# DATA ANALYSIS
## AT SCALE

PETRA ISENBERG (slides by WESLEY WILLETT)

VISUAL ANALYTICS    04 NOV 2014

# DATA ANALYSIS AT SCALE

CHALLENGES

ANALYSIS AND CLUSTER COMPUTING

INTERACTING WITH BIG DATA

PARALLELIZING HUMAN INTELLIGENCE

# CHALLENGES FOR ANALYZING LARGE DATA SETS

**SIZE** **SPEED**

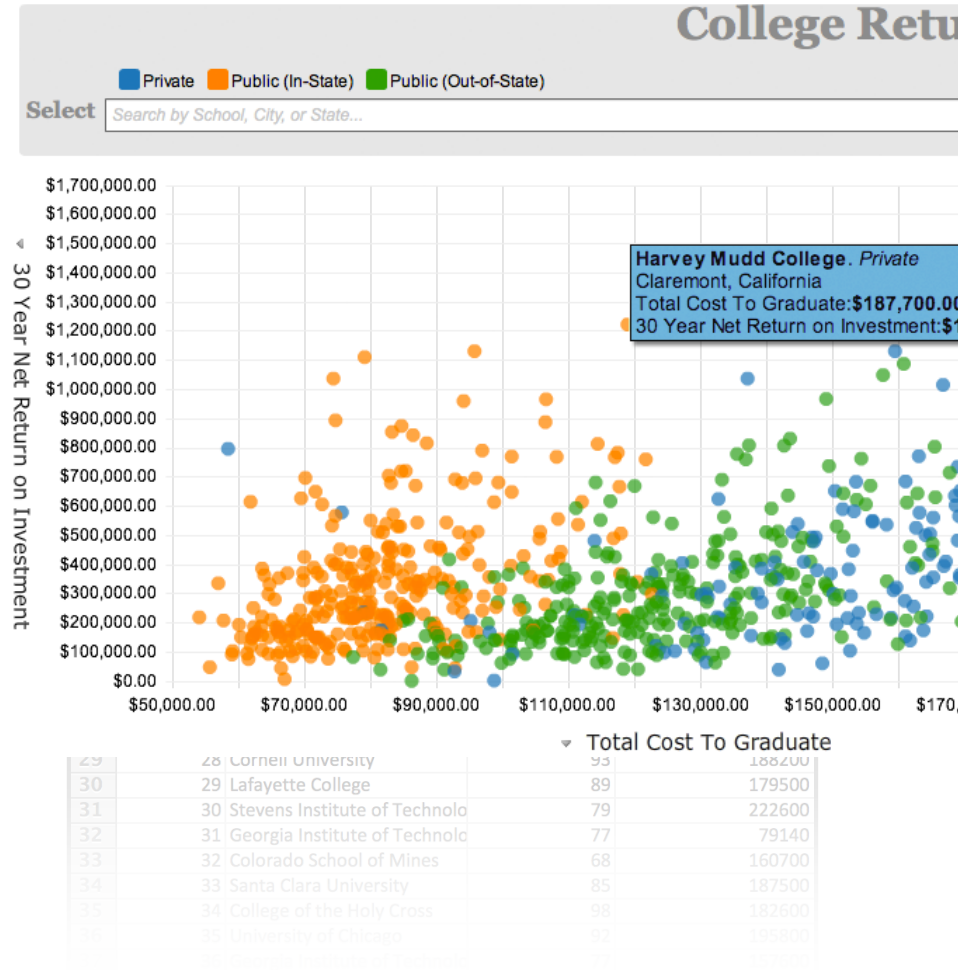**ATTENTION**

# SIZE

KILOBYTES OF DATA

MEGABYTES OF DATA

GIGABYTES OF DATA

TERABYTES OF DATA

PETABYTES OF DATA

...

# SIZE

**KILOBYTES OF DATA**

MEGABYTES OF DATA

GIGABYTES OF DATA

TERABYTES OF DATA

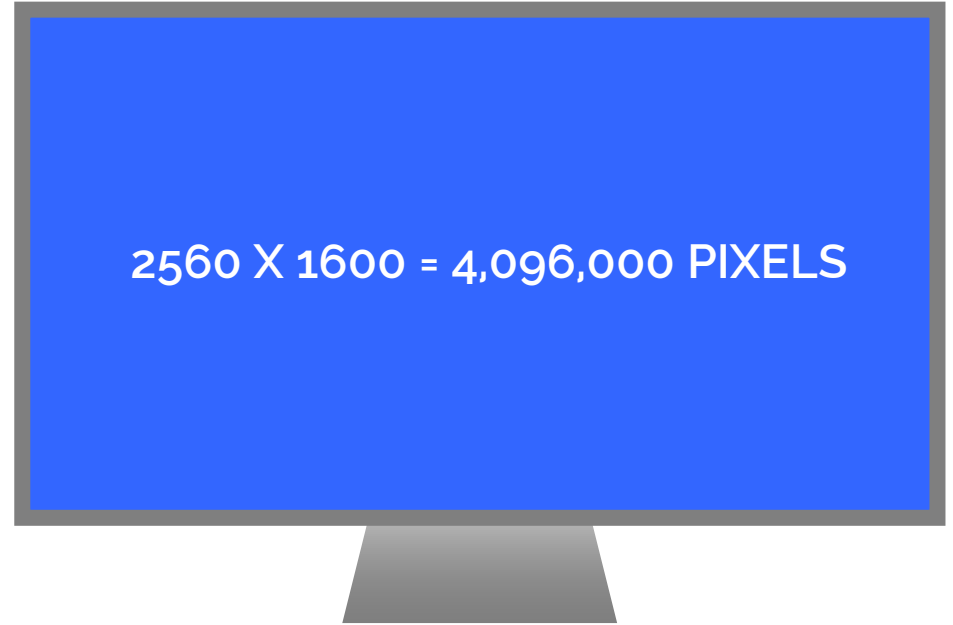PETABYTES OF DATA

...

# SIZE

**KILOBYTES OF DATA**

**MEGABYTES OF DATA**

GIGABYTES OF DATA

TERABYTES OF DATA

PETABYTES OF DATA

...

2560 X 1600 = 4,096,000 PIXELS

**EVEN A MEGABYTE IS MORE BITS OF DATA THAN THERE ARE <u>PIXELS ON A SCREEN</u>!**

# SIZE

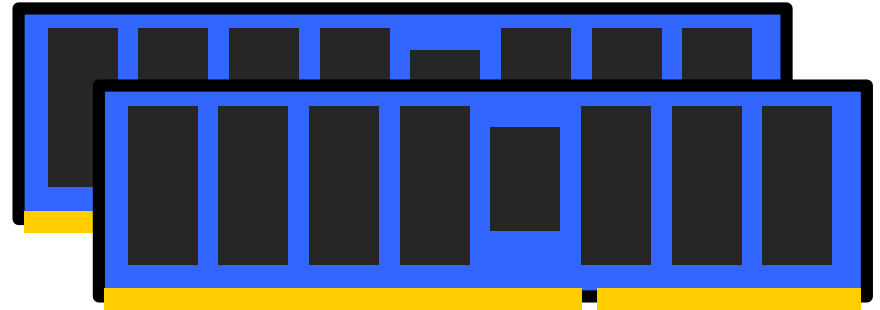**KILOBYTES OF DATA**

**MEGABYTES OF DATA**

**GIGABYTES OF DATA**

TERABYTES OF DATA

PETABYTES OF DATA

...

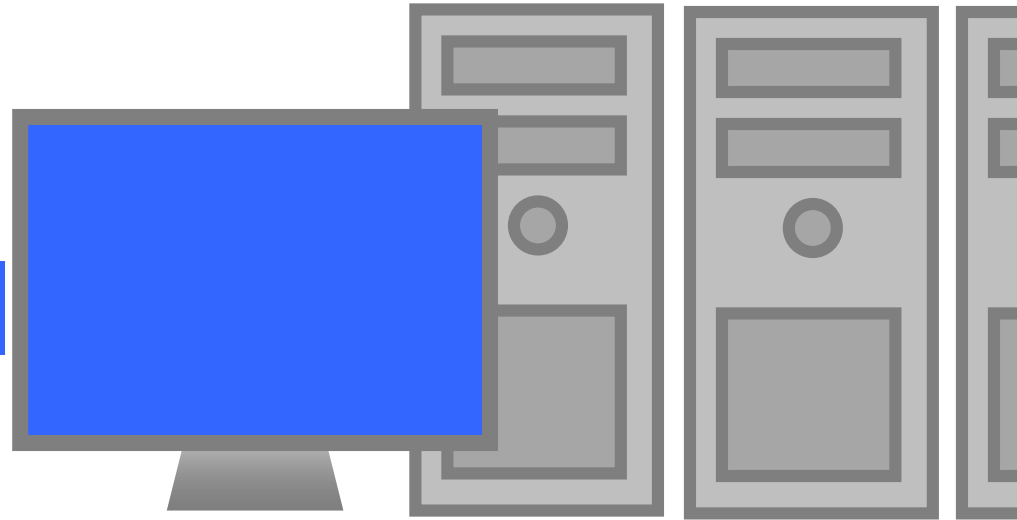# MORE DATA THAN CAN FIT <u>IN MEMORY</u>

# SIZE

KILOBYTES OF DATA
MEGABYTES OF DATA
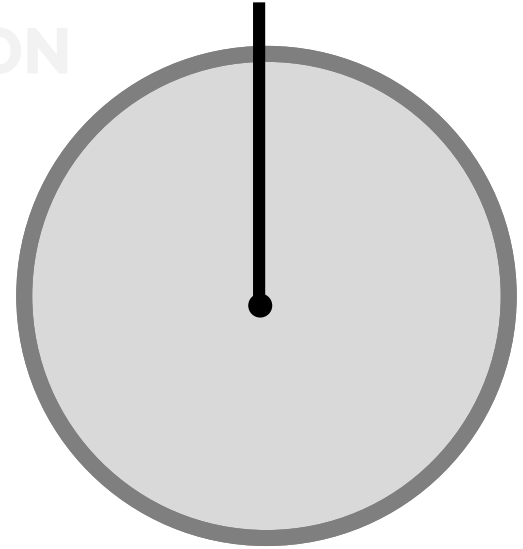GIGABYTES OF DATA
**TERABYTES OF DATA**
PETABYTES OF DATA

...

