I'm Grace Peng, one of the data specialists at the NCAR Research Data Archive aka as the RDA.

Many within NCAR know us as the people who stage data on /glade for you. Some people in this room may have never heard of us. For more than 40 years, the RDA has been collecting and disseminating weather and climate data to the research community. We host a growing collection of over 600 datasets.

Data specialists are part data curators, part data engineers, part software developers, and part subject matter experts. Among our data specialists, we have experts in meteorology, oceanography, mathematics, physics, chemistry and engineering.

We're witnessing a large increase in the use of wx and climate data in the commercial, government and educational sector. Our growing user base includes more non-experts so we find ourselves becoming educators about wx and climate data as well.

Data remains alive as long as people are using it. Working with students helps us learn what we need to do to today's data to make it useful for future users.

By the people come to me with their questions, they have usually begun their

# Questions

- Can this problem be answered with data?
- What type of data can help me answer this question?
- Where can I find the best data for the job?
- Does the data match the data documentation?
- How do I cite the data for reproducibility? DOI

NCAR UCAR | Data Thinking    *air • planet • people*  2

I'll try to get through these five fundamental questions in 25 minutes.  If I don't, please come to the Better than Free workshop tomorrow.

The first Q, Can this… sounds obvious, but it's actually difficult.  This is the part where you make sure you have a science project and not an Op-Ed piece.

The second (and first) Q together helps you formulate your research plan.

Data discovery, finding the best data for your project, is a difficult technical problem. Data search engines founded on metadata can help you.

If the data doesn't match what you expect (from the documentation), you could be experiencing technical difficulties reading the data in correctly, or you may have discovered an error in the data.

Science should be reproducible.  You've heard other talks this morning about data curation and digital object identifiers for data sets.  I'll only need to cover the mechanics of how to cite data so that others understand what you did.

Meta: Can this problem be answered with data? Should it?

- Is it measurable?
- Or related to something measurable?
- Is it technically possible?
- Ethics
  - Privacy
  - Attract attention from wrong places

NCAR UCAR | Data Thinking                    *air • planet • people* 3

Can this problem be answered with data? is a meta question that you need to think through before starting on any endeavor.

Is it measurable?
    If not directly measureable, is it related to something measurable?
    Has the data been collected already? May have to wait until data exists.

Should this problem be tackled?
    Don't collect sensitive information if you can't safeguard it.

We think of weather as benign data. But, as weather apps turn people into mobile weather stations, what are the ethical and safety considerations? If we know where someone lives (their night-time closest cell tower), and also where they work or go to school during the day, we can narrow it down to a few people. If that information ends up in the hands of advertisers, it is annoying and intrusive. If it ends up in the hands of a violent ex or a repressive government, it could be deadly.

In my own work, I help researchers obtain data for pollution studies in countries where questioning the environmental costs of development has landed people in jail. What if their government demands I hand over their data request records?

Data availability informs our approaches to research problems.

What kind of spatio-temporal coverage and resolution do you need?
Precision and accuracy (P&A) of measurements?

We all want the "best" data, with the most P&A. But, work through your error budget to find the minimum acceptable P&A. This will open up more data possibilities for your research.

Does it require data fusion; do you need to combine data from multiple sources? This is usually the case or else it wouldn't be research.

New data enables new inquiries. But, do not overlook the value of looking at old data in new ways. My favorite example is the paper that examined the prevalence of cirrus clouds in the SLC area before and after it became a Delta Airlines hub airport. They correlated jet fuel tax receipts and cirrus prevalence and made a very convincing case that jet contrails are the cause of the increasing prevalence of cirrus clouds near SLC.

## Data Discovery: Does this data exist?

- Prior work
  - Check their data citations
- Data papers
  - Describe how new datasets are prepared
- Data centers
  - Use their search features
    - Metadata: data about data
  - Read their "What's new" announcements

NCAR UCAR | Data Thinking                           *air • planet • people* 5

The practice of Data discovery is still in the juvenile phase. There isn't one killer app or dominant search engine that can search all data repositories yet.

This is a mixed bag with upsides. Consider the time spent in data discovery a learning opportunity.

Read research related to your project. What kinds of data do they use and where did they obtain it?

Data can be cited directly, as an entity in its own right, often with a DOI. This is the modern, preferred method.

