

Cloud Archiving and Data Mining: Operational and Research Examples

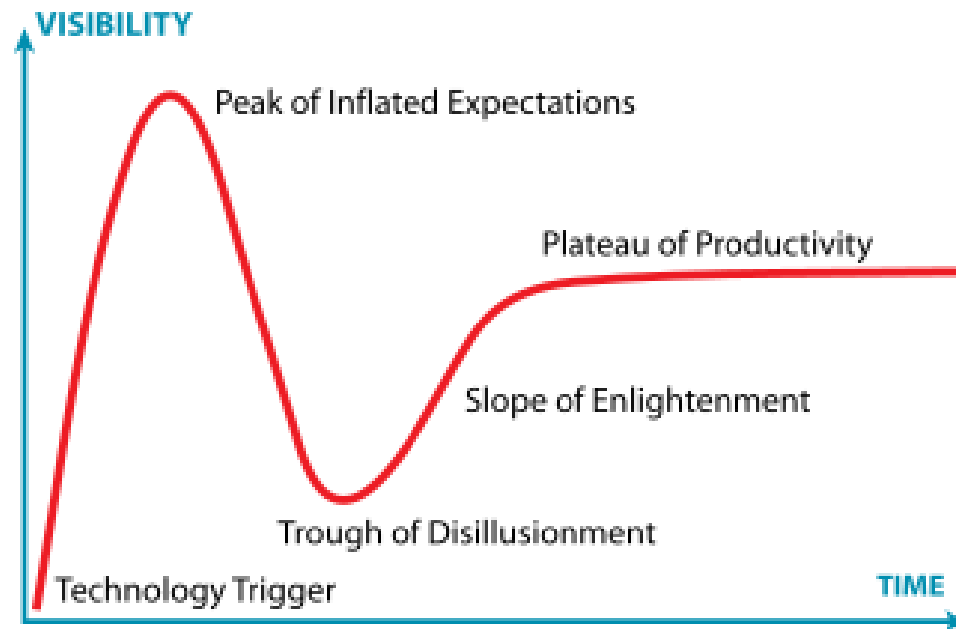
John Horel*, Brian Blaylock, Chris Galli*

Department of Atmospheric Sciences

University of Utah

john.horel@utah.edu

*Synoptic Data Corporation



*Amara's "law":
Overestimating the
effects of a
technology in the
short run and
underestimating the
effects in the long run*