## MORE DATA THAN CAN FIT ON <u>ONE MACHINE</u>!

# SIZE

**KILOBYTES OF DATA**
**MEGABYTES OF DATA**
**GIGABYTES OF DATA**
**TERABYTES OF DATA**
**PETABYTES OF DATA**

**MANY BIG DATA-DRIVEN QUESTIONS TODAY**

# SPEED

~0.1 SECOND  DIRECT MANIPULATION

~1 SECOND  INTERACTIVE
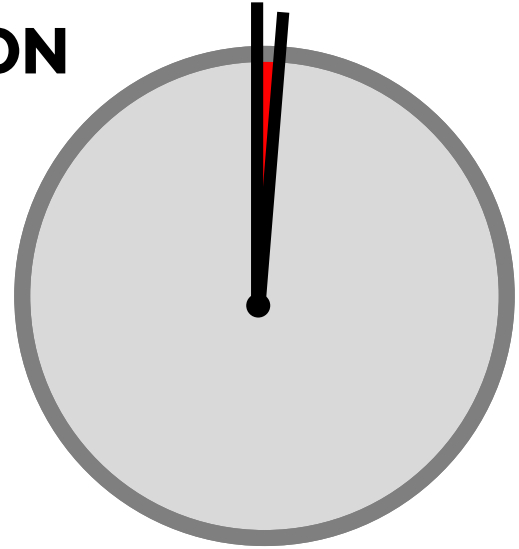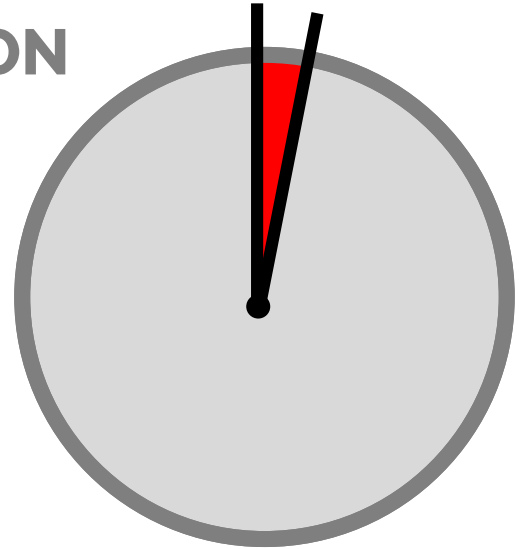
~10 SECONDS  QUERY / RESPONSE

MINUTES  ...

HOURS  BATCH PROCESSING
(VERY SLOW)

# SPEED

**~0.1 SECOND**    **DIRECT MANIPULATION**

~1 SECOND    INTERACTIVE
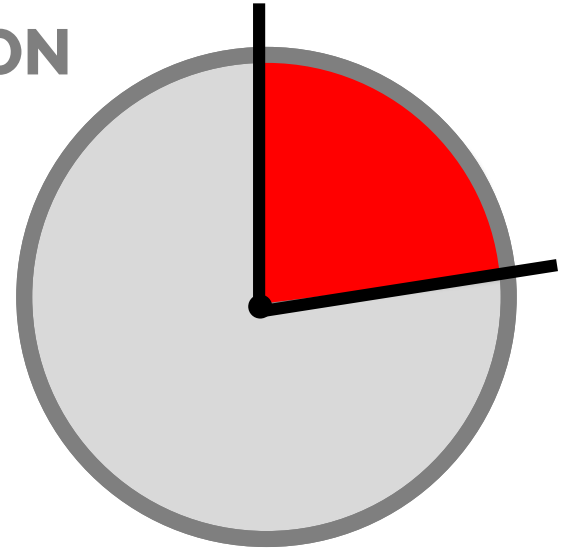
~10 SECONDS    QUERY / RESPONSE

MINUTES    ...

HOURS    BATCH PROCESSING (VERY SLOW)

# SPEED



~0.1 SECOND  DIRECT MANIPULATION
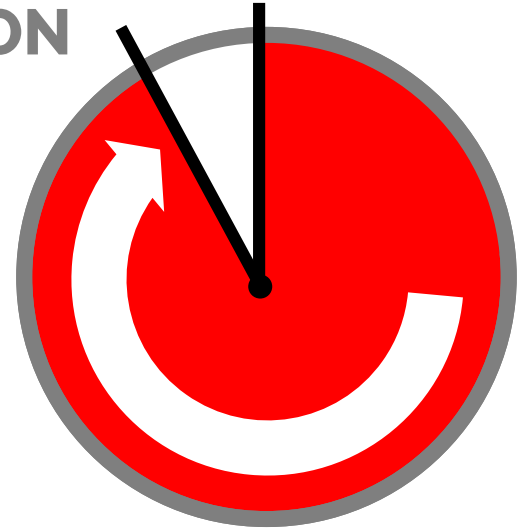
~1 SECOND  INTERACTIVE

~10 SECONDS  QUERY / RESPONSE

MINUTES  ...

HOURS  BATCH PROCESSING
(VERY SLOW)

# SPEED

| | |
|---|---|
| **~0.1 SECOND** | DIRECT MANIPULATION |
| **~1 SECOND** | INTERACTIVE |
| **~10 SECONDS** | QUERY / RESPONSE |
| MINUTES | ... |
| HOURS | BATCH PROCESSING (VERY SLOW) |

# SPEED

~0.1 SECOND    DIRECT MANIPULATION

~1 SECOND    INTERACTIVE

~10 SECONDS    QUERY / RESPONSE

**MINUTES**    ...

**HOURS**    **BATCH PROCESSING (VERY SLOW)**

# ATTENTION

EVERY PERSON ONLY HAS A FINITE
NUMBER OF WORKING HOURS

5-8 PERSON-HOURS PER DAY

1,489 PERSON-HOURS PER YEAR (FRANCE)

(1,388 GERMANY   2,163 IN S. KOREA   1,788 IN USA) [OECD STATS]

HOW LONG CAN YOU AFFORD TO SPEND FINDING EXAMPLES,
PROCESSING A DATASET, OR ANSWERING A QUESTION?

# ATTENTION

AN INDIVIDUAL ANALYST IS UNLIKELY TO BE ABLE TO SEE DATA FROM MANY PERSPECTIVES

**"MANY EYES FIND MORE BUGS"**

# DATA ANALYSIS AT SCALE

CHALLENGES

ANALYSIS AND CLUSTER COMPUTING

INTERACTING WITH BIG DATA

PARALLELIZING HUMAN INTELLIGENCE

# ANALYSIS & CLUSTER COMPUTING

BIG DATASETS ARE LIKELY TO BE
SPREAD OUT ACROSS A CLUSTER (OR
CLUSTERS)

ANALYSIS REQUIRES
DISTRIBUTED DATA PROCESSING

# HOW CAN WE PERFORM ANALYSIS ACROSS A CLUSTER?

How can we split work across machines?

# MAP-REDUCE

Jeffrey Dean and
Sanjay Ghemawat
(Google) 2004

# A SIMPLE EXAMPLE

**HOW TO COUNT NUMBER OF TIMES WORDS OCCUR IN A DOCUMENT?**
(IF THAT DOCUMENT IS SPREAD ACROSS MANY MACHINES)

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

I: 3
am: 3
Sam: 3
do: 1
you: 1
like: 1
…

# JUST A HASH TABLE

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

{}

# JUST A HASH TABLE

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

{I:1}

# JUST A HASH TABLE

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

{I:1,
am:1}

# JUST A HASH TABLE

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

{I:1,
am:1,
Sam:1}

# JUST A HASH TABLE

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?"

