

Strategies for supporting students' explorations of big data

Kim Kastens

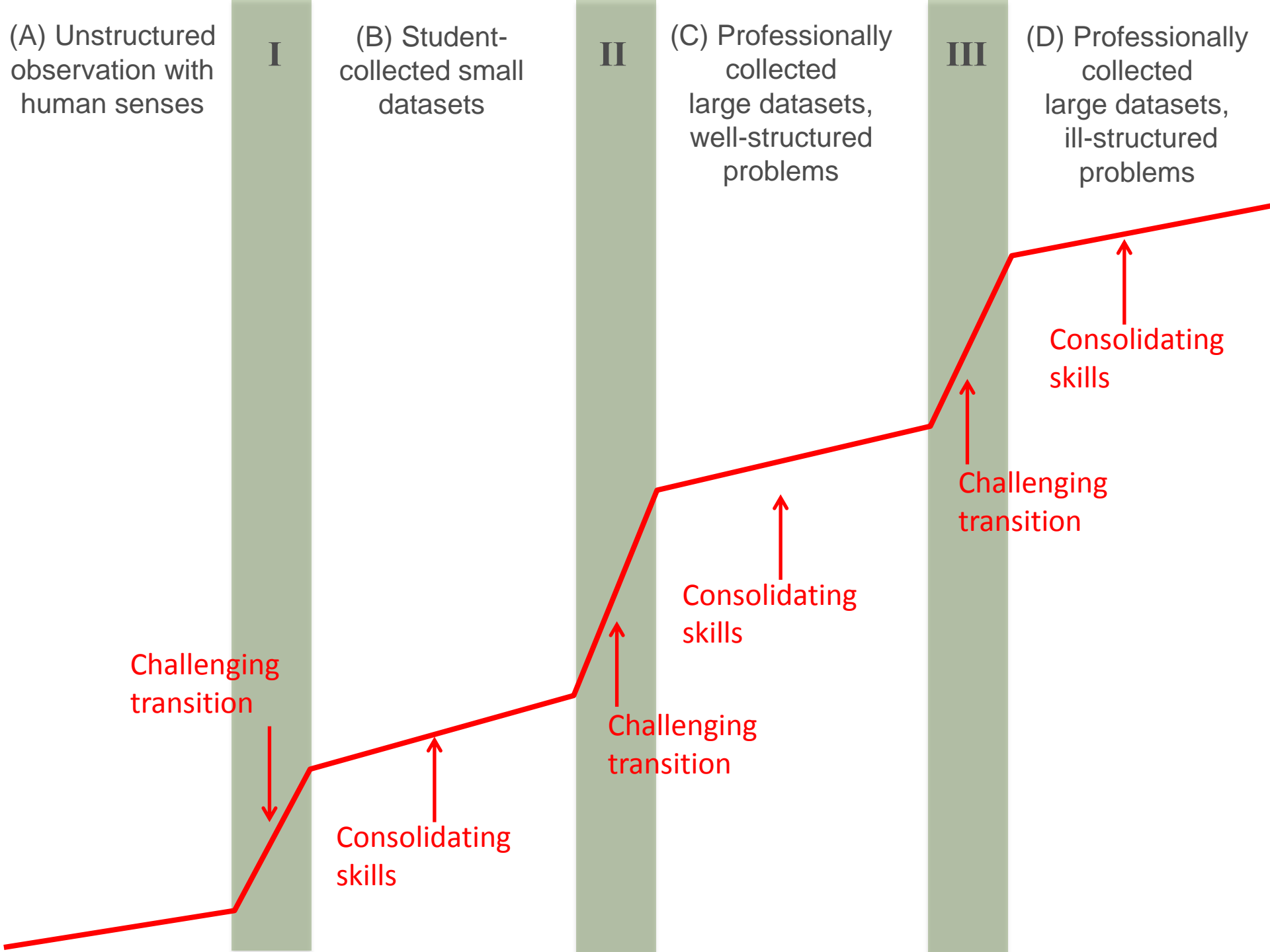
Education Development Center & Lamont-Doherty Earth Observatory

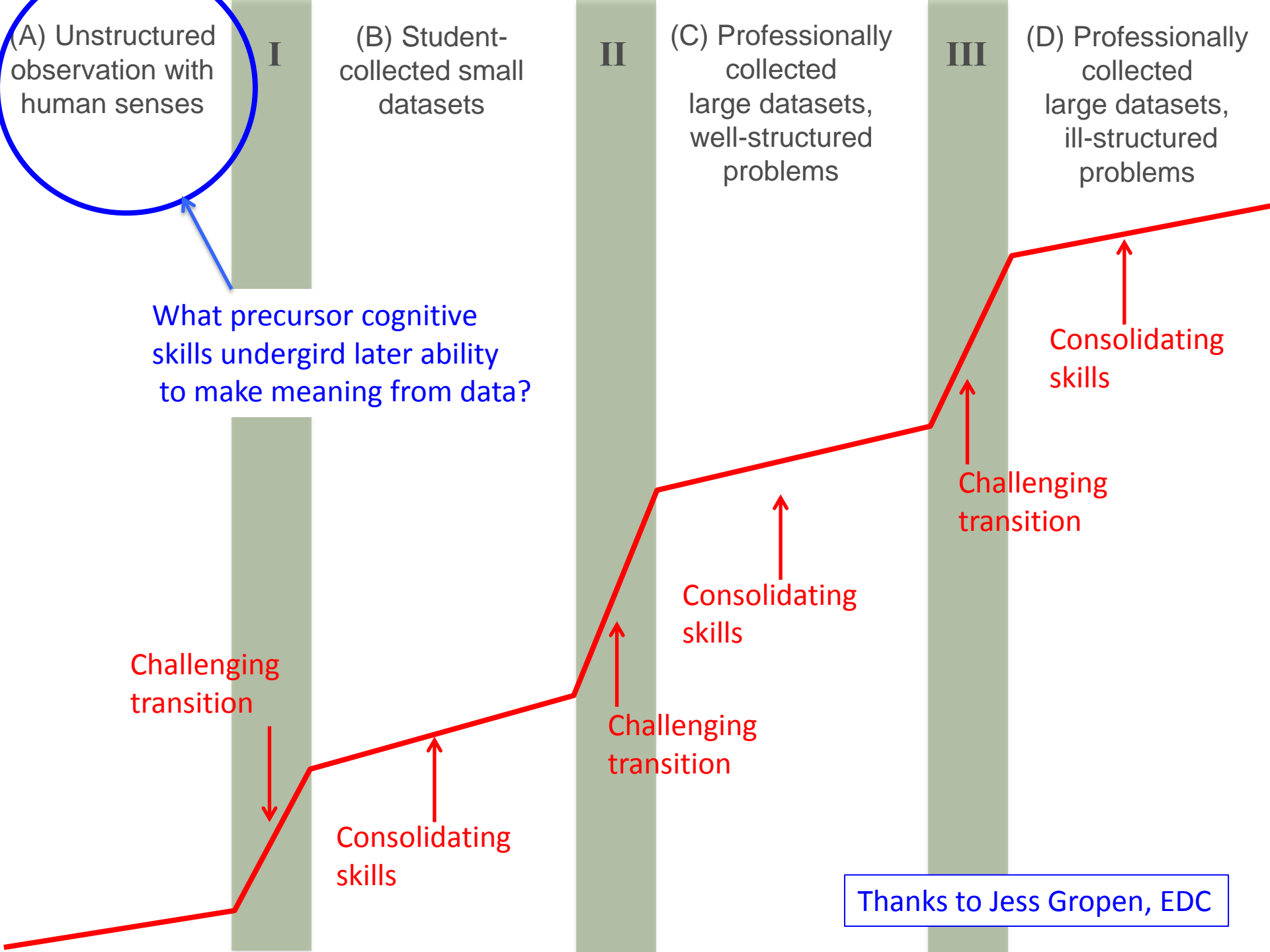
Presented at American Association for the Advancement of Science
San Jose, CA; February 15, 2015

Learning Science from Scientific Data: Why bother?