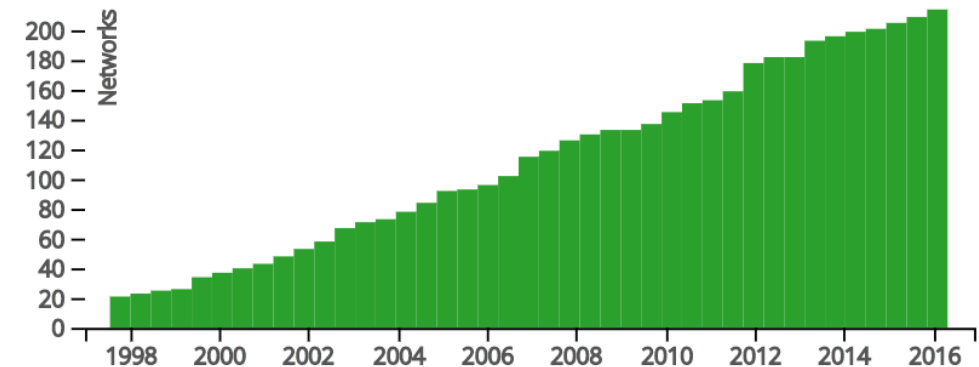
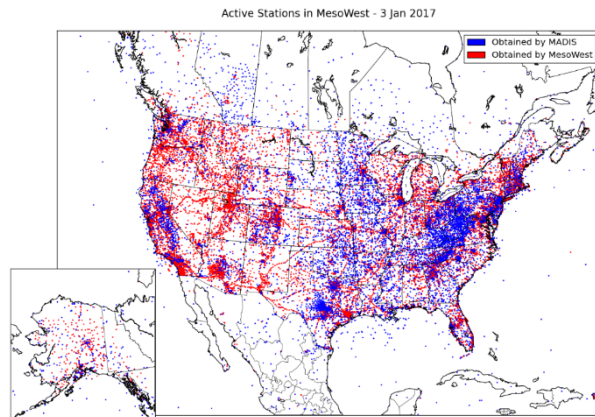
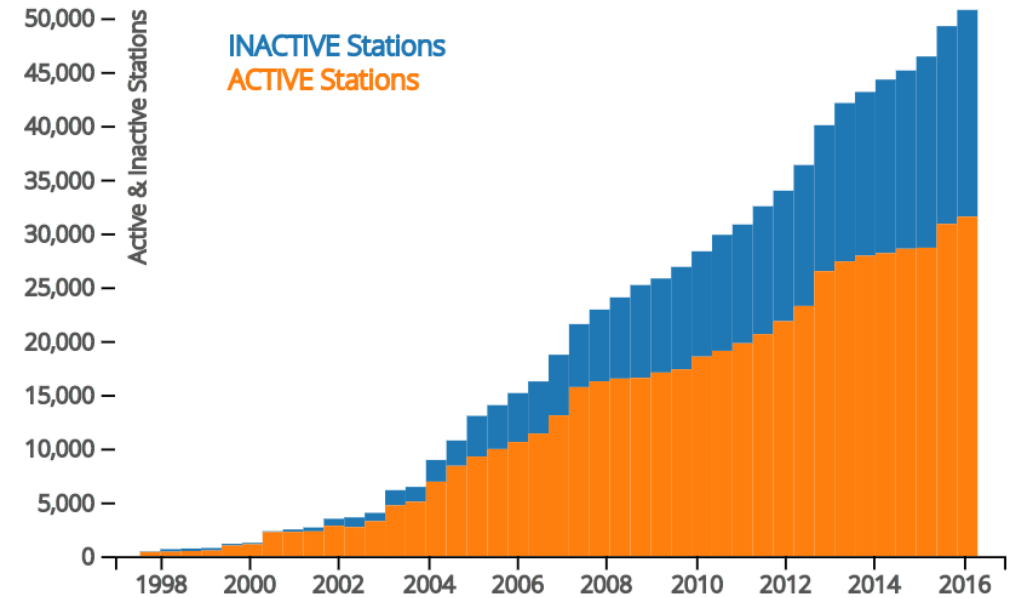
MesoWest/SynopticLabs

<https://synopticlabs.org/>

- Objective: access, archive, and disseminate publicly-accessible provisional environmental data
- Transitioned over last several years from University of Utah to cloud-based IT infrastructure



~40 billion observations





Public Cloud Infrastructure

- Cloud archiving:
 - highly efficient mySQL/TokuDB databases
- Data mining via API
 - over a million requests to download over 10 billion data values per day