{I:2,
am:1,
Sam:1}

# BUT YOU SAID THE DOCUMENT IS REALLY BIG?

# COMPUTE IN PARALLEL

"I am Sam
I am Sam
Sam I am

Do you like
Green eggs and ham?
I do not like them

Sam I am
I do not like
Green eggs and ham

Would you like them
Here or there?
…"

{I: 3,
am: 3,
Sam: 3

{Do: 2,
… }

{Sam: 1,
… }

{Would: 1,
… }

{I: 6,
am: 4,
Sam: 4,
do: 3
… }

# COMPUTE IN PARALLEL

"I am Sam
I am Sam
Sam I am
Do you like
Green eggs and ham?
I do not like them
Sam I am
I do not like
Green eggs and ham
Would you like them
Here or there?
…"

{I: 3,
am: 3,
Sam: 3

{Do: 2,
… }

{Sam: 1,
… }

{Would: 1,
… }

{I: 4,
am: 3,
… }

{I: 2,
do: 1,
… }

{I: 6,
am: 3,
… }

# COMPUTE IN PARALLEL

"I am Sam
I am Sam
Sam I am

Do you like
Green eggs and ham?
I do not like them

Sam I am
I do not like
Green eggs and ham

Would you like them
Here or there?
…"

{I: 3,
am: 3,
Sam: 3

{Do: 2,
… }

{Sam: 1,
… }

{Would: 1,
… }

{I: 6,
do: 3, …}

{am: 5,
Sam: 4, …}

{you: 2
…}

{Would: 1
…}

MAP | REDUCE

"I am Sam
I am Sam
Sam I am
--------------------------------
Do you like
Green eggs and ham?
I do not like them
--------------------------------
Sam I am
I do not like
Green eggs and ham
--------------------------------
Would you like them
Here or there?
…"

{I: 3,
am: 3,
Sam: 3

{Do: 2,
… }

{Sam: 1,
… }

{Would: 1,
… }

{I: 6,
do: 3, …}

{am: 5,
Sam: 4, …}

{you: 2
…}

{Would: 1
…}

# MAP-REDUCE

SPLIT DATA & SEND TO MULTIPLE MACHINES (IF NOT ALREADY THERE)

**MAP** FILTER, SORT, AND PROCESS DATA LOCALLY

**REDUCE** CONSOLIDATE AND SUMMARIZE

# MAP-REDUCE

## CAN BE SHORT, SELF-CONTAINED FUNCTIONS
(HERE AS PYTHON-ESQUE PSEUDO CODE)

**MAP**

```
function Map(Document document):
    for each Word w in document:
        EmitIntermediate(w, 1)
```

**REDUCE**

```
function Reduce(Word w, Iterator intermediates):
    int count= 0
    for each int value in intermediates:
        count += value
    Emit(w, count)
```

# MAP-REDUCE

BIG INSIGHT ISN'T
MAP / REDUCE METHODS,
BUT THEIR **SIMPLICITY**
AND THE **ARCHITECTURE
AROUND THEM**

PROVIDES **SCALABILITY**
AND **FAULT-TOLERANCE**
FOR BIG DATA
PROCESSING JOBS

# DEALING WITH ERRORS

## SERVER FAILURE

1 server fails every 3 years

➔ 10K nodes see 10 faults/day

## STRAGGLERS

Nodes are slow or unresponsive

# JUST LAUNCH A REPLACEMENT



"I am Sam
I am Sam
Sam I am

{I: 3,
am: 3,
Sam: 3

Do you like
Green eggs and ham?
I do not like them

{Do: 2,
… }

Sam I am
I do not like
Green eggs and ham

{Sam: 1,
… }

Would you like them
Here or there?
…"

{Would: 1,
… }

# APACHE HADOOP

**OPEN-SOURCE** DISTRIBUTED FILE SYSTEM + MAP REDUCE **AND MORE**

**INSPIRED** BY GOOGLE'S SYSTEMS

**MANY DATA PROCESSING** PIPELINES NOW BUILT ON HADOOP INFRASTRUCTURE

PIG (DATA FLOW LANGUAGE)

HIVE (DATA WARE-HOUSING)

SPARK (IN-MEMORY, MACHINE LEARNING, ETC.)

AND MANY, MANY, MORE...

**MAP REDUCE**

**HDFS** (DISTRIBUTED FILE SYSTEM)

# SOME OPTIONS FOR SPECIFYING BIG DATA PROCESSING OPERATIONS

**WRITE YOUR OWN** MAP-REDUCE METHODS

USE A QUERY LANGUAGE LIKE **APACHE PIG** THAT CAN COMPILE DOWN TO MAP REDUCE-STYLE DISTRIBUTED COMPUTATIONS

```
a = load '/documents';
b = foreach a generate flatten(TOKENIZE((chararray)$0)) as word;
c = group b by word;
d = foreach c generate COUNT(b), group;
store d into '/pig_wordcount';
```

# BENEFITS AND CHALLENGES

Data manipulation on clusters
is now a **big business.**

There is a **huge library of tools** for querying
and processing distributed data.

**BUT...** Most of these tools are **not**
real-time or interactive.

WHAT IF YOU NEED TO <u>INTERACTIVELY</u> EXAMINE OR VISUALIZE A BIG DATASET?

# DATA ANALYSIS AT SCALE

CHALLENGES

ANALYSIS AND CLUSTER COMPUTING

INTERACTING WITH BIG DATA

PARALLELIZING HUMAN INTELLIGENCE

# **STRATEGIES** FOR PROVIDING INTERACTIVITY WITH BIG DATA

**1. INTERACTIVITY VIA PRECOMPUTATION**
**(AGGREGATE AND <u>THEN</u> INTERACT)**

**2. VISUALIZATION AS QUERY SPECIFICATION**
**(LEAVE BIG DATA ON THE SERVERS)**

**3. SAMPLE INTERACTIVELY**
**(APPROXIMATE FIRST THEN REFINE)**

# STRATEGIES FOR PROVIDING INTERACTIVITY WITH BIG DATA

**PARALLELIZE**

1. INTERACTIVITY VIA PRECOMPUTATION
(AGGREGATE AND <u>THEN</u> INTERACT)

2. VISUALIZATION AS QUERY SPECIFICATION
(LEAVE BIG DATA ON THE SERVERS)

3. SAMPLE INTERACTIVELY
(APPROXIMATE FIRST THEN REFINE)

# SAMPLING FOR INTERACTION

# STANDARD QUERY

BIG
DISTRIBUTED
DATABASE

ALL RESULTS
RETURNED AT
THE END

RESULTS

# INTERACTIVE SAMPLING

**Q?**

QUERY

RESULTS
RETURNED
WHILE QUERY
IS STILL RUNNING

BIG
DISTRIBUTED
DATABASE

# INTERACTIVE SAMPLING



http://control.cs.berkeley.edu/papers.html

**CONTROL [HELLERSTEIN ET AL. 1999]**

# INTERACTIVE SAMPLING



**SAMPLEACTION [ FISHER ET AL. 2012 ]**

# INTERACTIVE SAMPLING

**BUT...**

MOST BACKENDS AREN'T DESIGNED TO
RETURN PROGRESSIVE RESULTS

WE NEED GOOD SAMPLING DISTRIBUTIONS FOR EACH
FIELD TO PRODUCE MEANINGFUL INTERMEDIATE RESULTS

HOW BEST TO VISUALIZE UNCERTAINTY?

HOW WELL CAN PEOPLE INTERPRET PARTIAL RESULTS?

**THIS IS STILL A <u>VERY</u> OPEN RESEARCH AREA!**

# HOW TO SHOW UNCERTAINTY?



[Olston & Mackinlay, 2002]

Figure 1: Error bars and ambiguation applied to some common chart types.

# HOW TO SHOW UNCERTAINTY?



[Streit, Pham, & Brown 2008]

# HOW TO SHOW UNCERTAINTY?



ERROR BARS CONSISTENTLY UNDERPERFORMED

- High uncertainty
- Low uncertainty

a.

[Sanyal, et al. 2009]

- High uncertainty
- Low uncertainty

c.

High uncertainty
Low uncertainty

d.

# HOW TO SHOW UNCERTAINTY?



[Boukhelifa, et al. 2012]

PEOPLE DON'T ALWAYS INTERPRET THESE AS SHOWING UNCERTAINTY

# A FEW INTERESTING RESEARCH PROTOTYPES

STILL USING SIMPLE VISUALS

NO UNCERTAINTY INFO

NOTE: ANALYSIS NOTEBOOKS AND PROVENANCE

**TEMPE [Microsoft Research 2014]**

# DATA ANALYSIS AT SCALE

CHALLENGES

ANALYSIS AND CLUSTER COMPUTING

INTERACTING WITH BIG DATA

PARALLELIZING HUMAN INTELLIGENCE

# HOW CAN WE LEVERAGE MULTIPLE PEOPLE TO EXPEDITE ANALYSIS?

Analyst

CollegeRankings2013.csv

Analyst

Crowd

MANY IMPORTANT ANALYSIS TASKS REQUIRE **HUMAN INTELLIGENCE** BUT LEND THEMSELVES WELL TO **PARALLELIZATION**

# MANY IMPORTANT ANALYSIS TASKS REQUIRE HUMAN INTELLIGENCE BUT LEND THEMSELVES WELL TO PARALLELIZATION

Sensemaking Loop

Foraging Loop

[Pirolli & Card 2005]

# MANY EYES



[Viégas, et al. 2007, 2008]

# GOOGLE BOOKS N-GRAMS

# CROWDSOURCING DATA ANALYSIS

DATA COLLECTION & CITIZEN SCIENCE

ANALYSIS COMPETITIONS

"MICROWORK" AND TASK MARKETS

COLLABORATION TOOLS FOR ANALYSTS

# CITIZEN SCIENCE

## DATA COLLECTION

## CREEK WATCH
[IBM]

4,000 Creek Watch users

in over 25 countries

## CHRISTMAS BIRD COUNT

CITIZEN SCIENCE

DATA PROCESSING

SETI@home
The Search for Extraterrestrial Intelligence

SETI@Home

Time (sec)

Frequency (Hz)

Power

Resolution    0.596  Hz

Base frequency:   1419804000

User  info

Name:
Team:
Total credit:   0.00

# CITIZEN SCIENCE

## HUMAN VISION & PROBLEM SOLVING

## FOLD.IT

# ANALYSIS COMPETITIONS

## KAGGLE



## NETFLIX PRIZE

# MICROWORK PLATFORMS

SITES WHERE **WORKERS** PERFORM SMALL PIECES OF WORK ("**TASKS**") - USUALLY IN EXCHANGE FOR **SMALL FINANCIAL REWARDS.**