Aside:
[DOI stands for Digital Object Identifiers are one to one mappings of data sets (DS) to DOI number (unlike DS names). Some places give new DOIs to new versions of a DS, while others give new DOIs for major new versions of DSs.

However, the practice of citing data papers, papers describing the preparation of a data set, is also in widespread use and has been used for much longer.

Before the advent of data papers, people would cite the first science paper that used

As an example, take a look at the RDA front door, rda.ucar.edu

You can do a free text search for data at the top.
Or, you can perform a faceted search, successively narrowing selections based upon metadata criteria.
These are just a few of the many options enabled at the RDA.

Check out our new datasets and other announcements on our blog.

Take a video tour of our home page to learn how to use the myriad features of our site.

## Data and Metadata Standards

- Global Change Master Directory
  http://gcmd.nasa.gov/learn/keyword_list.html
- CF (Climate and Forecast) NetCDF
  http://cfconventions.org/index.html
- WMO GRIB1, GRIB2, BUFR
- ANSI 19115-X
- Federal Geographic Data Committee (FGDC) 1998
  http://www.usgs.gov/core_science_systems/
  csas/metadata/standards.html

NCAR
UCAR | Data Thinking                    *air • planet • people*  7

The RDA's search features are powered by standardized metadata, most commonly Global Change Master Directory keywords.

NetCDF is closely associated with CF standards and NCAR.  It is self-describing and self-contained.

There are several other standards in wide use.  Some older standards such as FGDC may be superceded by newer ones.

However, some old standards such as WMO and GRIB remain in daily use worldwide. GRIB and BUFR are self-describing, but not self-contained.  They require external table to unpack the binary data correctly.  If you apply the wrong table, you read the values incorrectly.   GRIB and BUFR are widely used in operational transmission, but must be careful if used later.  If you are unsure about the correct GRIB table to apply, ASK!

A data file can be mapped to more than one data standard for interoperability.  In fact, older but still useful data sets should be mapped to new metadata standards to maintain their discoverability by both humans and machines.

GCMD keywords map out explicitly the contents of data files. This should eliminate guesswork, which can lead to misinterpreting data.

This looks excessively detailed, but this is how machine-to-machine communication works.

This also enables distributed search, where one portal can crawl other data sites for you.

Data producers can obtain help from data curators/metadata specialists to map out their data.
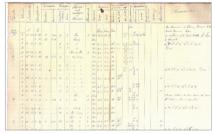
Figure 3. Page of the experimental universal logbook of the *Prince of Orange*, covering the period 15-19 May 1853. The logbook, designed by Maury and Jansen, was tested in practice during this trip of this ship. This happened in the preparation of the Brussels Conference in November 1853.
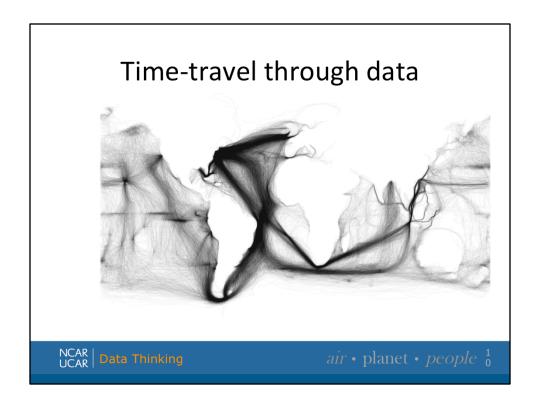
Standardization of data is not a new idea.  Once, whaling ship logbooks looked like the one at left.  In 1853,  Maury and Jansen created and tested an universal logbook format that is similar to what's used today.  Then they organized an international conference to promote the use of their universal logbook format.

Seafarers could still maintain idiosyncratic logbooks like those at left, but they were required to also log in their information in the universal logbook with standardized units.

Bauman Rare Books http://sappingattention.blogspot.com/2012/11/reading-digital-sources-case-study-in.html

The Brussels Conference and its legacy ftp://ftp.wmo.int/Documents/PublicWeb/amp/mmop/documents/JCOMM-TR/J-TR-27-BRU150-Proceedings/DOCUMENTS_JCOMM_27/Session_4/4_2_Konnen.pdf

Time-travel through data

NCAR UCAR | Data Thinking          air • planet • people

This data standardization enables us to create maps like this.

