388 ++ / 1,371 --

**codypiersall**
1 commit / 1 ++ / 1 --

**vanandel**
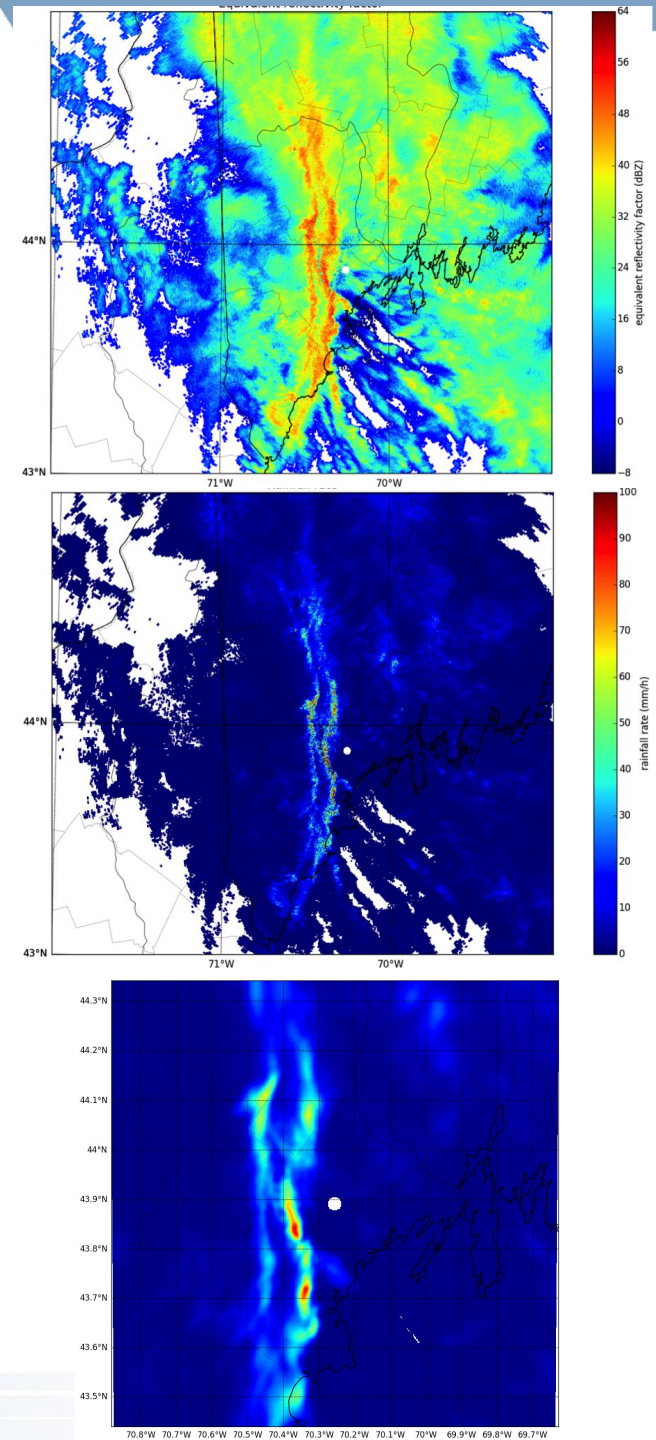1 commit / 0 ++ / 0 --

**ARM**
CLIMATE RESEARCH FACILITY

# IPython/Jupyter Cluster. Easy map reduce

- In the quest for reproducible, collaborative science, on thing the Jupyter project has done is to work with "detached kernels".
- This is how notebooks work.. But you can have an array of kernels waiting for 0MQ packets to do work and return results.
- The Jupyter/IPCluster system adds a management layer. It can communicate with nodes via a variety of mechanism and even multiple mechanisms lending it very nicely to distributed computing.



Python Kernel 1

Python Kernel 2

Python Kernel n

Cluster controller

Client/Analysis machine

```
from IPython.parallel import Client
My_Cluster = Client()
My_View = My_Cluster[:]
My_View.block = False
#on all engines do an import of Py-ART
My_View.execute('import matplotlib')
My_View.execute('matplotlib.use("agg")')
My_View.execute('import pyart')
#Map the code and input to all workers
result = My_View.map_async(do_grid_map_gates_to_grid, good_files)
#Reduce the result to get a list of output
rain_rates = result.get()
```

KGYX20100101_000147_V03
KGYX20100101_001131_V03
KGYX20100101_002114_V03
KGYX20100101_003058_V03
KGYX20100101_004041_V03
......
KGYX20151125_012802_V06
KGYX20151125_013744_V06
KGYX20151125_014727_V06
KGYX20151125_015710_V06

**Mapped**

Data / Result

File list

Command/Code

# The study: Rainfall around Portland, Maine.

- Portland is the target for a Department of Homeland Security Regional Resiliency Assessment Program study.
- Step 1 : Download A LOT of data.. This is the **most challenging step.** We downloaded 10 years worth of data from NCEI and staged it on an optimized storage setup.
- Step 2: Develop your technique for one file.
- Step 3: Google like mad to find out how to set up a Jupyter cluster on a PBS managed system. Fortunately this step stays done.. and you have a lasting resource!
- Step 4: Map the file list from step 1 to the technique in step 2.