Number & Types of Server Instances

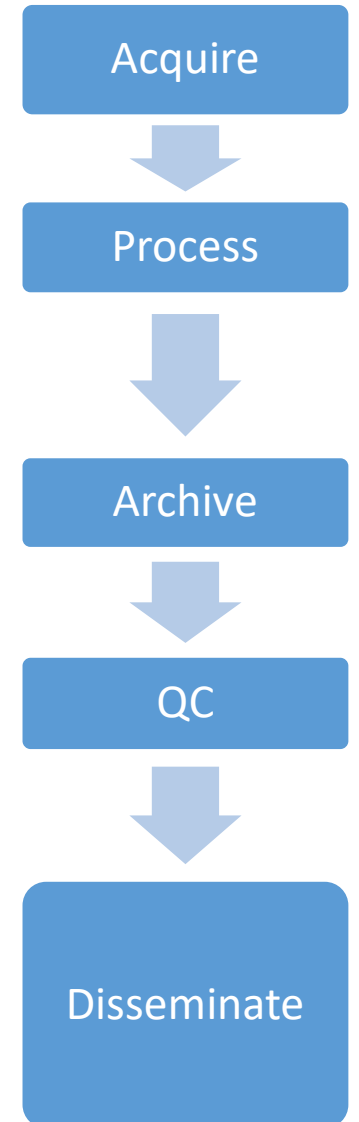
13 Ingest, Processing, Load Balancing

6 Database

2 Real-time data checking

8 Web services/API/alerts

Type	Quantity	Monthly Cost
Server Instances	29	\$1800
Disk	3.5 TB	\$350
Data Transfers	Egress	\$350
Total		\$2500



Issues with Cloud Archiving/Data Mining

- Critical to manage costs and only use resources that are absolutely necessary
- Balancing data storage vs. data access costs is challenging
- Cheaper solutions exist to store large amounts of data, but may result in large latency to retrieve that data
- Compute nodes required for data mining need to be close to data archive to avoid data transfer costs

Cloud Archiving and Data Mining of Forecast Model Output: High Resolution Rapid Refresh (HRRR)

No retrospective archive of HRRR model output at this time at NCEI, NCEP, or ESRL accessible externally

What we needed for WRF model initialization and HRRR model validation:

- Efficient and expandable archival storage for thousands of large GRIB2 files
- Fast retrieval of 2D fields within those files
- Ability to make data publicly accessible to other researchers

Solution: Object storage is an affordable, useable, and reliable long-term archive approach

See poster by Brian Blaylock & access archive via: <http://hrrr.chpc.utah.edu/>

Manuscript submitted: Blaylock, B., J. Horel, and S. Liston, 2017: Cloud Archiving and Data Mining of High Resolution Rapid Refresh Forecast Model Output. *Computers and Geosciences*. In review.