We know global shipping routes and traffic over time from these logbooks.  We agreed to common units of measurement (latitude and longitude, Greenwich time, Celsius scale, frequency of recorded measurements) http://sappingattention.blogspot.com/2012/11/reading-digital-sources-case-study-in.html

The data of the past is available to us today because of preservation efforts of data compilers and curators of yesterday.  It's understandable to us only because they handed down not just the data, but the knowledge about the data in the form of metadata and adherence to well-described data conventions.  These are known as data and metadata standards.  The corollary is that our adherence to data standards today will enable your data of today to time travel and aid future scientists.

This is currently the RDA's most popular dataset; it is used to initialize Numerical Weather Prediction (NWP) models such as WRF.

If we archive a data set, we also archive documentation about how it was prepared and software to read it.
If the data set is archived in one of the major standards, such as NetCDF or GRIB, we don't need to say much more about it.

But, if it is in a more idiosyncratic format, we will provide specialized readers and other help.

Note that each data set is assigned to a data specialist whose contact information is at the top right. Data specialists serve as a conduit between data users and data producers.

- Access data with tools
- Standardized data
- Reusable tools
- Reusable data
- Happy users
- Data infrastructure
- Data that remains in use remains alive

NCAR | UCAR    Data Thinking            air • planet • people ½

Users access data through tools.  One-off (ad-hoc) data requires one-off tools.

Data and metadata standardization allows us to create reusable tools.

This ecosystem of data standards, tools, search and delivery tools makes up our shared data infrastructure.  When infrastructure works smoothly, the magic becomes mundane and invisible.

Users are not charged the full cost of developing or maintaining infrastructure.  But please keep in mind that infrastructure is not cheap and delaying maintenance can cause catastrophic failures in the future.

# Veracity

- Is this data what you expect?
- Is it as described in the documentation?
  - Format
  - Field contents
- Are the values physical?
- If any of these answers are no, consult with the data provider and/or data specialist. rdahelp@ucar.edu

Speaking of magic.  When you open up those data files, take some time to interrogate its contents.  If you have time, explore beyond the values that you originally wanted to use to see if something else might open up new possibilities for your research.

If you don't get what you expect to see, either in the file contents or values that don't make physical sense, stop.  Consult with the data provider.  If you got the data from us, send a message to rdahelp@ucar.edu and cc it to the data specialist for that data set.

# Reproducible Science: Data Citation

- Provider/creator
- Dataset name
- DOI if it has one

- Source
- Revision, access date
- RDA citation widget

**How to Cite This Dataset:**
RIS
BibTeX

National Centers for Environmental Prediction/National Weather Service/NOAA/U.S. Department of Commerce. 2000, updated daily. *NCEP FNL Operational Model Global Tropospheric Analyses, continuing from July 1999*. Research Data Archive at the National Center for Atmospheric Research, Computational and Information Systems Laboratory. http://dx.doi.org/10.5065/D6M043C6. Accessed† dd mmm yyyy.
†Please fill in the "Accessed" date with the day, month, and year (e.g. - 5 Aug 2011) you last accessed the data from the RDA.
Bibliographic citation shown in [ Federation of Earth Science Information Partners (ESIP) ⇕ ] style
Get a customized data citation

NCAR | UCAR | Data Thinking          *air • planet • people* ¼

---

For the sake of reproducible science, use sound data citation practices.   Different journals have different data citation styles. They share common required elements.

Who created the dataset?

How is it called?.  Give the DOI if it has one.

Where did you get it? Data can be replicated in multiple places that serve different subsets of the data.  Citing the data source is important.

Data is sometimes reprocessed or corrected.  Put down the version number and data access date.

If you make a data request through the RDA web interface, the request is automatically logged in our database.  Should anything happen, we can regenerate the data request for you or anyone else who wants to repeat your work.  This can save you storage space as you need only save locally what you are currently using.

Our database records also allow us to make custom data citations for you with this widget found on all of our data set home pages.  The widget generates citations in AMS, AGU and several other popular styles.  You can also create RIS and BibTeX

Questions? rdahelp@ucar.edu
http://ncarrda.blogspot.com/p/recommended-data-reading.html

NCAR UCAR | Data Thinking            air · planet · people  15

If you are interested in thinking meta about data, I put up a list of my favorite data books.  We also maintain a blog where we make announcements, give tips and post tutorials.

Thank-you.