**amazon mechanical turk**™
Artificial Artificial Intelligence

**CrowdFlower**     **mobileworks**

# MICROWORK

## USING APIS – DEVELOPERS CAN WRITE PROGRAMS THAT INCORPORATE HUMAN JUDGEMENT

"HUMAN COMPUTATION"

# APPLYING MICROWORK TO DATA ANALYSIS

## CROWDSOURCING LOW-LEVEL ANALYSIS

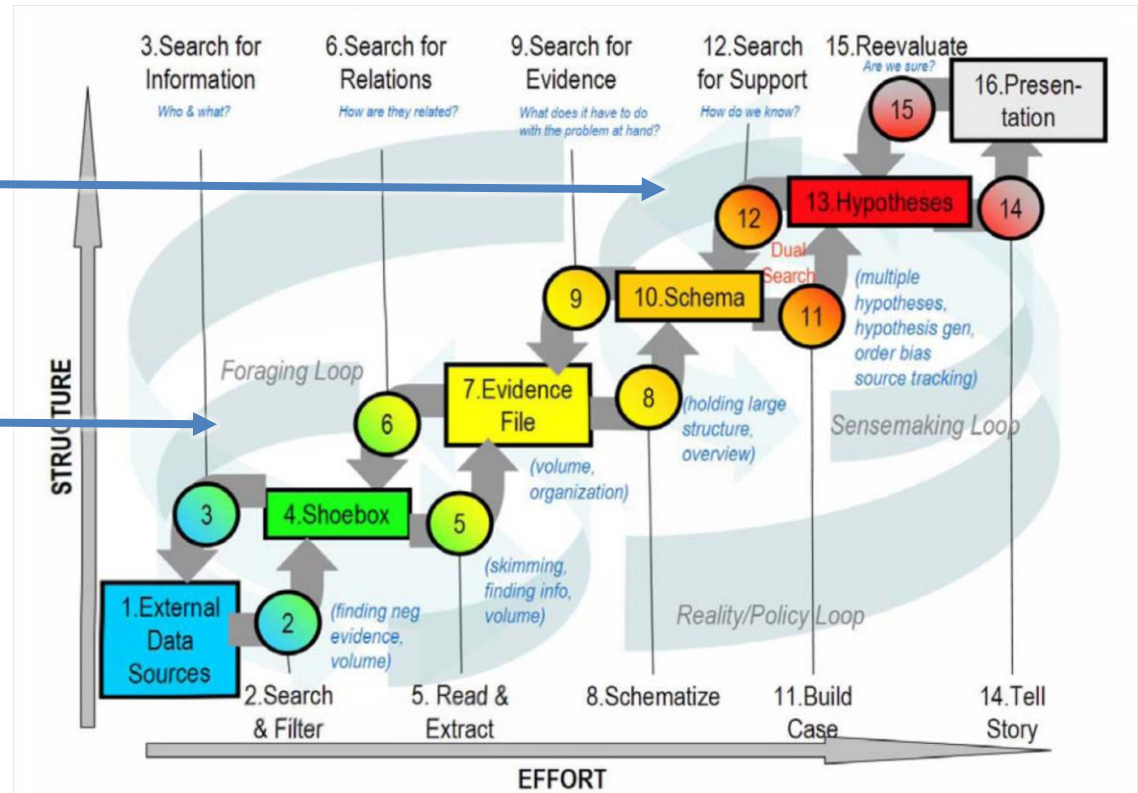DATA COLLECTION AND DATA ENTRY

LABELING

DATA CLEANING

SENTIMENT ANALYSIS

# MANY IMPORTANT ANALYSIS TASKS REQUIRE HUMAN INTELLIGENCE BUT LEND THEMSELVES WELL TO PARALLELIZATION
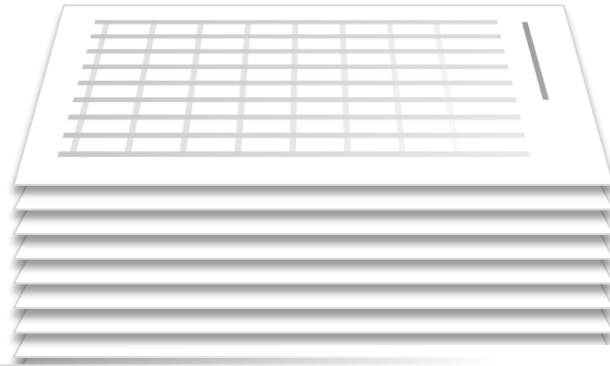
Sensemaking Loop

Foraging Loop

[Pirolli & Card 2005]

# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS



Data    Analyst    Crowd

[Willett et al. CHI 2012, VAST 2013]

# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS

Data

Analyst

Crowd

[Willett et al. CHI 2012, VAST 2013]

# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS

Data

Analyst

Crowd

[Willett et al. CHI 2012, VAST 2013]

# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS



Select Charts with Trends & Patterns

Data

Analyst

Crowd

[Willett et al. CHI 2012, VAST 2013]

proxy.commentspace.net/explainTask?studyName=Oil–Demo&assignmentId=TestAI&workerId=Exper... ☆ » ⚒

Each of the charts in this HIT shows the **average amount of oil produced per day** by one or more countries over the past 50 years



This chart shows **Oil Produced (Thsnds. of Barrels/Day)** by **Year**. The view is filtered by **Country** to show only **"Iran"**.

1.  Explain **why** the strong *peak or valley* highlighted in the chart might have occurred.
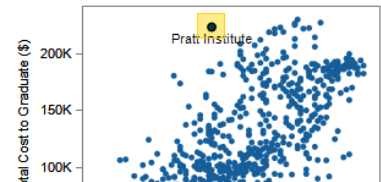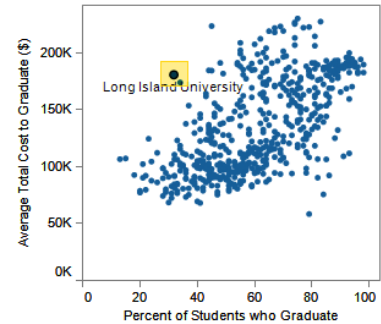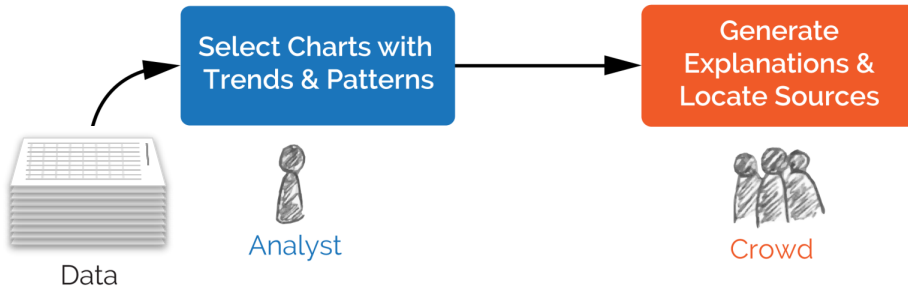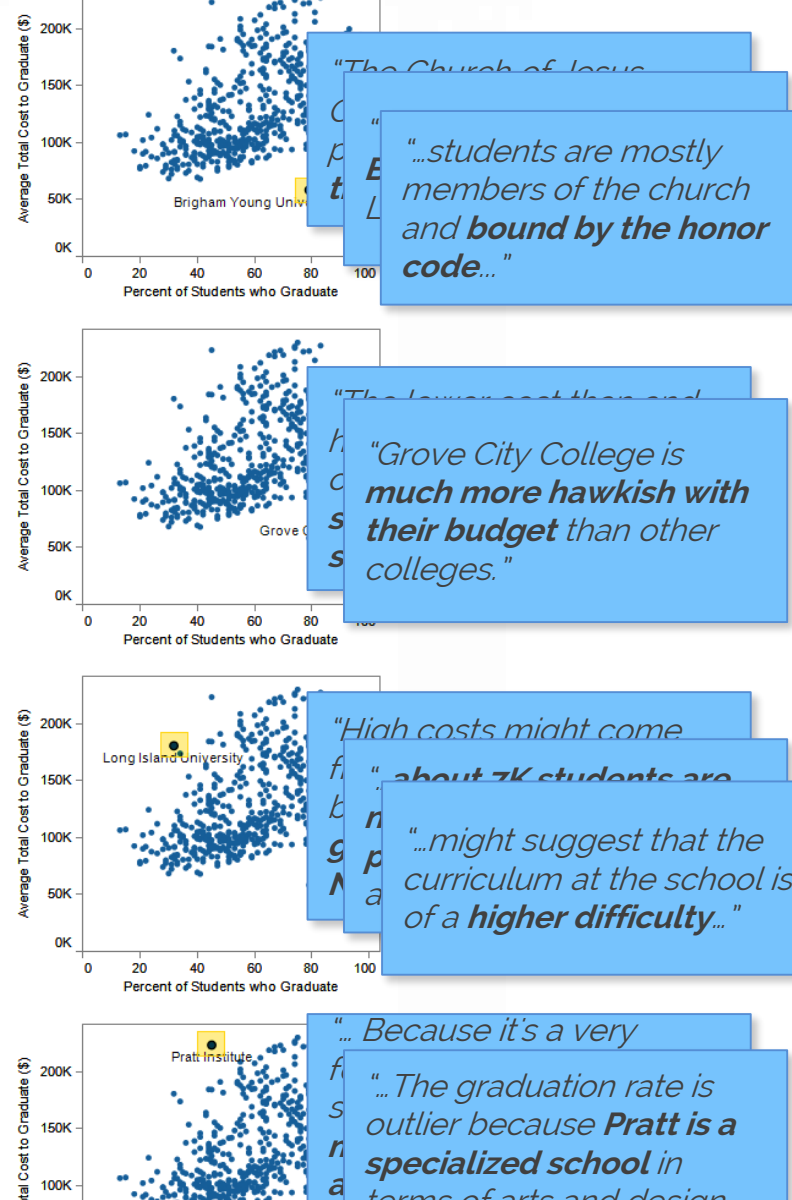
Submit Task