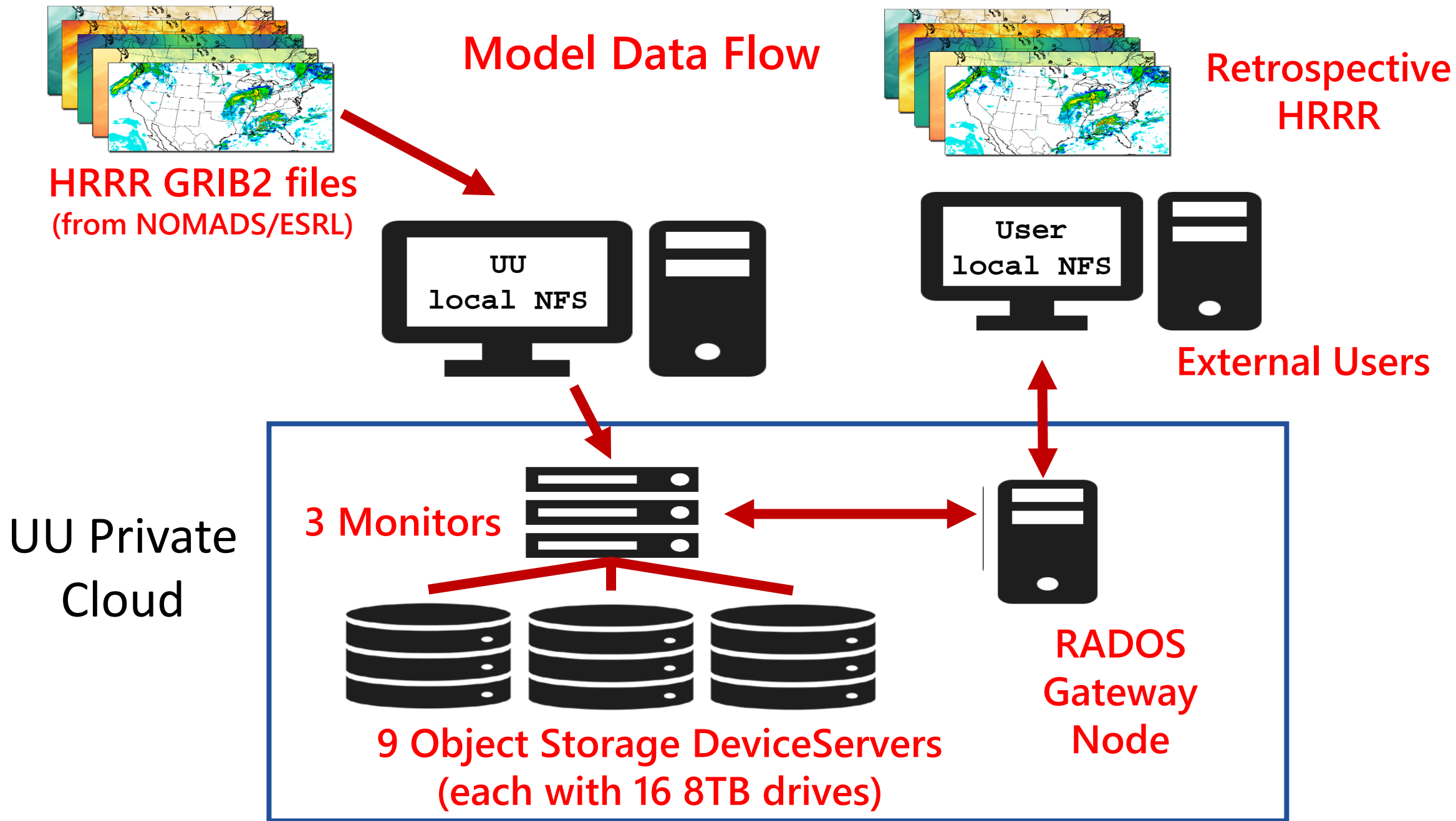
Object Store Options

- Public cloud (*Amazon Web Services, Microsoft Azure, etc.*)
- Private cloud (*University of Utah Data Center*)

Our choice: Private Cloud @ University of Utah

- Disk-based S3-like object storage at the University of Utah's Data Center managed by the Center for High Performance Computing (CHPC)
- Configured with 6+3 erasure coding: Objects are broken into 9 pieces—6 data pieces and 3 redundancy pieces
- No data loss, even if every disk fails on three servers
- New hardware and expanded storage can easily be added or replaced over time

Cloud Storage Service	Cost over 5 Years
CHPC Pando	\$120/TB
Amazon Glacier	\$540/TB + \$30/TB download
Microsoft Azure	\$600/TB + \$10/TB download

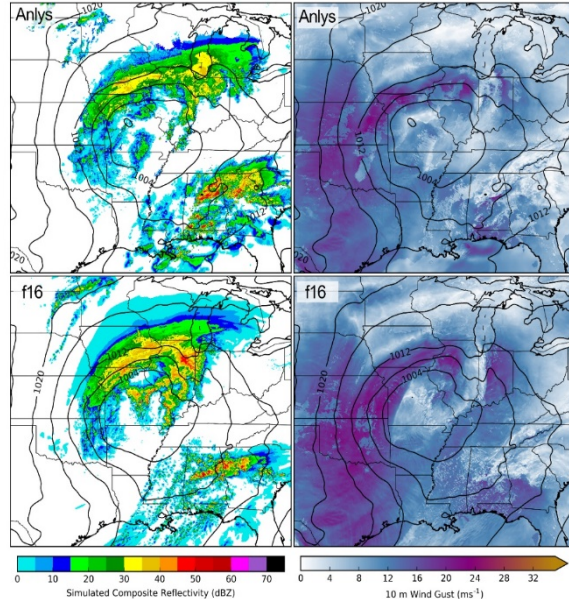


HRRR Output Being Archived (~70 GB/day)