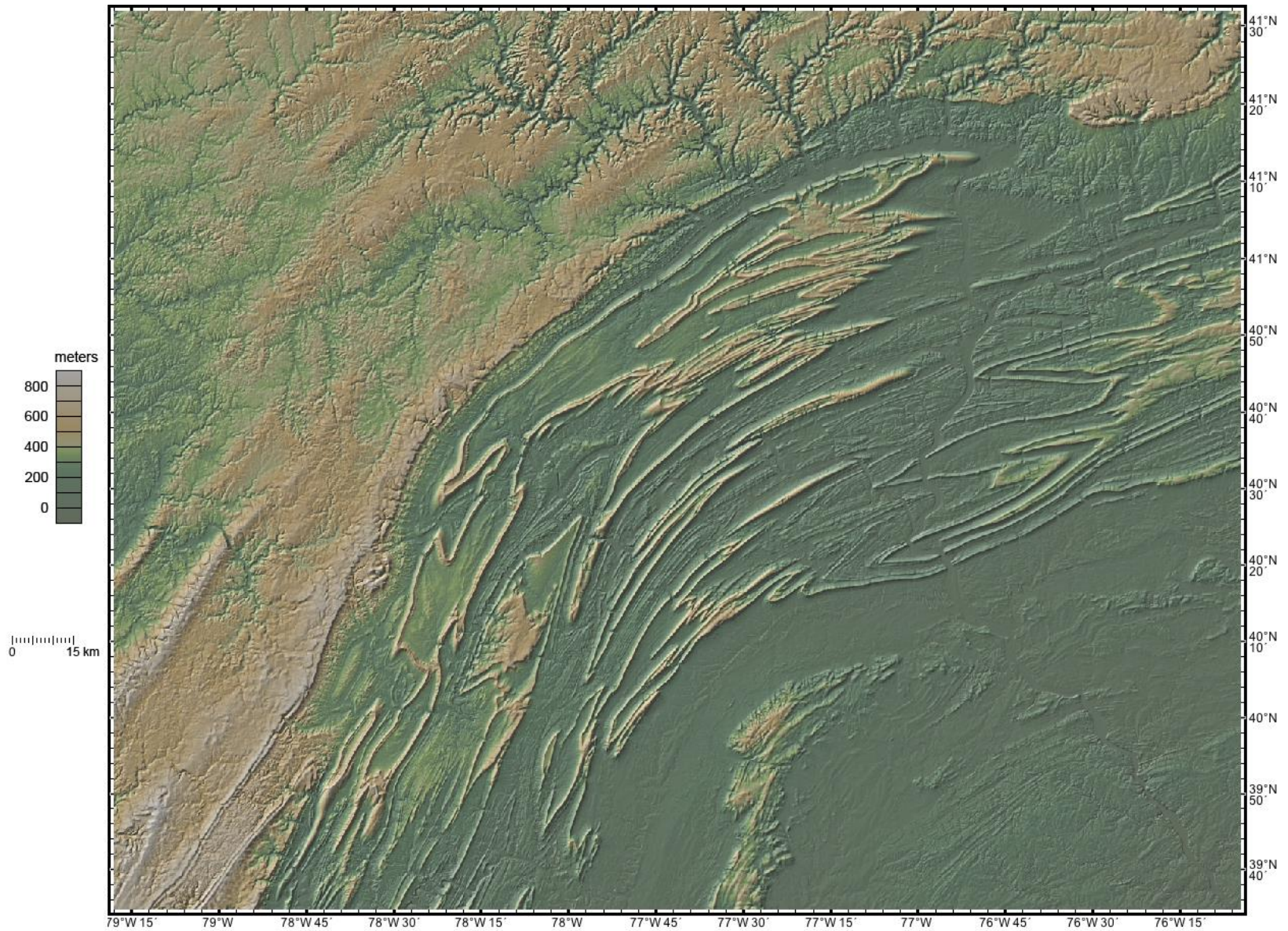
Reason #1: Students can grasp the evidence base that underlies the big ideas of science, rather than having to take these ideas on authority.

Reason #2: The world faces tough decisions and society is making some bad decisions. We want to raise up a generation who have the skills and disposition to make decisions based on evidence.





Events leave traces, and by looking at the traces we can sometimes make inferences about those events.



Abduction or abductive inference

Given: a set of specific observed facts

Find: one or more explanations that are consistent with the observed facts, using knowledge of the system to constrain the hypothesis space, and making plausible assumptions.



No causal inference:
*"don't know" or
"can't tell" or
orthogonal comment*

Single concrete
adamant hypothesis:
"The cat knocked it over"

Two concrete
alternative hypotheses:
*"Maybe the cat or maybe
my baby brother"*

The "alternative
working hypo-
theses" of the
historical sciences

Hypothesis with
bounded variable:
*"...the thing that knocked it
over was bigger than a
feather and ..."*

The organizational
strategy of certain
scientific computer
models, including
climate models.

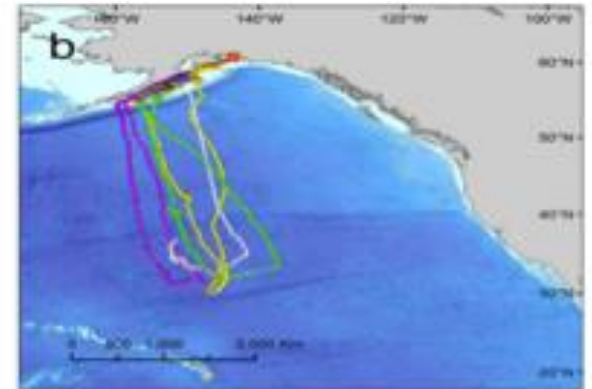
Everyday life



Middle school curriculum



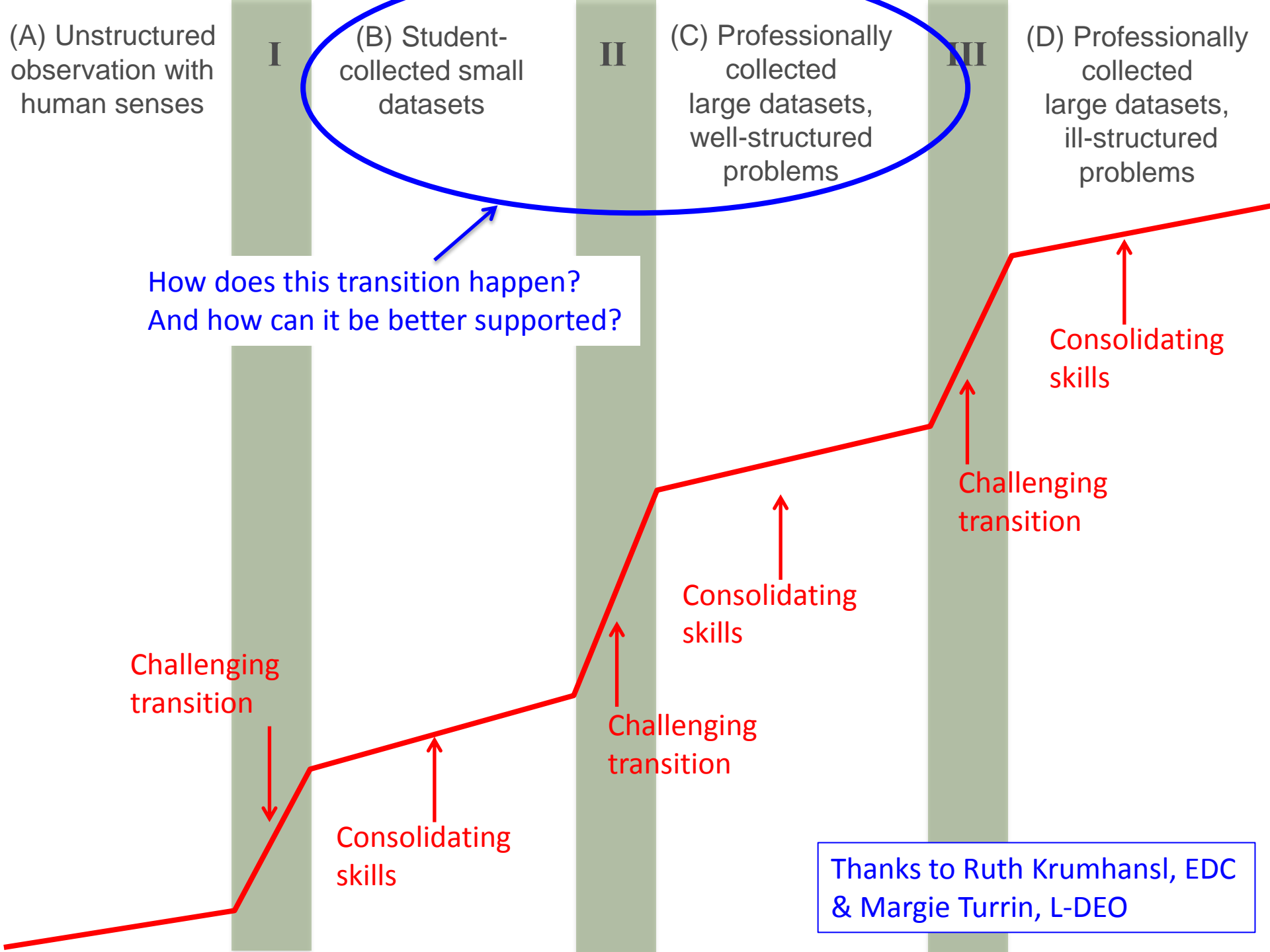
Scientists' data



- more elements & complexity
- farther removed from direct experience
- requiring more learned knowledge

Precursor understandings supporting data interpretation

- Events leave traces; by looking at traces you can sometimes make inferences about the events
- Form follows function (sometimes); Form reflects formative processes (sometimes)
- Sequence constrains causality: If A happened before B, A can have caused or influenced B, but not vice versa
- others?