# Demo 1: Very simple rainfall grid with Py-ART

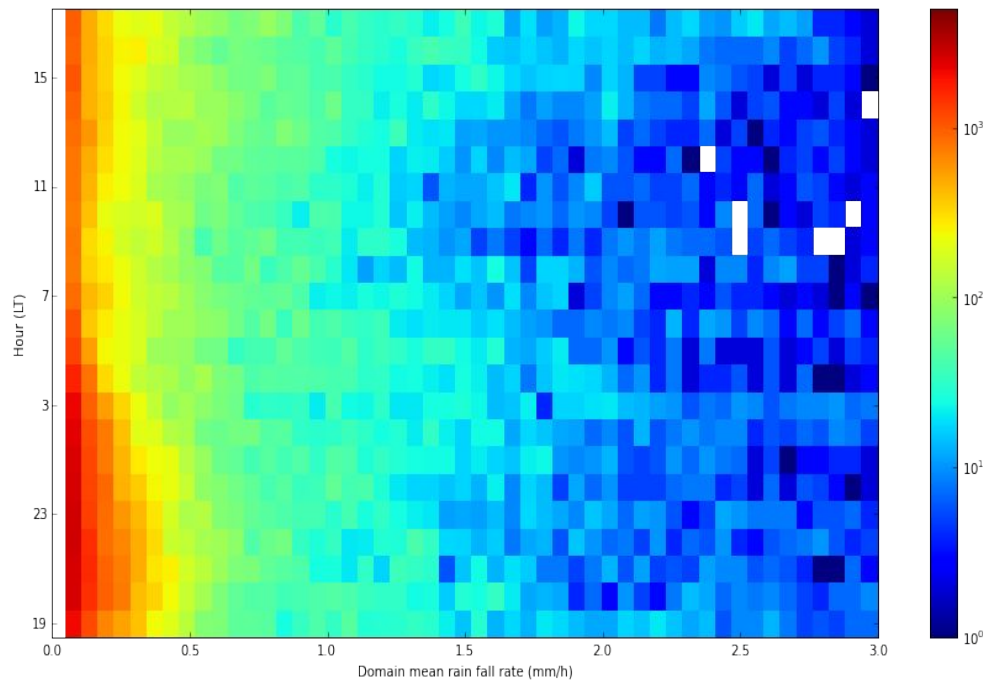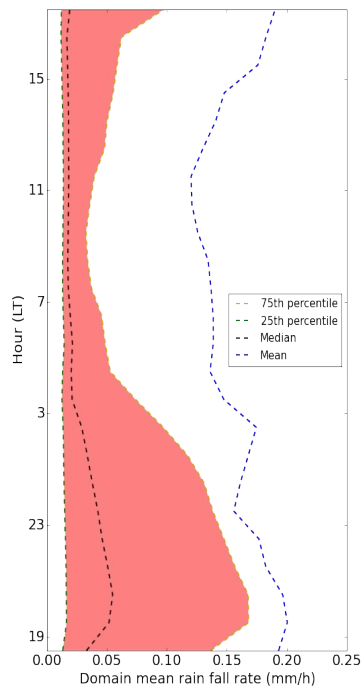# Demo 2: Simple demo of the Jupyter Cluster.

# Some Results

- 700,000 files, 10+TB was reduced to grids and domain mean rainfall calculated.
- We saved these off as a 100MB times series for easy analysis and data indexing.
- We then, because we can, stratified by hour of the day and monthly.

# Some concluding thoughts

- Ability to push single "granule" code to a distributed environment is key.
- This allows very quick revisit times for development and makes decadal+ processing possible.
- Good, thread safe, well tested open source software,**even if it is written without the explicit desire to be parallel** is key.
- As is a good system to map the problem to many cores.
- This was, well.. very messy. Good data tools are a key need.. would love it if NEXRAD data came with some descriptive inventory data.
- Py-ART has grown in strides. We have benefited greatly from unfunded contributions. Key in encouraging this is using funds to grow out accepted formats base. This is a key communique to programs: open mindedness pays dividends..
- We are now working with Amazon to see how far we can push this.



https://github.com/scollis/high_resolution_hydrology