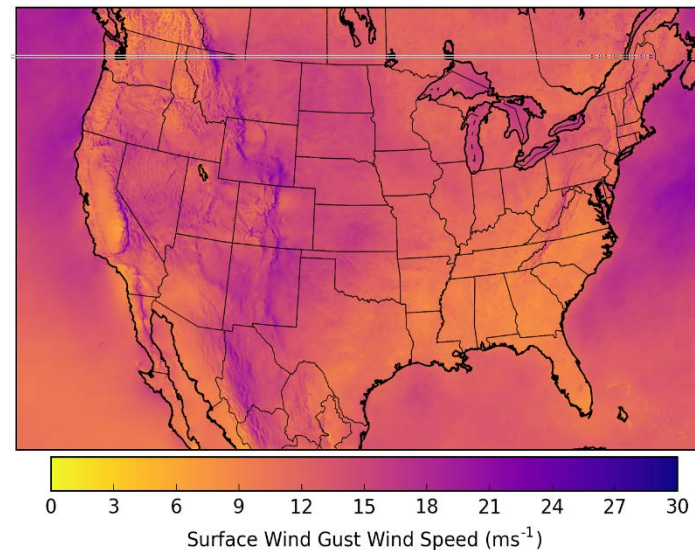
Model Type	File Type		First Date Available
Operational HRRR	2D fields	f00-f18	Analyses: April 18, 2015 Forecasts: July 27, 2016 Subhour: May 11 2017
	3D fields	f00	
	<i>Subhourly fields</i>	<i>f00-f18</i>	
Experimental HRRR	2D fields	f00	December 1, 2016
Experimental HRRR Alaska	2D fields	f00-f36	September 1, 2016
	3D fields	f00	

Data mining of HRRR output

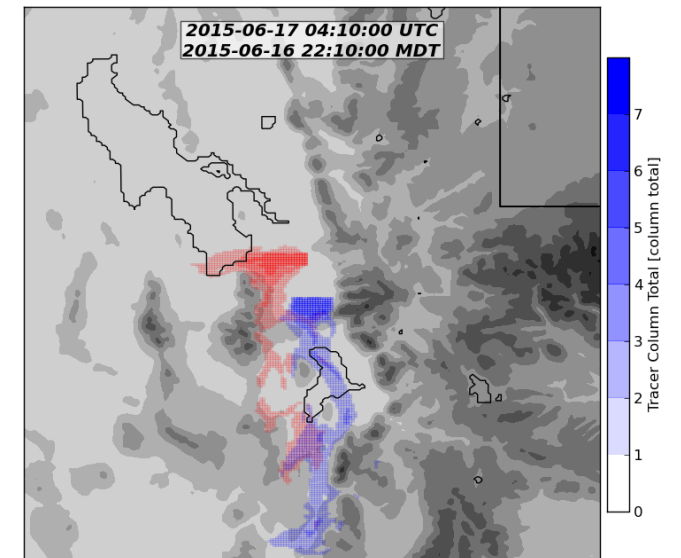
- Users may download the entire field with `wget` or `cURL` or download a single 2D surface if the byte range is known
- Possible to use multi-thread python to simultaneously return thousands of grids in small amount of time