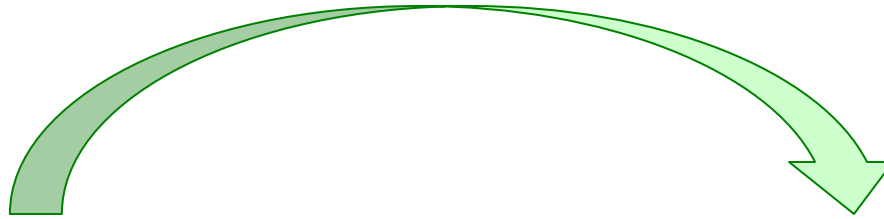
Most data experiences have been with small, student-collected data sets



Chapter 6

Science Content Standards

What is involved in this transition?



Student-collected data

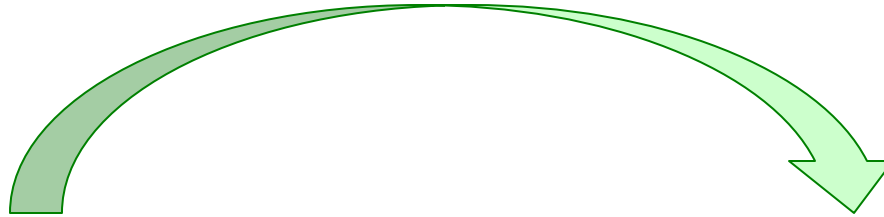


Day in the Life of the Hudson

Professionally-collected data



Kim aboard *Joides Resolution*, Leg 107



Embodied, experiential
grasp of the natural setting
and data collection methods



(from School in the Forest powerpoint,
<http://www.blackrockforest.org/docs/about-the-forest/schoolintheForest/>)

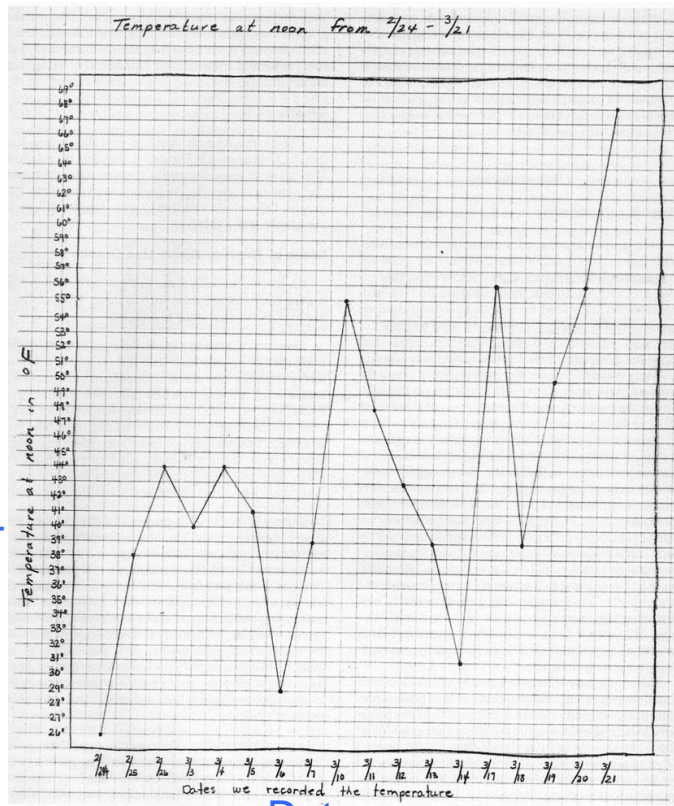
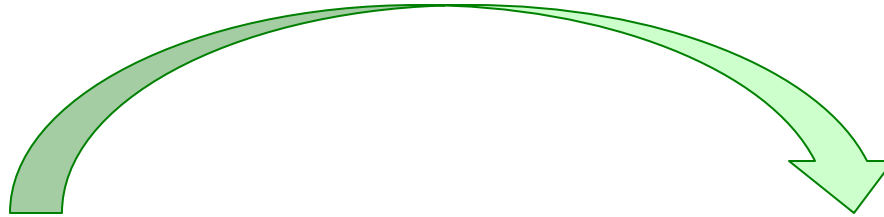
Metadata



(from Using a Digital Library to Enhance Earth Science Education,
Rajul Pandya, Holly Devaul, and Mary Marlino)

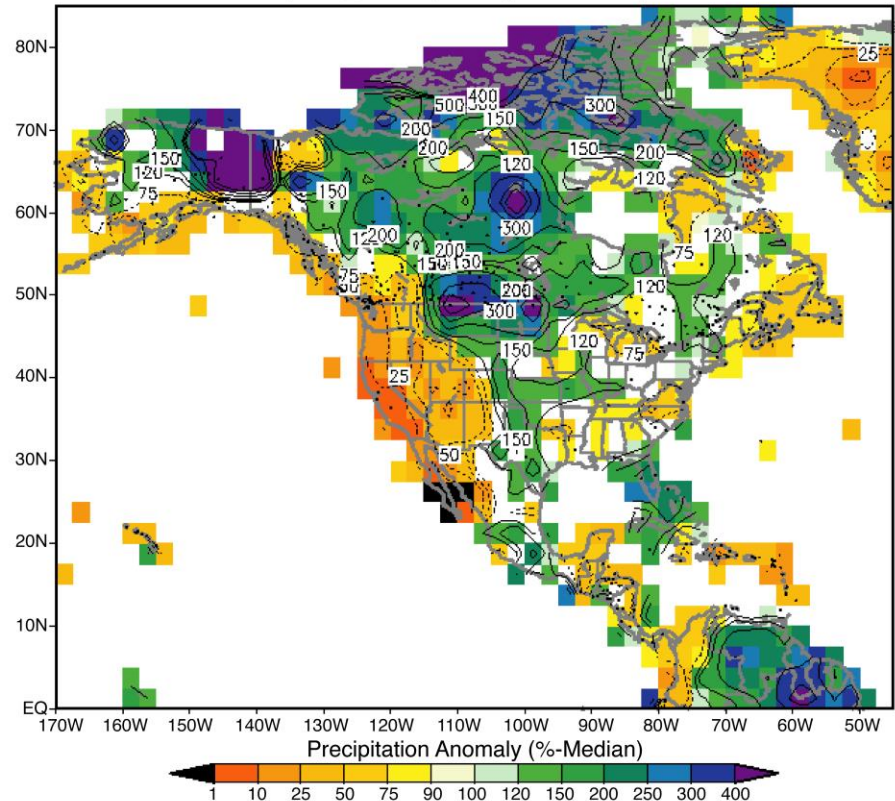
Dozens of data points

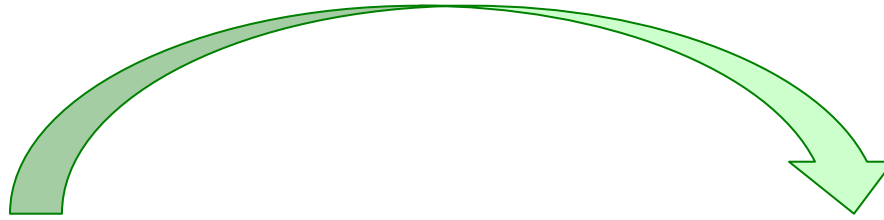
Megabytes



(from Clement, 2002)

Observed Precipit. Anomaly OND 2002
Shaded ONLY for "ABOVE-Normal" & "BELOW-Normal"
[CAMS_OPI data, courtesy of NCEP/CPC]