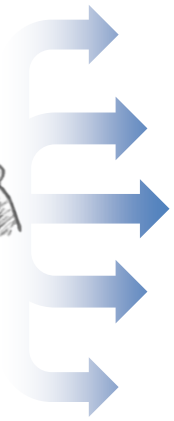# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS



[Willett et al. CHI 2012, VAST 2013]

# "COULD THIS CREATE MORE WORK FOR THE ANALYST?"
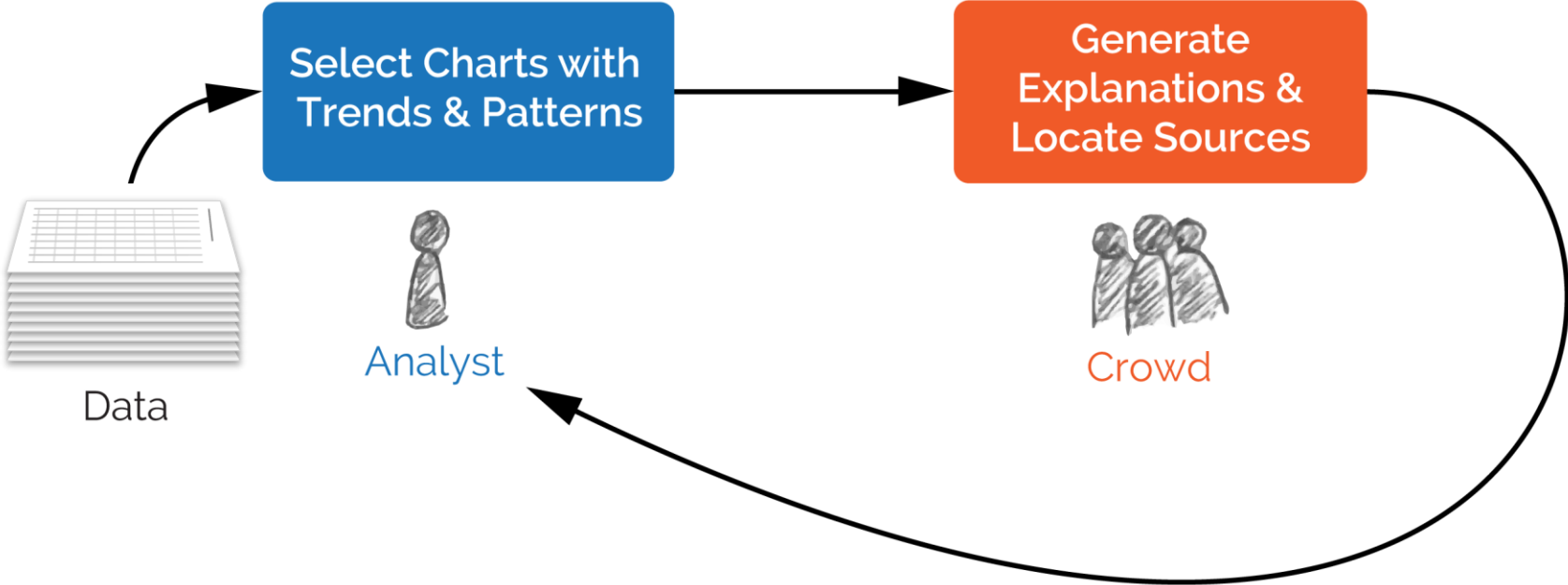
# "COULD THIS CREATE MORE WORK FOR THE ANALYST?"

"High costs might come f... b... g... N..."

"...about 7K students are n... p... a..."

"...might suggest that the curriculum at the school is of a **higher difficulty**..."

"The lower cost than and hig... co... sm... sc..."

"Grove City College is **much more hawkish with their budget** than other colleges."
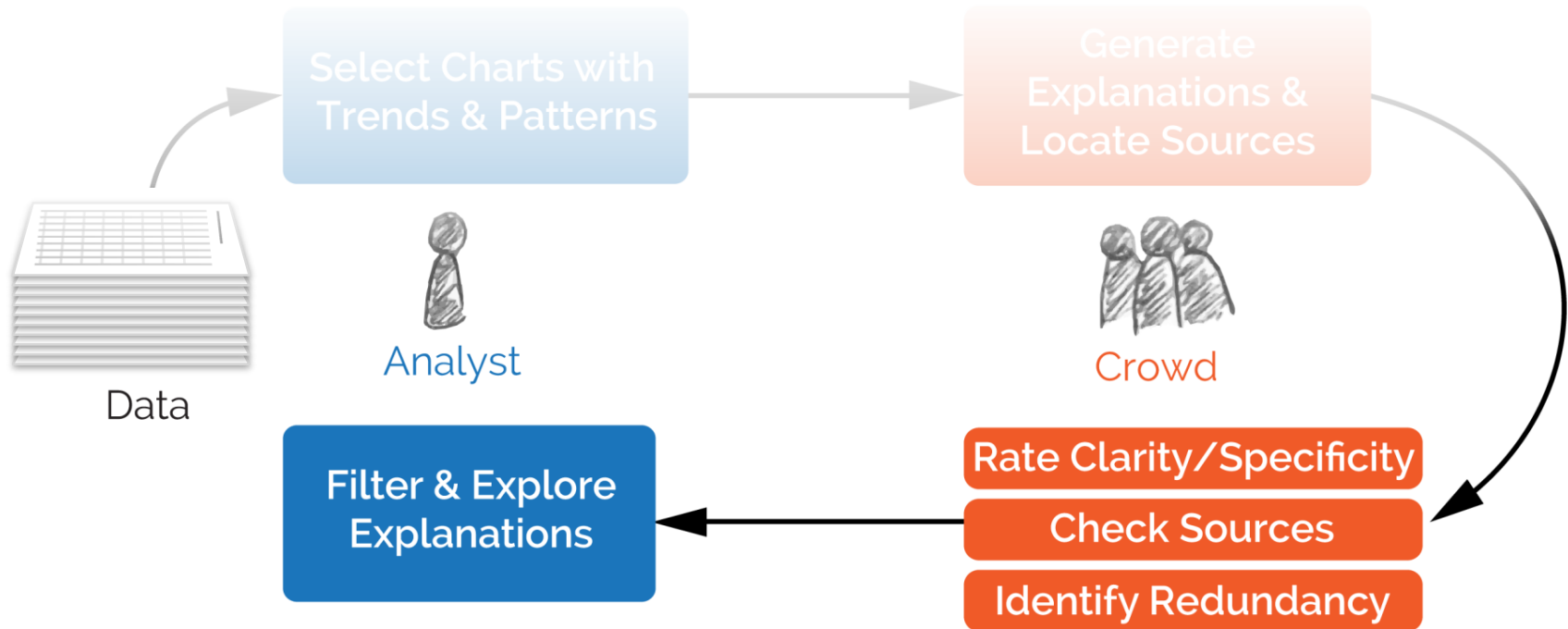
"The Church of Jesus C... p... **E... t...** L..."

"...students are mostly members of the church and **bound by the honor code**..."

"... Because it's a very f... s... **r... a... t...**"

"...The graduation rate is outlier because **Pratt is a specialized school** in terms of arts and design and students..."

# A WORKFLOW FOR CROWDSOURCING DATA ANALYSIS

# CROWD-ENABLED EXTENSIONS FOR PROCESSING AND MANAGING RESULTS

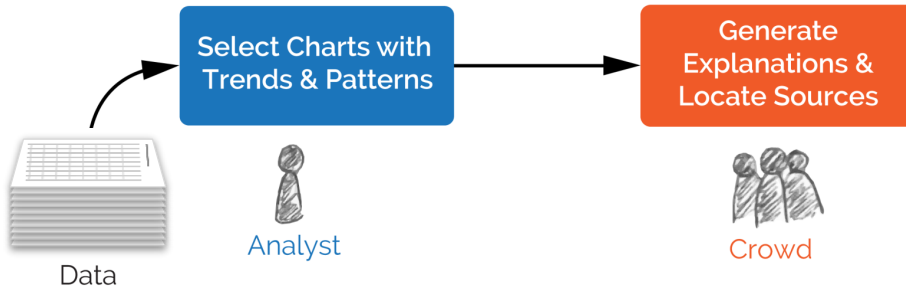# THREE CRITERIA FOR PLAUSIBLE EXPLANATIONS

## CLARITY AND SPECIFICITY

## PROVENANCE

## REDUNDANCY

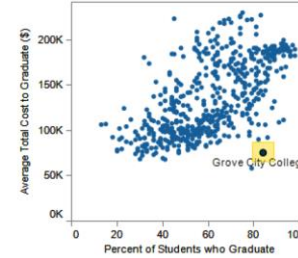+ AN INTERFACE FOR MANAGING CROWDSOURCED EXPLANATIONS

**CLARITY & SPECIFICITY**

# CLARITY AND SPECIFICITY

Data

Analyst

Crowd

Select Charts with Trends & Patterns

Generate Explanations & Locate Sources
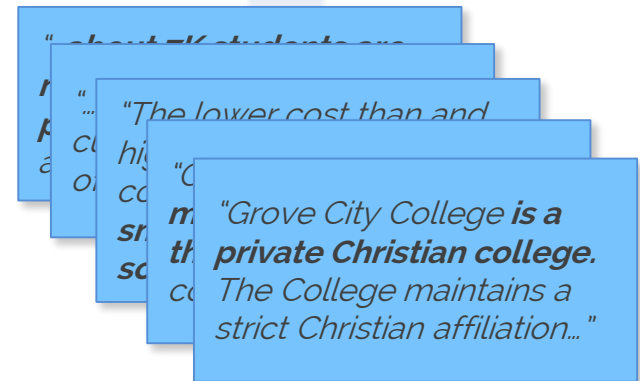
Show Instructions

Each of the charts in this hit compares the graduation rate (x-axis) and the total cost to graduate (y-axis) for 554 top US colleges and universities (as ranked by Bloomberg Businessweek in 2010). Each point represents a single college or university.
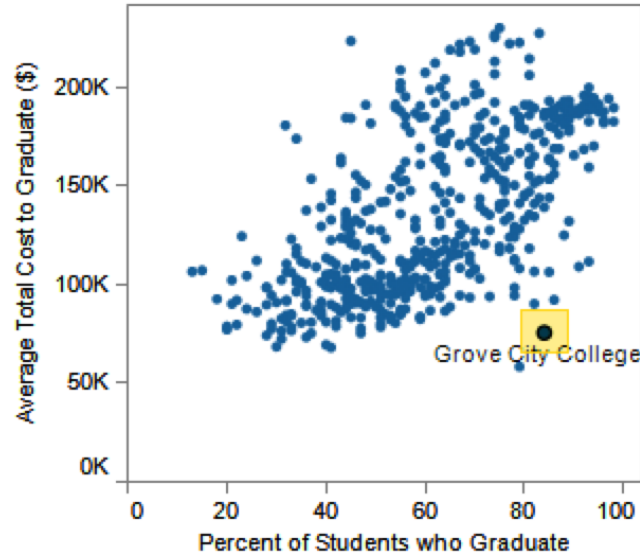
**Prompt:** Explain **why** the **outlier** highlighted in the chart might be different from the other items. (Give **one** specific, well-justified answer.)

**Response R2:** " *Grove City College is a private Christian college. The College maintains a strict Christian affiliation, in contrast to many institutions whose religions affiliations have become merely historical in nature. This Christian identity, as well as a heavily politically Conservative identity, on campus may likely attract superior students who would not choose to attend otherwise comparable institutions lacking this culture.*"(Reference: http://www.discoverthenetworks.org/Articles/Conservative%20Colleges.htm )

1. Does this response provide an explanation for **why** the highlighted outlier in the chart might have occurred?    ◯ **Yes**  ◯ **No**  ◯ **None Present**

2. How **clear** and **specific** is the response?    (Not Clear/Specific) ← ◯ 1 ◯ 2 ◯ 3 ◯ 4 ◯ 5 → (Very Clear/Specific)

"Grove City College **is a private Christian college.** The College maintains a strict Christian affiliation..."

Each of the charts in this hit compares the graduation rate (x-axis) and the total cost to graduate (y-axis) for 554 top US colleges and universities (as ranked by Bloomberg Businessweek in 2010). Each point represents a single college or university.



**Prompt:** Explain **why** the *outlier* highlighted in the chart might be different from the other items. (Give **one** specific, well-justified answer.)

**Response R2:** " *Grove City College is a private Christian college. The College maintains a strict Christian affiliation, in contrast to many institutions whose religions affiliations have become merely historical in nature. This Christian identity, as well as a heavily politically Conservative identity, on campus may likely attract superior students who would not choose to attend otherwise comparable institutions lacking this culture.*"(Reference: http://www.discoverthenetworks.org/Articles/Conservative%20Colleges.htm )

1. Does this response provide an explanation for **why** the highlighted outlier in the chart might have occurred?
        ○ **Yes**    ○ **No**    ○ **None Present**

2. How **clear** and **specific** is the response?     (Not Clear/Specific) ← ○ 1 ○ 2 ○ 3 ○ 4 ○ 5 → (Very Clear/Specific)
Clear/Specific)

PROVENANCE

# PROVENANCE



Explanation Task

What are our



workers doing?

# PROVENANCE



Explanation Task

# INSTRUMENTING EXPLANATION TASKS



Examine a line chart showing employment change in a US city and briefly explain it.
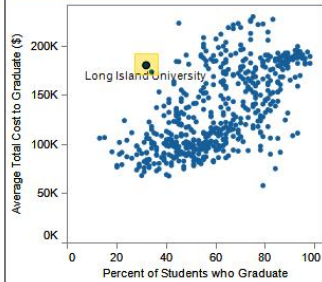
Requester: visualizationlab.ucb
Qualifications Required: Location is US

Reward: $0.40 per HIT    HITs Available: 10    Duration: 30 minutes

Each of the charts in this hit compares the graduation rate (x-axis) and the total cost to graduate (y-axis) for 554 top US colleges and universities (as ranked by Bloomberg Businessweek in 2010). Each point represents a single college or university.

Average Total Cost to Graduate ($)

200K
150K
100K
50K
0K

Long Island University

0    20    40    60    80    100
Percent of Students who Graduate

1. What college or university is highlighted in this chart?

2. Explain **why** the strong *outlier* highlighted in the chart might be different from the other items. (Try to give **one** specific, well-justified answer per text box.)

   If there are multiple explanations, enter each one in a separate text box.

   Using the browser to the right, find text on a web page that justifies each explanation. Select the text and click the "mark as source" button to add it.
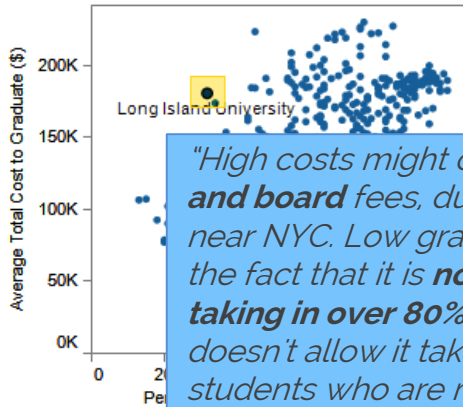
**Explanation 1**

Source:

+ Add Another Explanation    --

Finished with this HIT?    Let someone else do it?
Submit HIT    Return HIT

# PROVENANCE

## Paragraph-level citations

## Visitation logs

"High costs might come from it's **high room and board** fees, due to its geographic location near NYC. Low graduation rates come from the fact that it is **not a very selective school, taking in over 80% of applicants**, which doesn't allow it take many top ranked students who are more academically motivated."

PROVENANCE

DID THE **FACTS** AND **INFERENCE** COME FROM THE **SOURCE** OR DID THE **WORKER** ADD THEM?
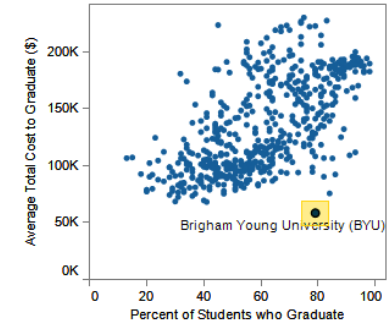
# SOURCE-CHECKING MICROTASKS



A second group of workers verifies links and attributes explanations to the source or the worker.( 75% accurate in our preliminary tests )

REDUNDANCY

# REDUNDANCY

Many explanations provided by workers are redundant.



"The Church of Jesus Christ of Latter Day Saints pays a significant part of the tuition costs..."

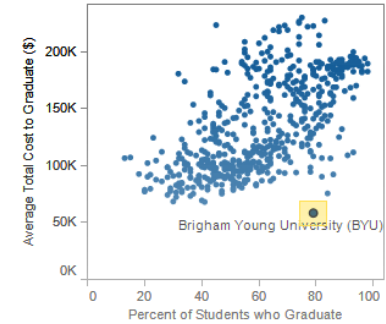"The cost of attendance at BYU is subsidized by the LDS church."

"98% of their students are members of LDS and they have lowered tuition..."

# REDUNDANCY

Many explanations provided
by workers are redundant.



— Duplicate results for analysts to examine.

"The Church of Jesus
Christ of Latter Day Saints
pays a significant part of

✚ Redundancy can signal high support
and corroborating sources.

"The cost of attendance at
BYU is subsidized by the
LDS church."

"98% of their students are
members of LDS and they
have lowered tuition

# REDUNDANCY



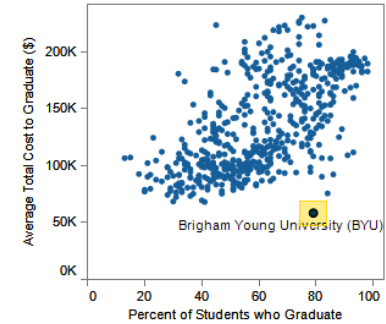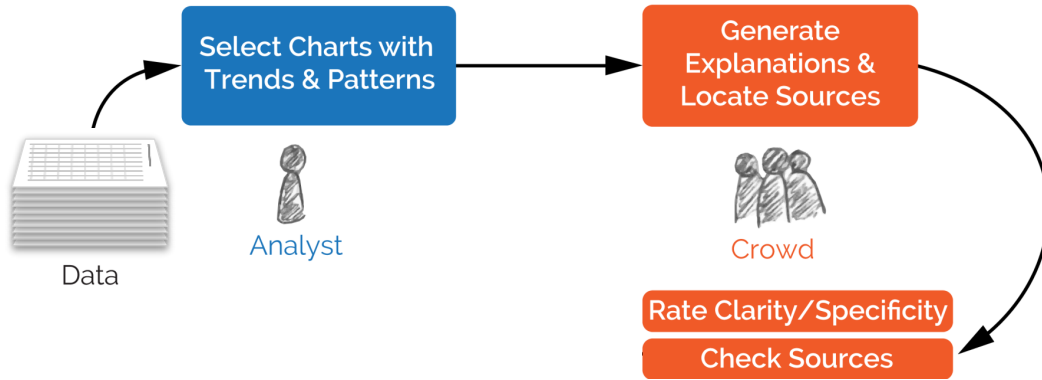Automated text similarity methods don't deal well with these kinds of content.

"The Church of Jesus Christ of Latter Day Saints pays a significant part of the tuition costs…"

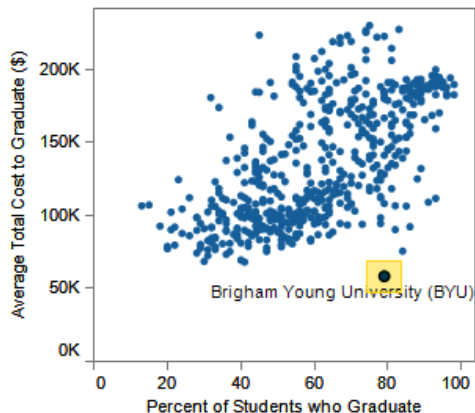"The cost of attendance at BYU is subsidized by the LDS church."

"98% of their students are members of LDS and they have lowered tuition…"

# REDUNDANCY



**Select Charts with Trends & Patterns**

Analyst

Data

**Generate Explanations & Locate Sources**

Crowd

**Rate Clarity/Specificity**

**Check Sources**

Can we crowdsource redundancy detection?

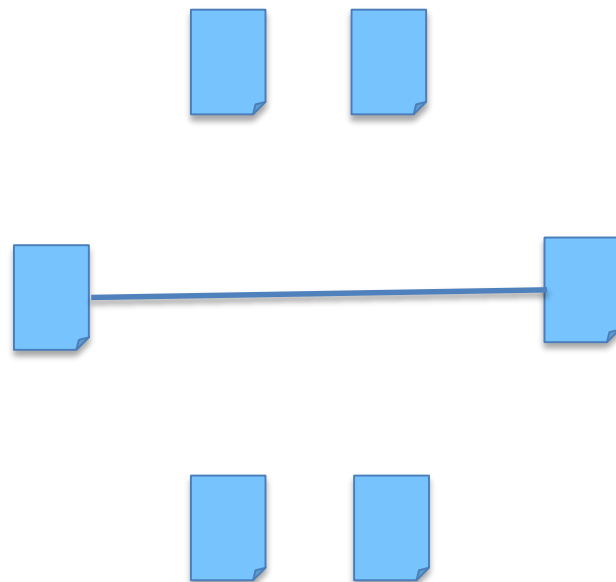# CLUSTERING VIA DISTRIBUTED COMPARISON



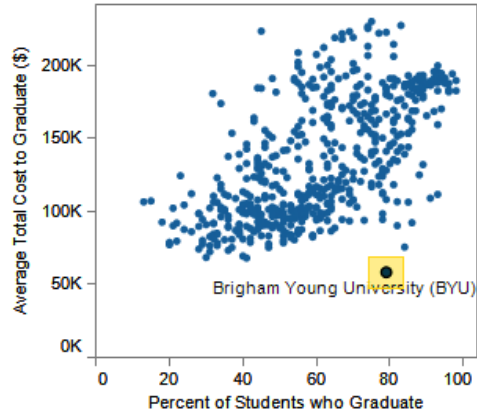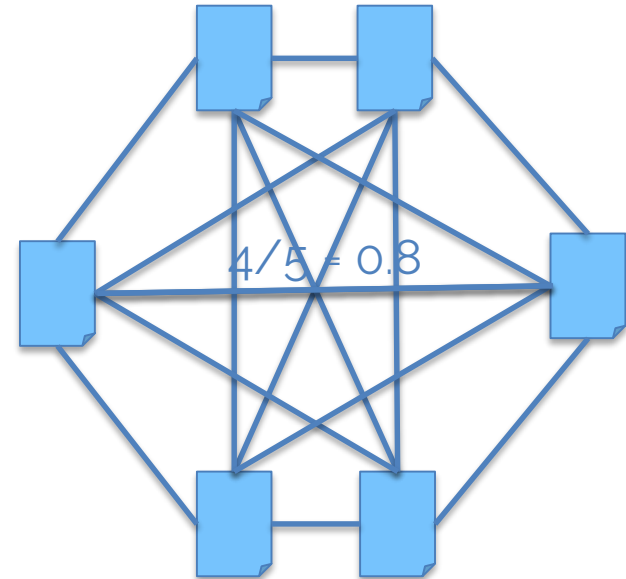"98% of their students are members of LDS and they have lowered tuition..."

"The cost of attendance at BYU is subsidized by the LDS church."

"...students are mostly members of the church and bound by the honor code..."

"The Church of Jesus Christ of Latter Day

# CLUSTERING VIA DISTRIBUTED COMPARISON



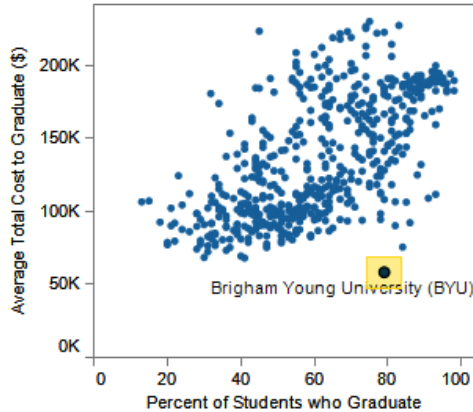"98% of their students are members of LDS and they have lowered tuition..."

"The cost of attendance at BYU is subsidized by the LDS church."

Do these two responses give the same general explanation for the peaks and valleys in the chart?
○ **Yes. Both responses give the same general explanation.**
○ **No. The responses do not give the same explanation.**

# CLUSTERING VIA DISTRIBUTED COMPARISON



"98% of their students are members of LDS and they have lowered tuition..."

"The cost of attendance at BYU is subsidized by the LDS church."

Do these two responses give the same general explanation for the peaks and valleys in the chart?
○ Yes. Both responses give the same general explanation.
○ No. The responses do not give the same explanation.

**+** Simple tasks for workers

Scales poorly
Sensitive to clustering method
Workers have little context

# CLUSTERING VIA COLOR-CODING

**Prompt:** Explain **why** the strong **peak or valley** highlighted in the chart might have occurred.

**Response R2:** "A new medical school is providing jobs"(Reference: newspapertree.com/opinion/3561-the-el-paso-stimulus )

**Response R7:** "The Medical Center of the Americas opened a new medical school and in 2008 construction on a new series of projects began at the University of Texas El Paso. "(Reference: newspapertree.com/opinion/3561-the-el-paso-stimulus )

**Response R3:** "Expansion of Fort Bliss"(Reference: newspapertree.com/opinion/3561-the-el-paso-stimulus )

**Response R1:** "Increase of construction jobs."(Reference: newspapertree.com/opinion/3561-the-el-paso-stimulus )

**Response R4:** "It would appear that the marked growth in jobs up until 2008 coincides with growth of businesses in the area. Notable amongst these businesses are the three school districts that service the city and growth in the health services industry."(Reference: www.google.com/search?q=el paso employers 2007 )

**Response R5:** "The high peak in 2008 was during the time when the economy was overheated. After that time the economy slipped into a recession which caused the employment status of many people to change. This is why after 2008 the graph shows a sharp drop in employment. " (Reference: www.google.com/url?q=http://en.wikipedia.org/wiki/Late-2000s_recession&sa=U&ei=ae5qT6yoBMaosQKGIOCWCA&ved=0CBQQFjAB&usq=AFQjCNGuzI5xk-iIEUTtOjK4C8Gi6DP0FQ )

MULTIPLE WORKERS **INDEPENDENTLY CLUSTER** THE WHOLE SET.

USE **COMPUTATIONAL SIMILARITY METRICS** TO SELECT THE BEST, CONSISTENT CLUSTERING.

**FINDING THE RIGHT BALANCE OF HUMAN AND AUTOMATED EFFORT**

# MANAGING THE CROWD'S WORK

# MANAGING THE CROWD'S WORK

# EXPLANATION MANAGEMENT INTERFACE

HUMAN COMPUTATION CAN BE A USEFUL COMPLEMENT TO AUTOMATED PROCESSING

EVEN MORE INTERESTING WITH EXPERTISE



**cheap low-skill crowds**
VS.
**more knowledgeable trusted ones**

UNDERSTANDING HOW TO PARALLELIZE **ANALYSIS PROCESSES** MAY BE AS IMPORTANT AS PARALLELIZING COMPUTATION HAS BEEN.

# DATA ANALYSIS AT SCALE

CHALLENGES

ANALYSIS AND CLUSTER COMPUTING

INTERACTING WITH BIG DATA

PARALLELIZING HUMAN INTELLIGENCE

# UP NEXT

**AFTER THE BREAK**
APPLICATION AREAS

**THIS AFTERNOON**
DINO FUN WORLD PRESENTATIONS
(OPEN LAB)

**DECEMBER 8th-19th**
INFORMATION VISUALIZATION LECTURES
AT UNIVERSITÉ PARIS SUD

# BONUS MATERIAL

MORE DETAILS ON CROWDSOURCED DATA ANALYSIS

# CLUSTERING VIA COLOR-CODING



Individual workers cluster the whole set.

**+** Workers have **complete context**

**—** **Individual workers** can cluster badly

**Hard to integrate clusterings** from multiple workers

# HOW TO INTEGRATE COLOR-CLUSTERINGS?



- A single worker's clustering is preferable to a combination of multiple clusterings.

- Clusters reproduced by multiple independent workers are likely to reflect actual redundancy.

- Errors tend to be either noisy or easy to catch.

# HOW TO INTEGRATE COLOR-CLUSTERINGS?



Selecting the Most-Representative Clustering

# HOW TO INTEGRATE COLOR-CLUSTERINGS?

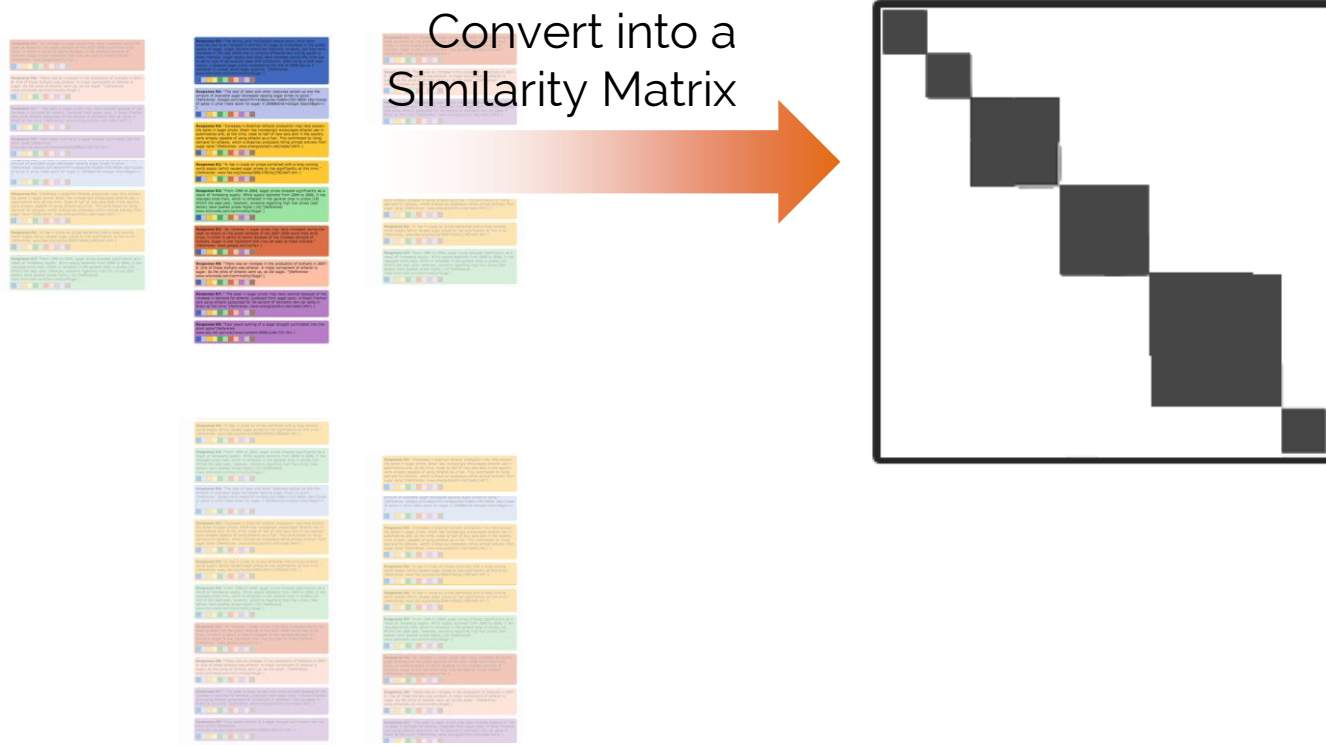# SELECTING THE MOST-REPRESENTATIVE CLUSTERING



Convert into a Similarity Matrix

# SELECTING THE MOST-REPRESENTATIVE CLUSTERING

# SELECTING THE MOST-REPRESENTATIVE CLUSTERING



Average

# SELECTING THE MOST-REPRESENTATIVE CLUSTERING

# EVALUATING REDUNDANCY-DETECTION

Does color clustering with most-representative selection produce good clusterings?

Our Explanation Dataset

12 charts (4 each from 3 different data sets)

10 workers explained each chart

➡ 93 Workers produced 156 explanations (Avg=13 per chart)

# EVALUATING REDUNDANCY-DETECTION

Does color clustering with most-representative selection produce good clusterings?

10 Workers used color clustering to group the explanations for each chart. (120 total clusterings)

We used most-representative selection to pick the best clustering for each chart. (12 clusterings)

# EVALUATING REDUNDANCY-DETECTION

Baseline - Expert clustering ( x 3 )

To score a clustering, we use the F-measure to compute similarity to each expert, then average.

(completely dissimilar) $[0 \longleftrightarrow 1]$ (identical)
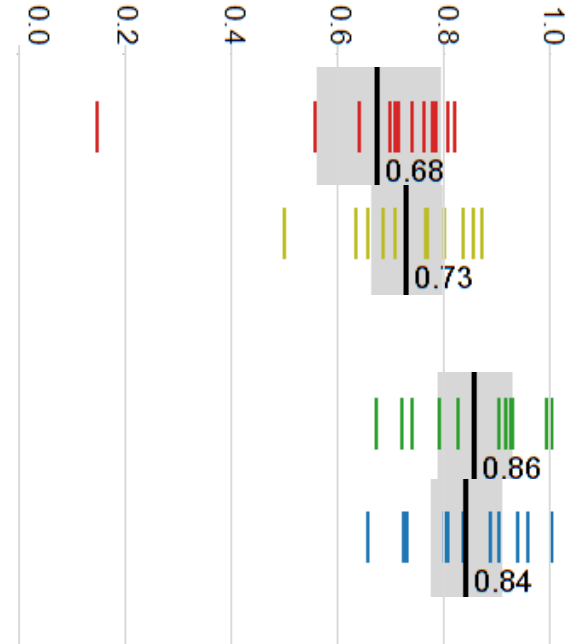
# EVALUATING REDUNDANCY-DETECTION



Average F-measure Score (vs. Experts)

Unclustered Results F=0.68
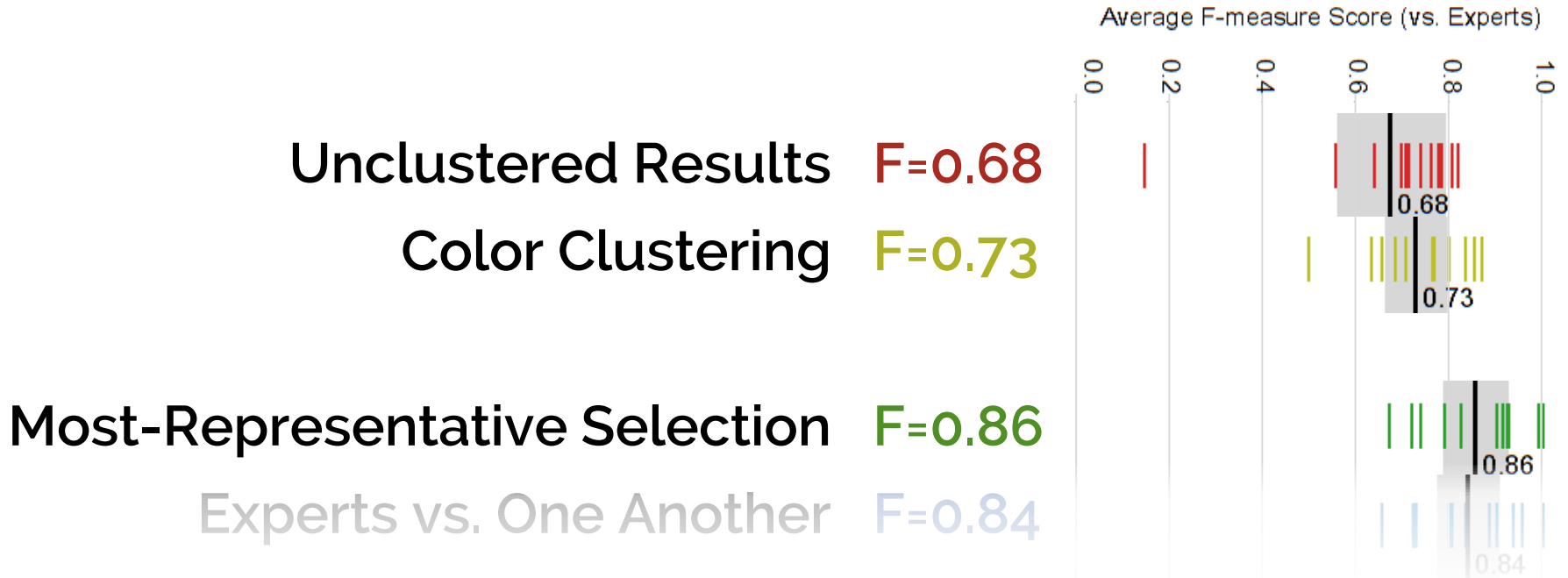
Color Clustering F=0.73

Most-Representative Selection F=0.86

Experts vs. One Another F=0.84

# EVALUATING REDUNDANCY-DETECTION



Average F-measure Score (vs. Experts)

Unclustered Results  F=0.68

Color Clustering  F=0.73

Most-Representative Selection  F=0.86

Experts vs. One Another  F=0.84

T-tests showed our most-representative results were significantly closer to experts than color clustering or unclustered were. (both p < 0.01)