HRRR model validation



Wind Gust 95th Percentile over 2 Years

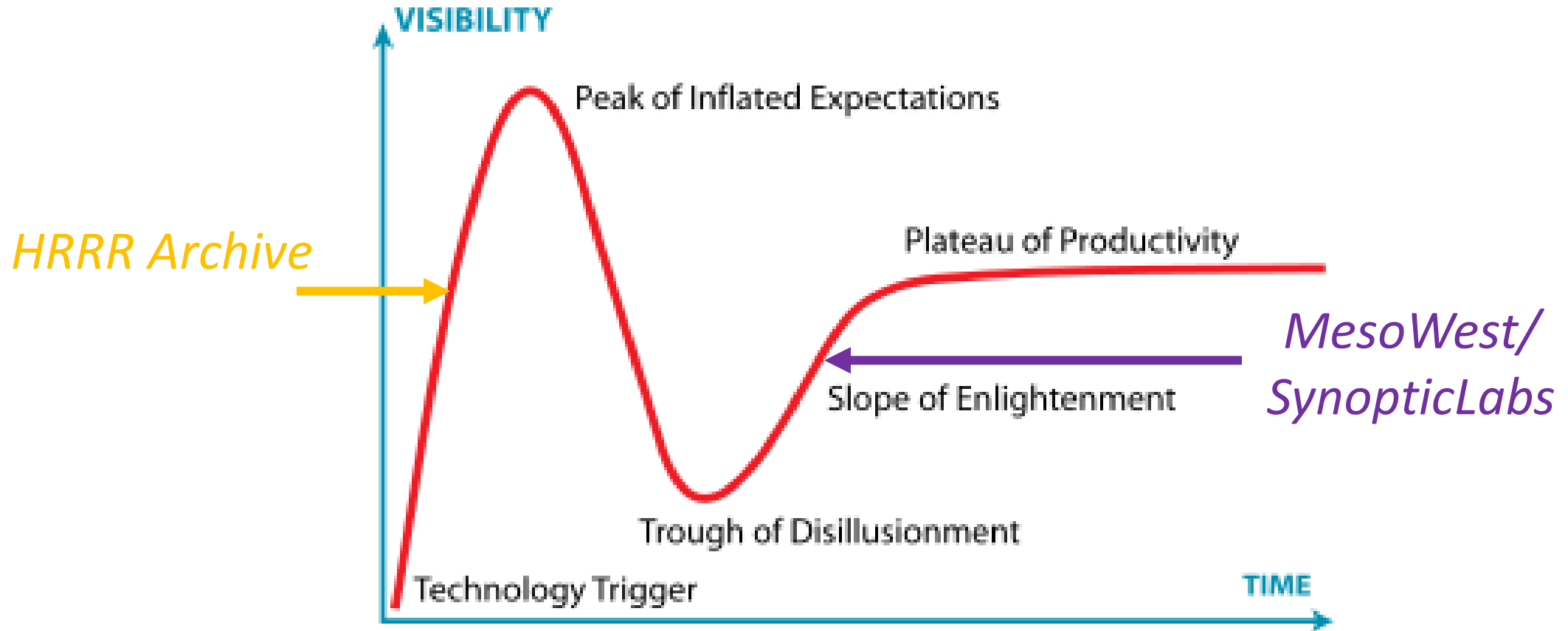


WRF Initialization

Issues with Cloud Archiving/Data Mining

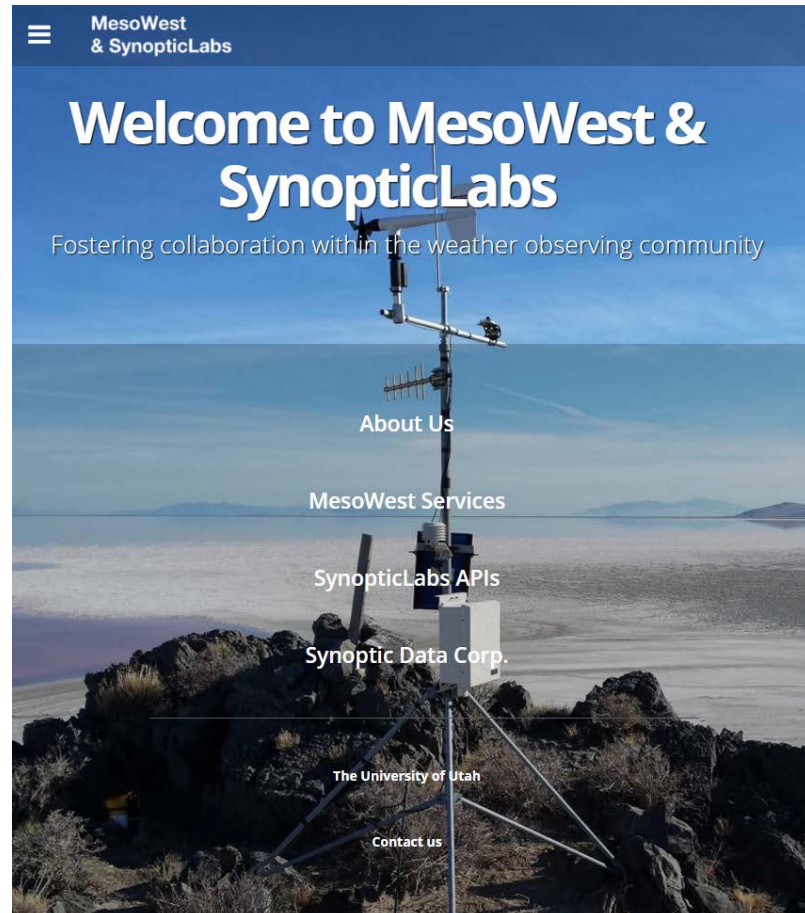
- S3-type objects must be downloaded to a local disk to process the data
 - Incur data download costs from public cloud provider
 - Red Hat now supports *POSIX* compliant Ceph File System to handle objects in cloud
- To avoid excessive data downloading, GRIB2 format allows selecting by byte range and returning only the fields within an object
 - Other file formats such as HDF5 may eventually allow subsetting of S3 objects by variable, region, single grid point, all vertical levels at a point, etc.
 - Siphon, NetCDF subsetting service
- NSF & other funding agencies require data management plans
 - While geoscience data repositories exist, they have strict standards
 - Academic institutions need to meet data stewardship requirements
 - Will institutions subsidize the costs to maintain large archives? *“UU Hive” has 500 GB limit*

Cloud Archiving and Data Mining: Operational and Research Examples



Questions?

<https://synopticlabs.org/>



<http://hrrr.chpc.utah.edu/>

HRRR Archive at the University of Utah

Frequently Asked Questions

If you are looking for the most current HRRR output (last two days) you should download directly from the [NOMADS](#) server.