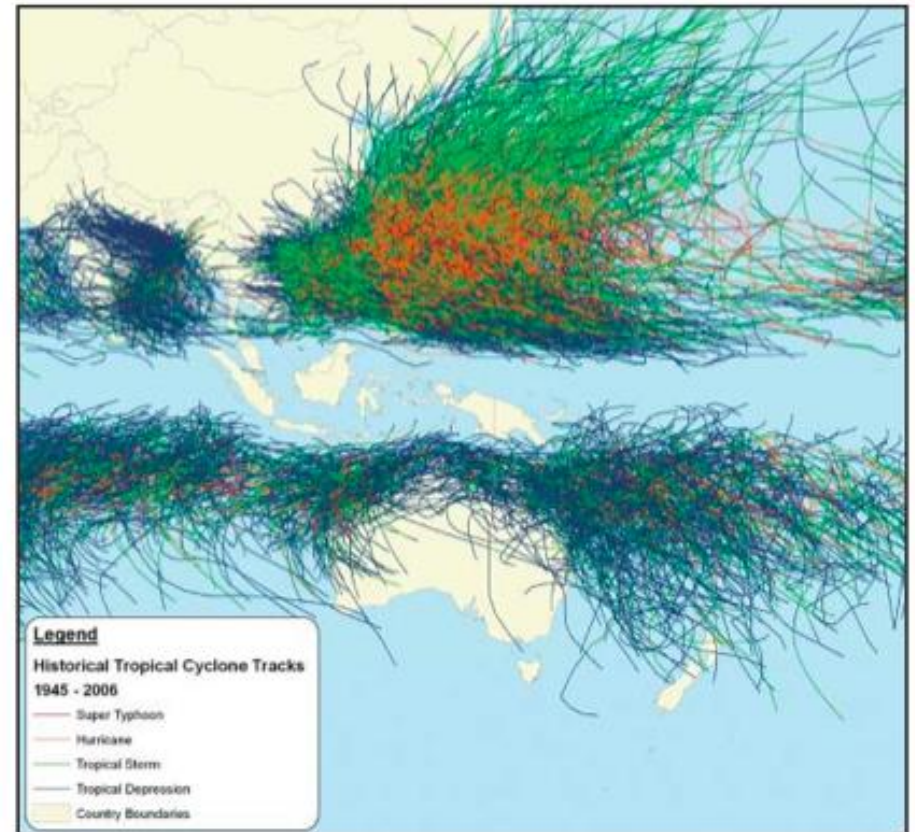
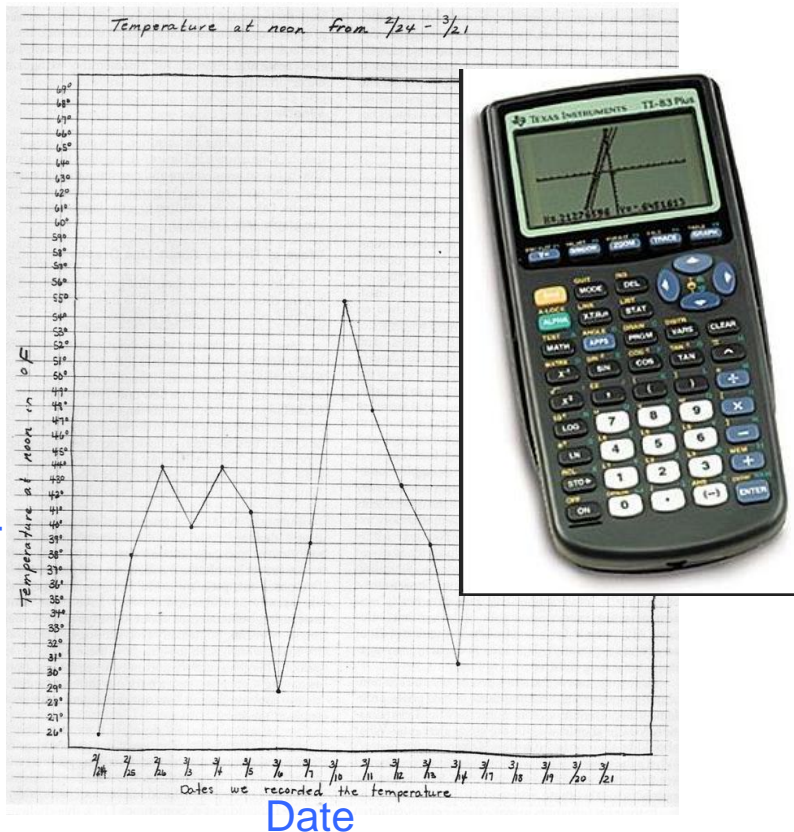


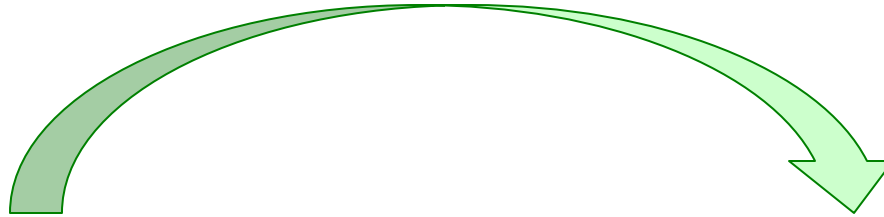


Simple, transparent
tools and techniques

Sophisticated tools &
techniques

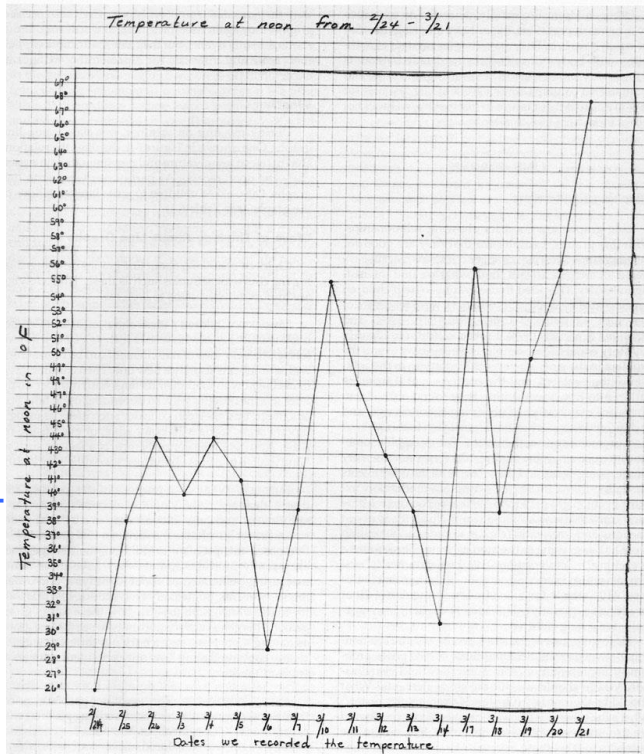
Air temperature at noon



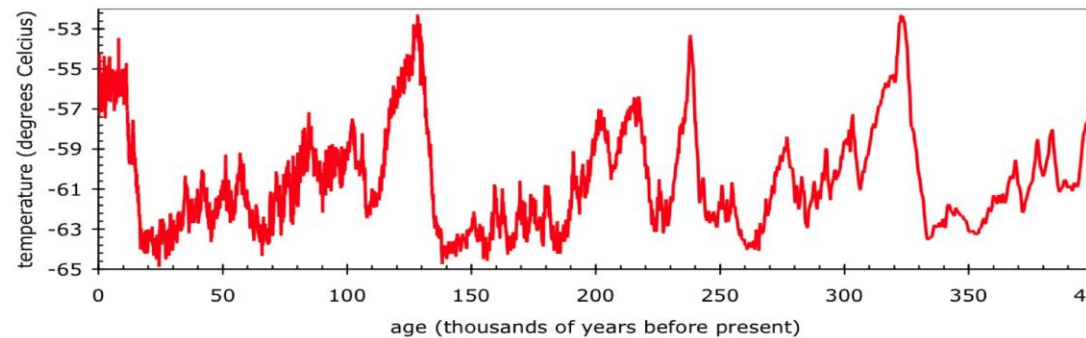
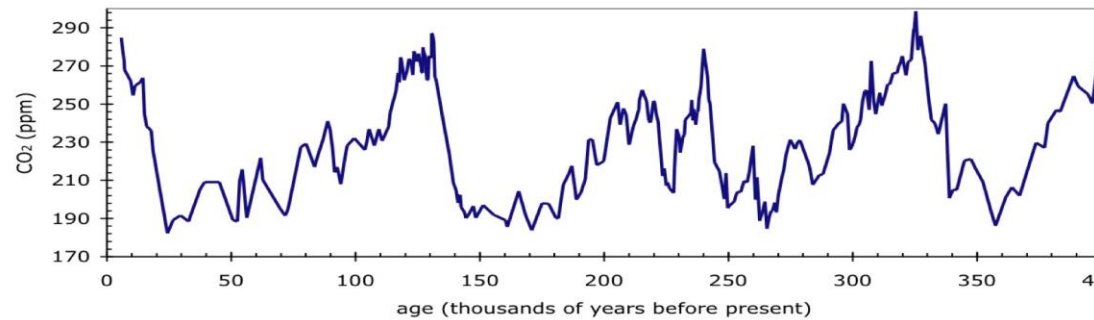


Interpret one data set at
a time

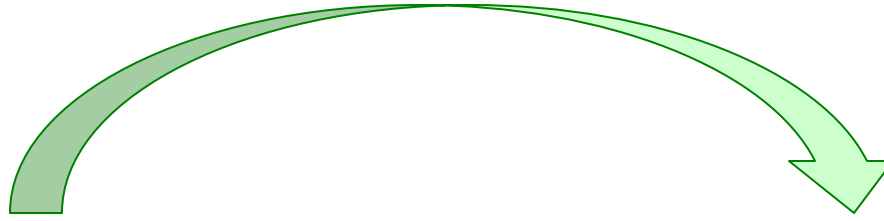
Multiple data sets with interactions;
varying data types



(from Clement, 2002)

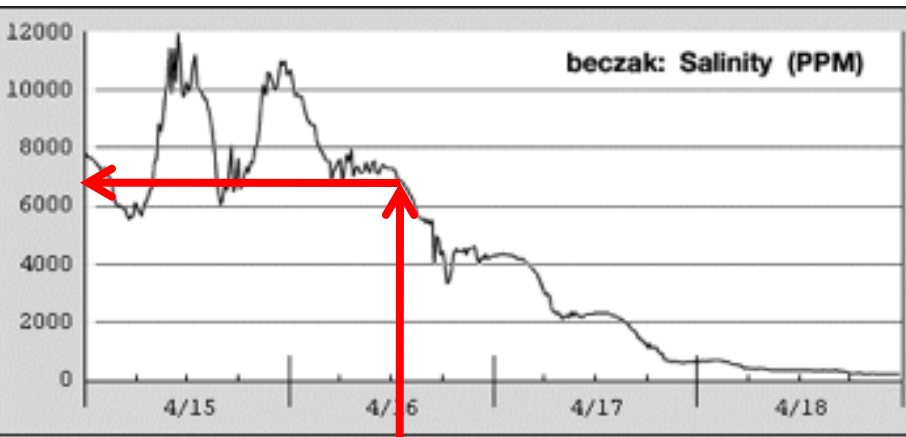


Date

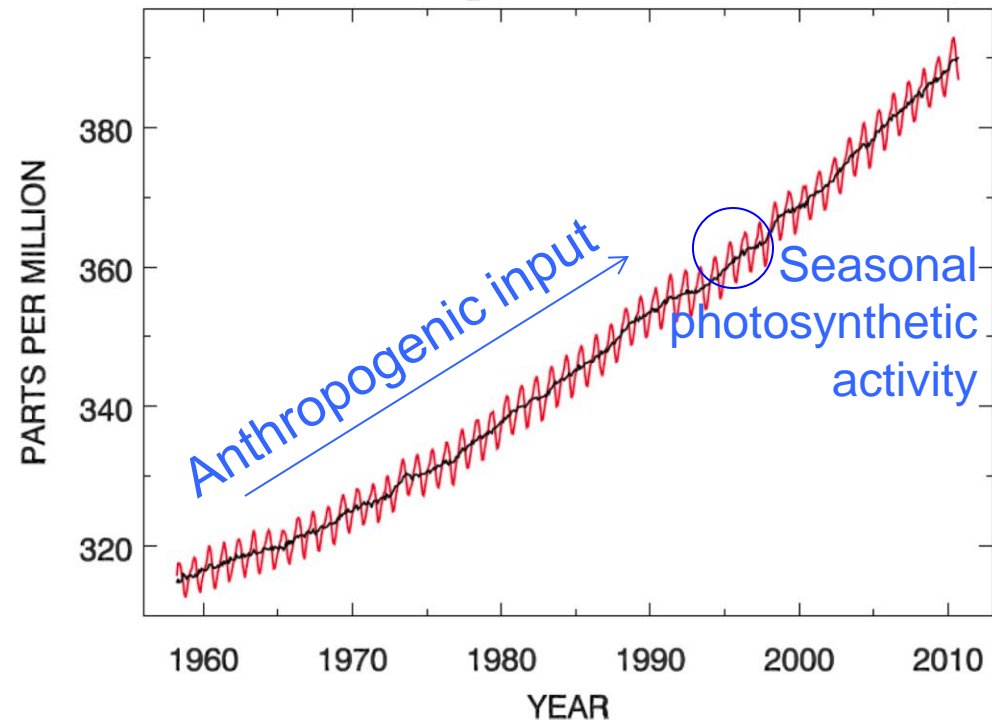


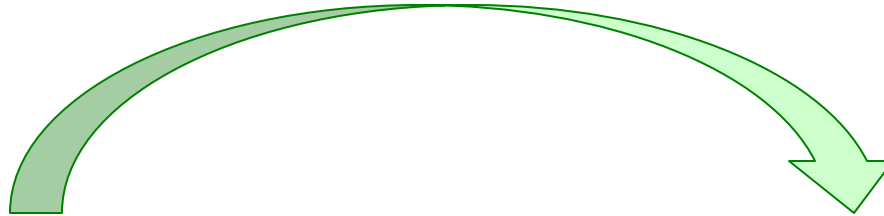
Looking up values

Seeing and interpreting patterns



Atmospheric CO₂ at Mauna Loa Observatory



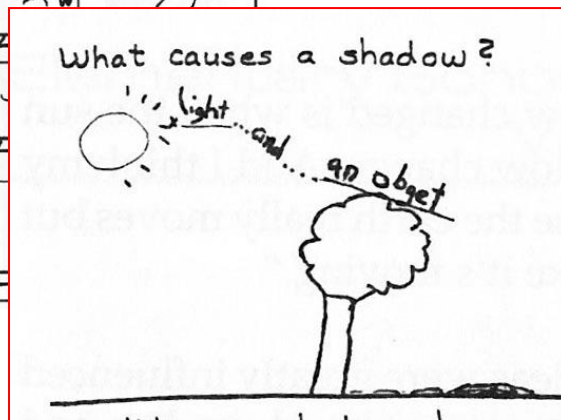


Common sense
lines of reasoning

Spatial, temporal, statistical
reasoning. Multi-step
chains of reasoning

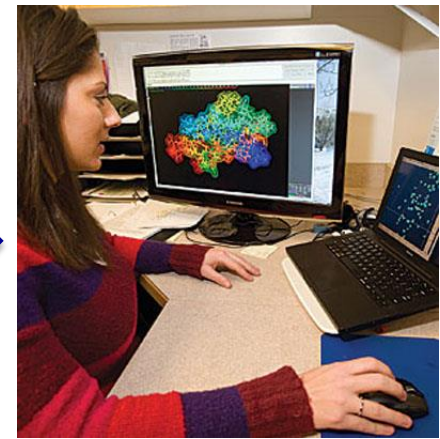
page time place

Time	Shadow Length	Position of Sun	Position of Shadow
9:15	129 inches		
11:00	78 inches		
12:15	68		
1:20	67		
2:30	76		



Using logic diagrams to
organize information





- Student-collected data
- Embodied, experiential sense of circumstances
- Dozens to hundreds of data points
- Simple, transparent tools & techniques
- Interpret one data set at a time
- “Common sense” lines of reasoning
- Single step causal chains

- Professionally-collected data
- Sense of circumstances from metadata
- Megabytes
- Complex tools & techniques; black boxes
- Multiple data sets and their interactions
- Temporal, spatial, quantitative and other lines of reasoning
- Multi-step lines of reasoning

Ways to scaffold students' transition from small, student-collected datasets to large, professionally-collected data bases

- *“Data puzzles”*: Use pre-selected snippets of high insight:effort ratio data
- *Nested Datasets*: Position a small student-collected dataset within a larger dataset.
- *Prediction*: Ask students to commit to a prediction of what they will see before they start making data visualizations.
- *Hypothesis array*: Provide an suite of candidate hypotheses; students seek the one best supported by the data.

Oceans of Data Institute Instructional Sequence Template #2

Nested data sets

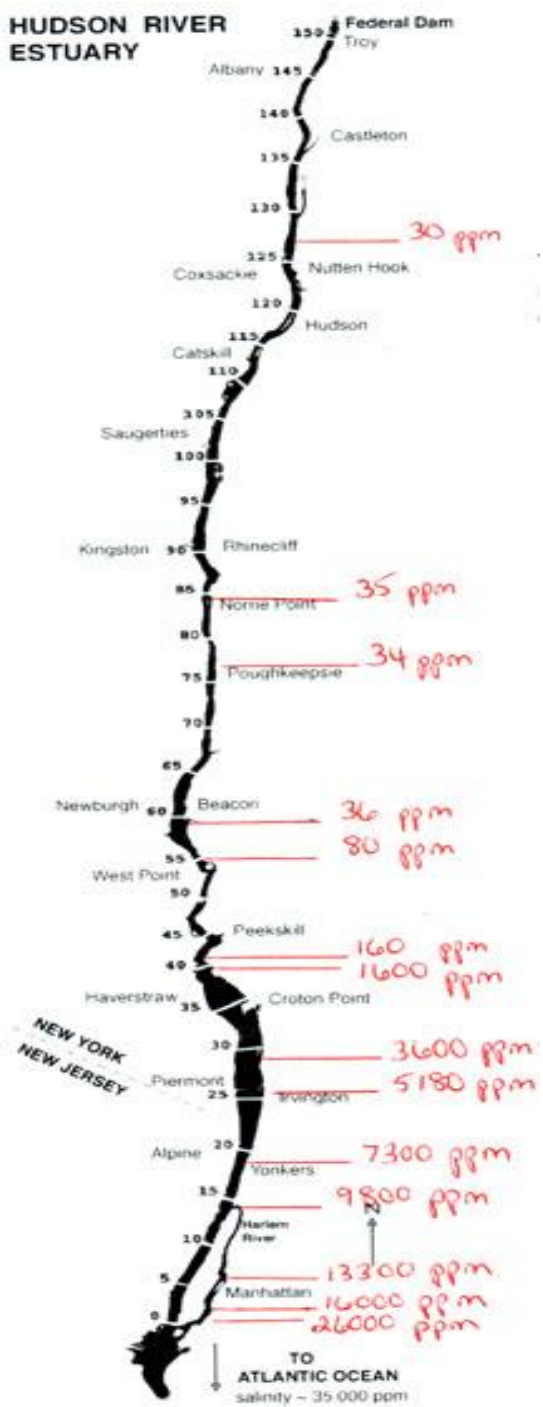
Procedure:

- 1) Students collect and interpret a local data set.
- 2) (optional) Students from multiple schools combine similar datasets to aggregate a larger sample or span a larger area.
- 3) Students interpret larger professionally collected dataset(s) which encompass and expand beyond the circumstances of their self-collected dataset.

A Day in the Life of the Hudson:

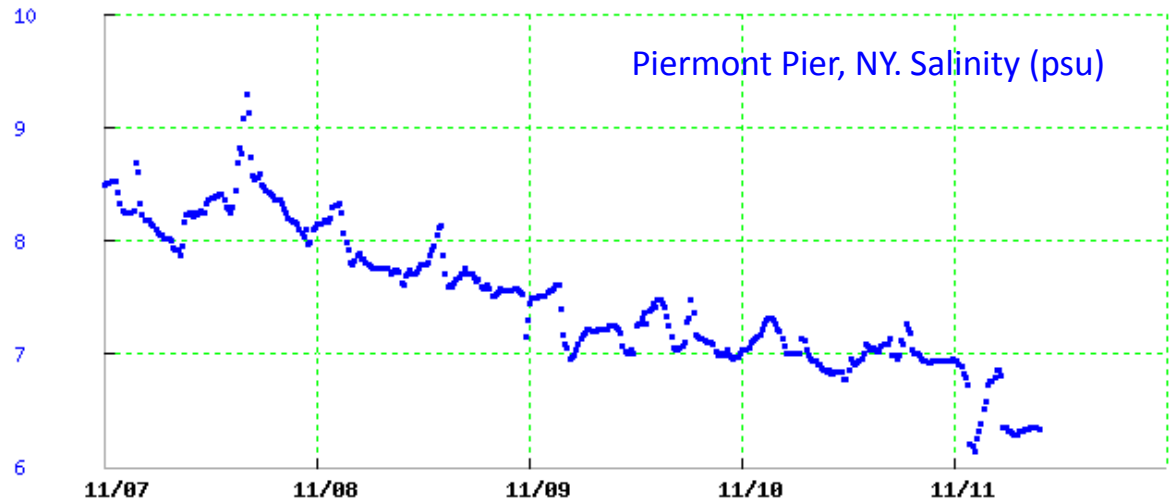


HUDSON RIVER ESTUARY



Combine with other school groups' data to explore variation across space.

Combine with professionally collected data to explore changes through time.

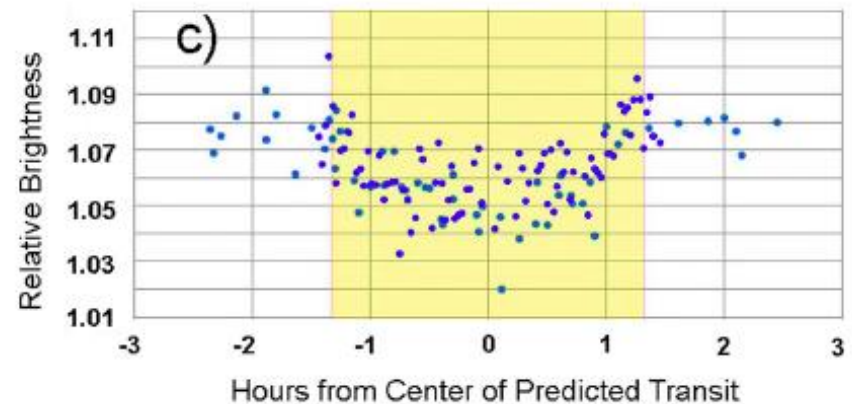
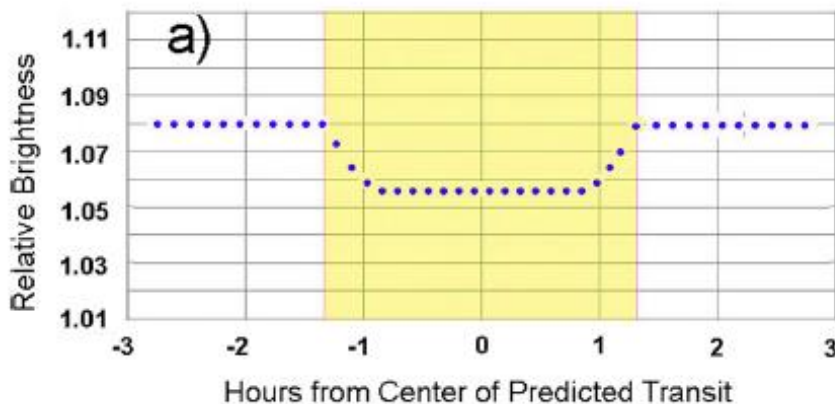
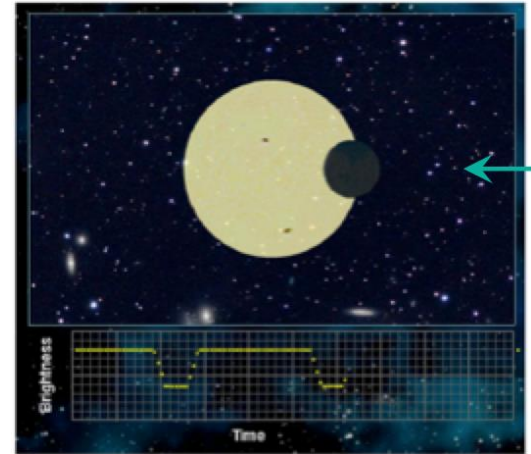


(from Turrin, M., & Kastens, K. A. (2010). In *Earth Science Puzzles: Making Meaning from Data* and <http://www.hrecos.org/>)

Prediction

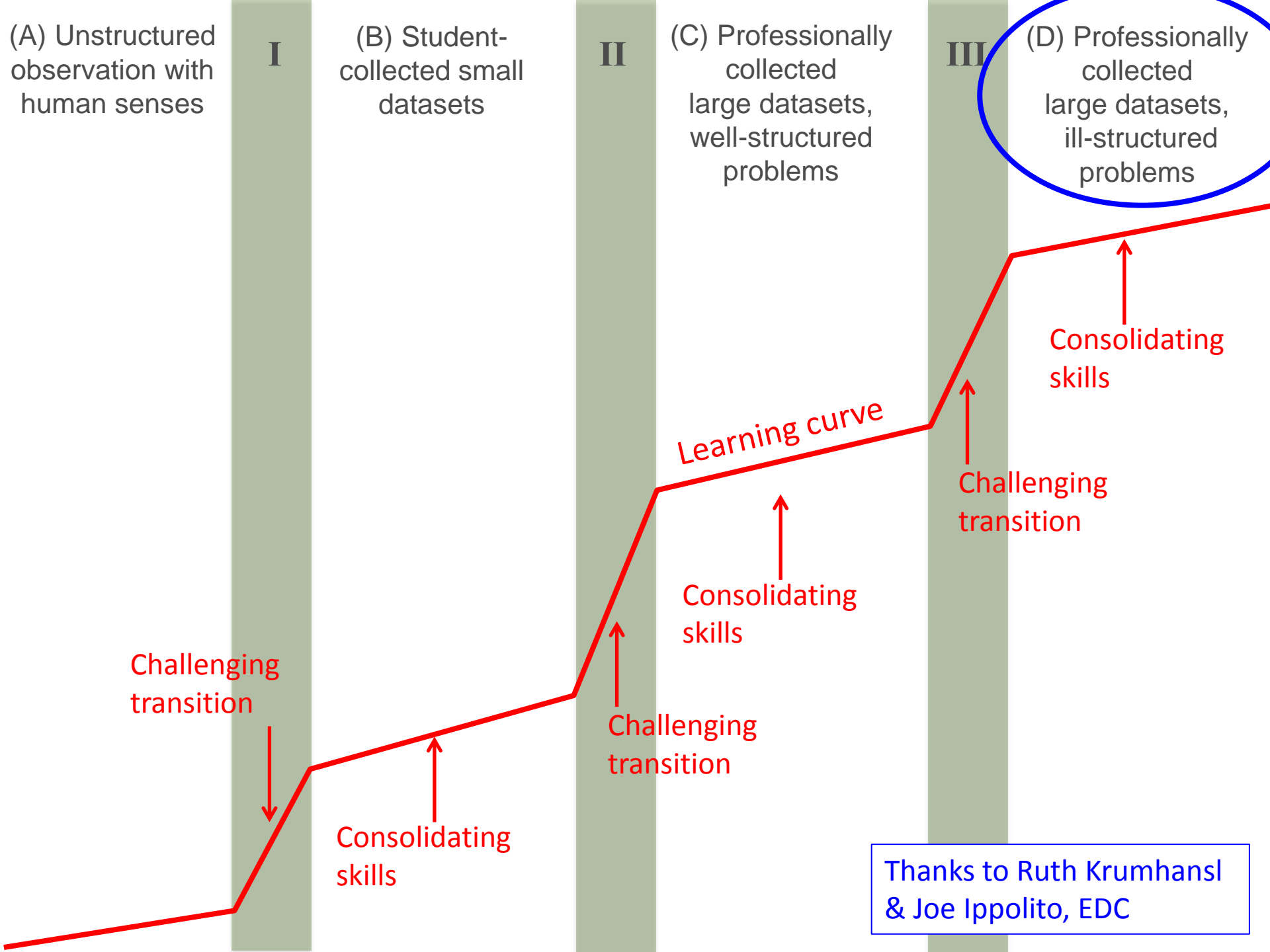
Procedure:

- 1) Based on either a conceptual model, physical model or computational model, students predict what data from the system under consideration would look like under various conditions.
- 2) Students examine professionally collected data taken under a range of conditions, looking for the presence or absence of predicted patterns.



Instructional templates or “design patterns”

- Can be reused in different contexts with different content, building capacity in teachers and students
- Can serve as the basis for a research agenda
 - research affordances and pitfalls of the strategy
 - rather than evaluating each bit of instructional materials separately



DACUM: a process for Developing A CurriculUM

A well-established methodology for occupational analysis, modified by EDC for emergent professions

Premise: experienced and respected practitioners can best define and describe their job or profession

Product:

- Definition of the job/career/profession
- Duties & Tasks
- Knowledge, Skills, Tools & Behaviors

Expert Panel:
Aug 14-15, 2014

Kartik Shah
Strategix Solutions

Ryan Kapaun
Eden Prairie Police
Department

Shannon McWeeney
Oregon Health & Science
University

Juan Miguel Lavista Ferres
Bing/Microsoft

Tim Chadwick
Dynamic Network
Services, Inc.

Steve Ross
Broadband Communities
Magazine

Randy Bucciarelli
Scripps Institution of
Oceanography
UC San Diego

Benjamin Davison
Google

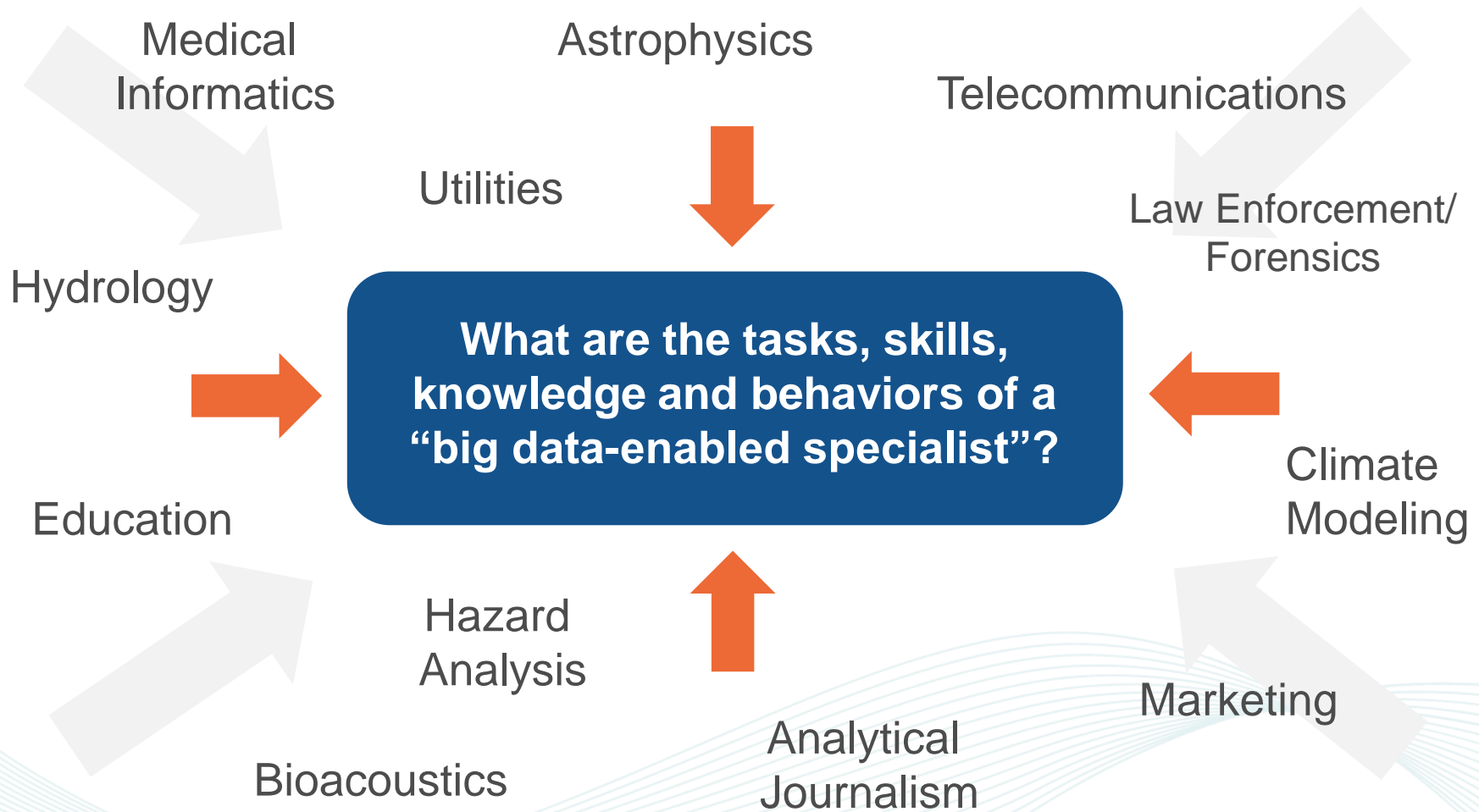


Lucy Drotning
Columbia University

Jay Parker
Jet Propulsion Laboratory
California Institute of
Technology

Kirk Borne
George Mason University

Developing an Occupational Profile



Joseph Ippolito, EDC

Data

- writes software

- data visualizations

- relevant data based on

descriptions

- ability control of what is

data

- data from different sources

- is misleading/questionable

data

- mis. the data

- between problem statement & data

- creates a data dictionary

- stores data

- designs workflow

- implements workflow

- stores the data

- organizes the data

- conducts data explor

- identifies tools that

needed for projects

- cleans data for dect

- detects discrepancies of

- maps heterogeneous dat

- standardizes heterogene

- collects data

- secures data/results

- protects data/results

- creates meta data and

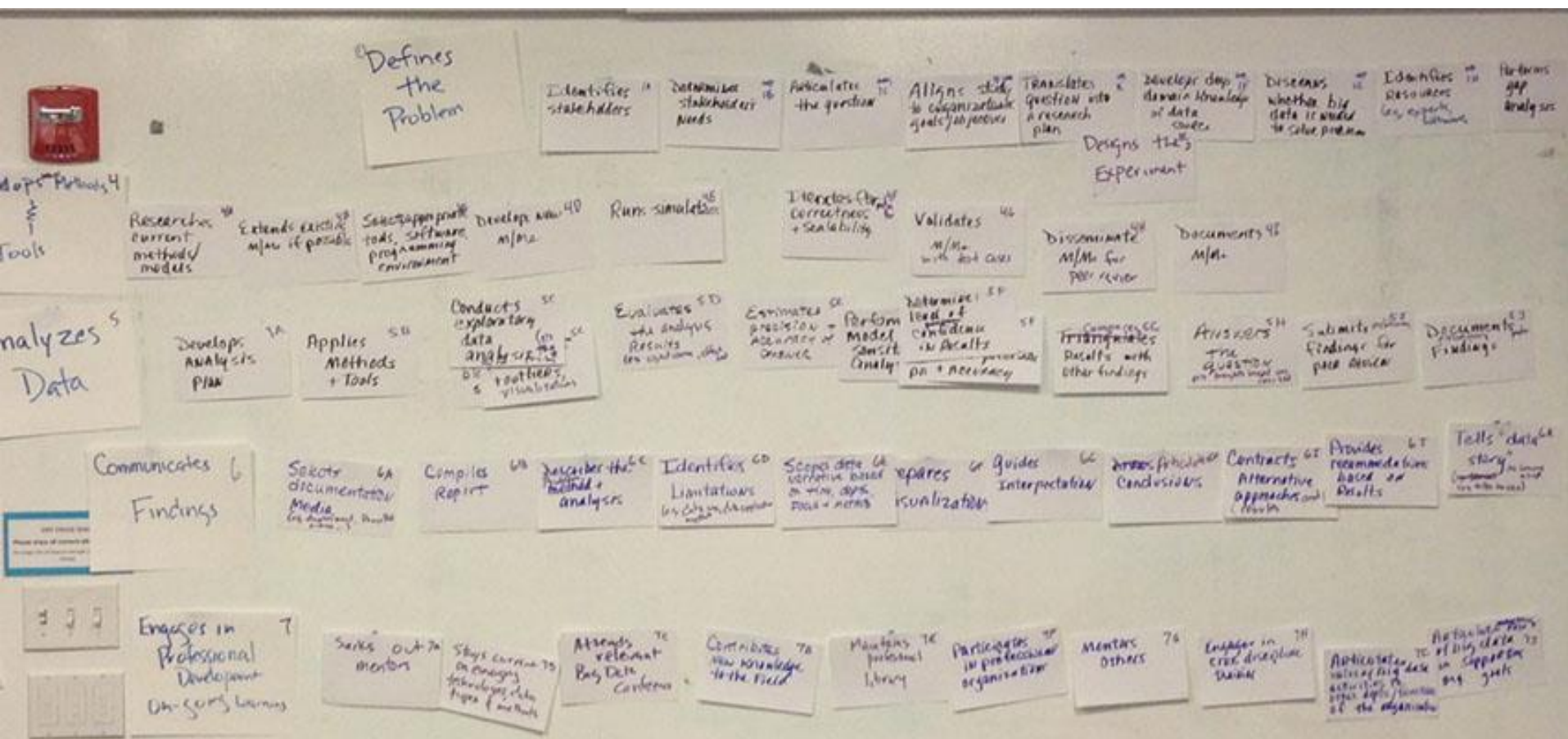
that describe the d



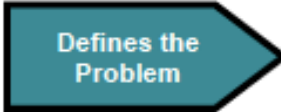



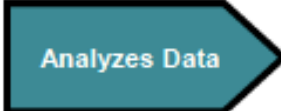
Occupational Definition

The Big-Data-Enabled Specialist is an individual who wrangles and analyzes large and/ or complex data sets to enable new capabilities including discovery, decision support and improved outcomes.

Duties & Tasks



Occupational Profile

DUTIES		TASKS					
1.  Defines the Problem	1A. Identifies stakeholders	1B. Determines stakeholders' needs	1C. Articulates the question	1D. Aligns study to organizational goals and objectives	1E. Translates question into a research plan	1F. Designs the experiment	1G. Develops deep domain knowledge of data source
	1M. Negotiates plan, including deadlines and budgets	1N. Creates requirement document (sign-off)					
2.  Wrangles Data	2A. Performs data exploration	2B. Identifies data	2C. Creates the data dictionary	2D. Collects data	2E. Assesses the extent/methods to clean the data	2F. Maps data across heterogeneous sources	2G. Identifies outliers and anomalies
	2M. Writes software to automate tasks	2N. Documents the process					
3.  Manages Data Resources	3A. Manages data life cycle	3B. Conducts capacity planning of resources	3C. Complies with legal obligations	3D. Applies ethical standards	3E. Identifies tools that may be needed for purchase or modification	3F. Protects data and results	3G. Determines access to the data
4.  Develops Methods and Tools	4A. Researches current methods/models	4B. Extends existing methods/models, if possible	4C. Selects tools/software/programming environment	4D. Develops new methods/models	4E. Runs simulations	4F. Iterates correctness and scalability of methods/models	4G. Validates methods/models with test cases
5.  Analyzes Data	5A. Develops analysis plan	5B. Applies methods and tools	5C. Conducts exploratory analysis (e.g., identifies anomalies, outliers, bias in sampling; visualizes)	5D. Evaluates results of the analysis (e.g., significance, effect, size)	5E. Estimates precision and accuracy of answer	5F. Determines level of confidence in results	5G. Compares results with other findings

Etc.

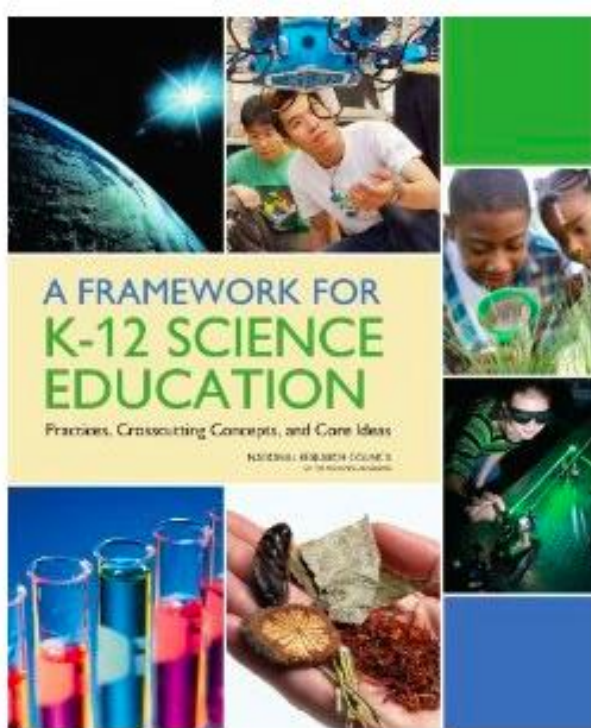
Etc.

Major work responsibilities- Duties

- 1) Defines the problem
- 2) Wrangles data
- 3) Manages data resources
- 4) Develops methods and tools
- 5) Analyzes data
- 6) Communicates findings
- 7) Engages in professional development

Gap analysis

How well is the current education system doing at preparing students for the tasks and duties of the big-data-enabled specialist?



- Disciplinary Core Ideas
- Cross-cutting Themes
- Practices of Science & Engineering
 - #4: Analyze & interpret data

(aspirational!)

Comparison of ODI occupational profile tasks with NGSS Performance Expectations

DUTIES		TASKS										
5.	Analyzes Data	5A. Develops analysis plan	5B. Applies methods and tools	5C. Conducts exploratory analysis (e.g., identifies anomalies, outliers, bias in sampling; visualizes)	5D. Evaluates results of the analysis (e.g., significance, effect, size)	5E. Estimates precision and accuracy of answer	5F. Determines level of confidence in results	5G. Compares results with other findings	5H. Answers the question (e.g., insights drawn from results)	5I. Submits preliminary findings for peer review	5J. Documents preliminary findings	
		6A. Selects documentation media (e.g., dashboard, PowerPoint, e-mail)	6B. Compiles report	6C. Describes problem, method, and analysis	6D. Identifies limitations (e.g., data use, data application methods)	6E. Scopes data narrative based on time, depth, and method	6F. Prepares visualizations	6G. Guides interpretation	6H. Articulates conclusions	6I. Contrasts alternative approaches and past results	6J. Provides recommendations based on results	6K. Tells "data story" to convey insight (e.g., talks to CEO)
7.	Engages in Professional Development	7A. Seeks out mentors	7B. Stays current on emerging technologies, data types, and methods	7C. Attends relevant big data conferences	7D. Contributes new knowledge to the field	7E. Maintains professional library	7F. Participates in professional organizations	7G. Mentors others	7H. Engages in cross-discipline training	7I. Articulates value of big data activities to other departments/ functions of organization	7J. Articulates evolving role of big data in supporting organizational goals	

Abundant in NGSS
 Potentially (implicitly) abundant in NGSS
 Sparse in NGSS
 Absent from NGSS

Occupational Profile tasks that are *well-represented* in NGSS

1. Defines the Problem

- 1B. Determines stakeholders' needs
- 1C. Articulates the question
- 1E. Translates question into a research plan
- 1F. Designs the experiment
- 1G. Develops deep domain knowledge of data source

2. Wrangles Data

- 2D. Collects data

5. Analyzes Data

- 5A. Develops analysis plan
- 5B. Applies methods and tools
- 5D. Evaluates results of the analysis (e.g., significance, effect, size)
- 5H. Answers the question (e.g., insights drawn from results)

Occupational Profile tasks that are *absent from* NGSS

2. Wrangles Data

- 2A. Performs data exploration
- 2G. Identifies outliers and anomalies
- 2N. Documents the process

3. Manages Data Resources

- 3D. Applies ethical standards
- 3F. Protects data and results

4. Develops Methods and Tools

- 4F. Iterates correctness ... of ... models

5. Analyzes Data

- 5F. Determines level of confidence in results

6. Communicates Findings

- 6D. Identifies limitations (e.g., data use, data application methods)

Bottom line:

- It's a long, complicated pathway to grow a populace that has the skills and disposition to use data as part of their tool-kit when confronted with a difficult question or problem.
- There are many interesting challenges along the way.
- Science education is where it's happening.