[Have you Registered?](#)

[Best Practices](#)

[HRRR FAQ](#)

[Web Download Page](#)

[About this HRRR archive](#)

[Who archives HRRR?](#)

[Where is the archive?](#)

[What days are available?](#)

[Who uses the HRRR archive?](#)

[Gallery](#)

[Tips for CHPC users](#)

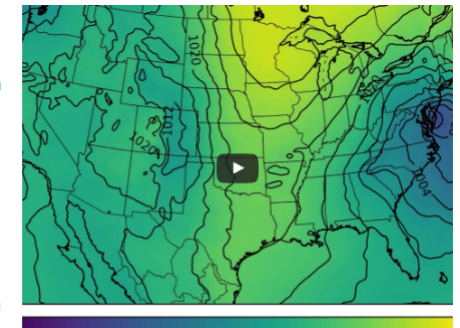
About this HRRR archive

What is the HRRR archive?

The HRRR archive is a collection of output from NCEP's High Resolution Rapid Refresh model. This is a model developed by NOAA ESRL and is run operationally every hour at NCEP. It continues to be developed by ESRL.

The HRRR generates hourly forecasts gridded at 3 km for 18 hours over the continental United States making it the highest spatial and temporal resolution forecast system run by NCEP.

HRRR analyses and forecasts are exceptionally valuable to the research community. However, an official HRRR archive does not exist. We began archiving HRRR data in April 2015 to support research efforts at the University of Utah. Instead of downloading all available files, we only download files most useful to accomplish our research efforts. We have made the archive publicly accessible for research purposes.



MesoWest/SynopticLabs

<https://synopticlabs.org/>

- Objective: access, archive, and disseminate publicly-accessible provisional environmental data **rapidly**
- Drivers:

