# Human-Powered Database Operations: Part 1

PENN STATE
1855

Dongwon Lee

Penn State University, USA

`dongwon@psu.edu`

Slide available @ **http://goo.gl/4pNUhB**
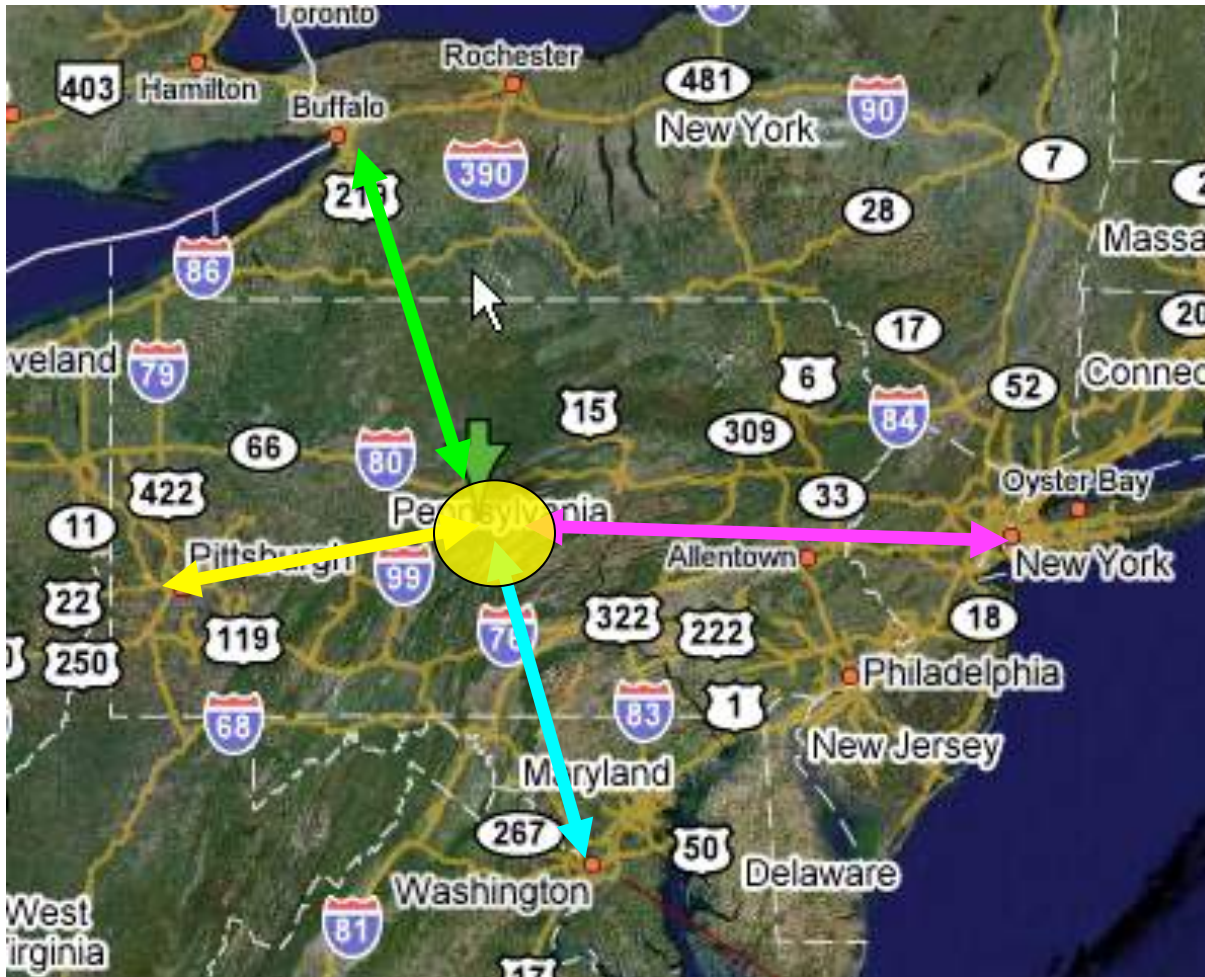
**SBBD 2014 Tutorial**

# Where Am I From?

# Penn State University



- ## State College, PA
  - Out of nowhere, but close to everywhere

- ## West: 2.5 hours to **Pittsburgh**
- ## East: 4 hours to **New York**
- ## South: 3 hours to **Washington DC**
- ## North: 3 hours to **Niagara Fall**

# Penn State *i*-School

- College of Information Sciences and Technology (IST)
    - http://ist.psu.edu/
- 40+ tenure-track faculty on diverse areas
    - CompSci & EE
    - MIS & LIS
    - Design
    - Law
    - Psychology
    - Medical Infomatics

# Other Tutorials on Crowdsourcing

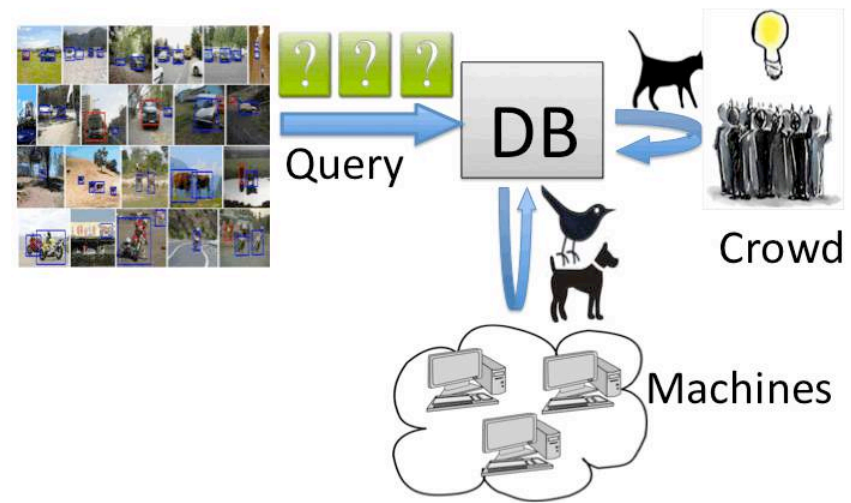| Year | Sub-field | Venue | Title |
|------|-----------|-------|-------|
| 2013 | Crowdsourcing | SDM | Crowdsourcing & Human Computation |
| 2013 | Crowdsourcing | HCOMP | Incentives in Human Computation |
| 2012 | HCI | AAAI | Crowdsourcing using MTurk for HCI Research |
| 2012 | Crowdsourcing | SBP | Crowdsourcing, Human Computation, and Collective Intelligence |
| 2012 | IR | SIGIR | Human Computation and Crowdsourcing |
| 2012 | DB | SIGMOD | Designing a Scalable Crowdsourcing Platform |
| 2011 | Crowdsourcing | AAAI | Human Computation: Core Research Questions and State of the Art |
| 2011 | IR | CLEF | Crowdsourcing for IR Experimentation and Evaluation |
| 2011 | ML | ICML | Collective Intelligence and Machine Learning |
| 2011 | Social Science | EC | Conducting Behavioral Research using AMT |
| 2011 | Multimedia | MM | Frontiers in Multimedia Search |
| 2011 | DB | VLDB | Crowdsourcing Application and Platforms |
| 2011 | Crowdsourcing | WWW | Managing Crowdsourced Human Computation |
| 2010 | Vision | CVPR | Mechanical Turk for Computer Vision |
| 2008 | IR | CIKM | Crowdsourcing for Relevance Evaluation |

# The Focus of This Tutorial

- Part 1 on basics of crowdsourcing

- Part 2 on DB operations that exploit crowdsourcing

CrowdSourcing

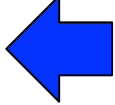When to Ask Whom

Query → DB → Crowd

Machines

http://istc-bigdata.org/index.php/crowdsourcing-big-data/
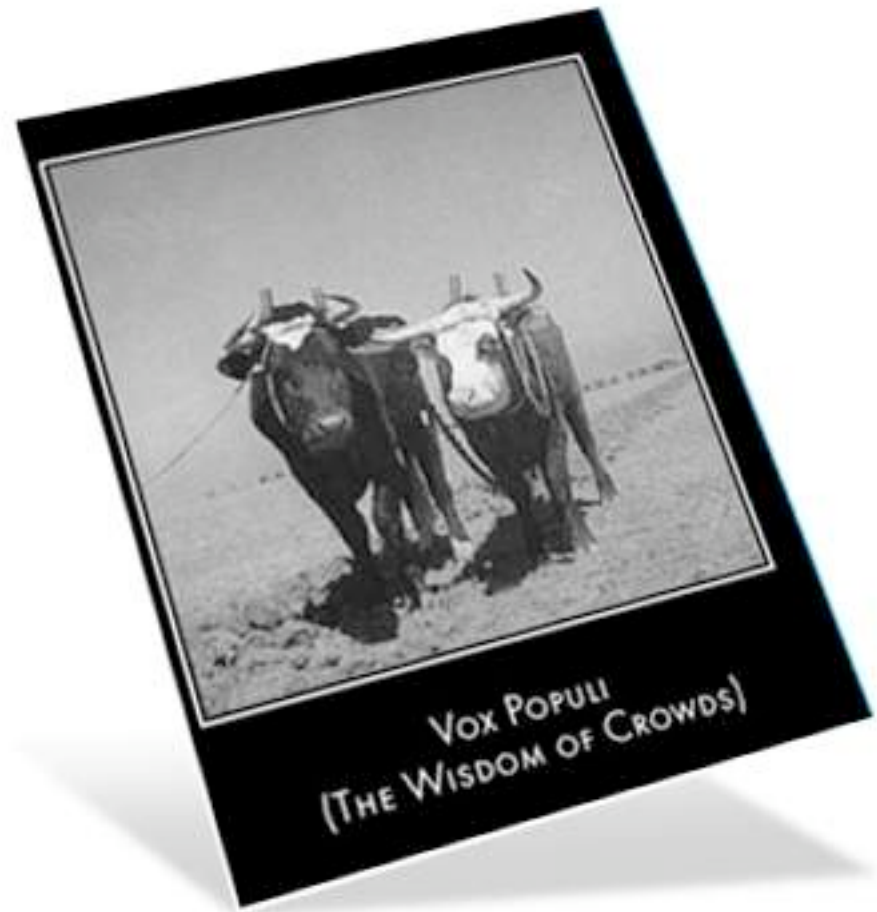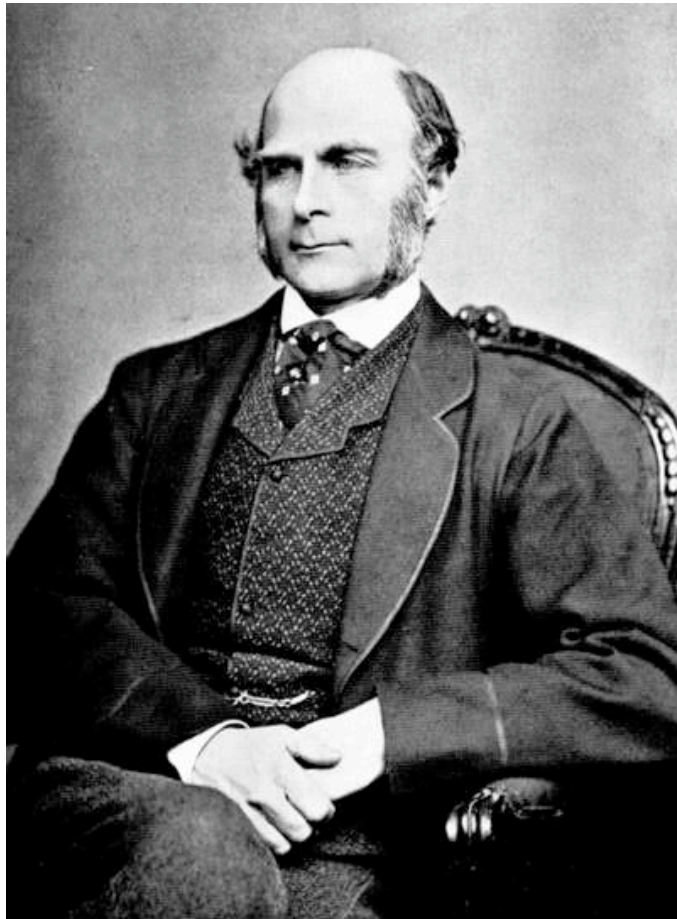
# Part 1: Crowdsourcing Basics

- **Examples** ⬅
- Definitions
- Marketplaces
- Computational Crowdsourcing
  - Preliminaries
  - Transcription
  - Sorting
- Demo

# Eg, Francis Galton, 1906

**Weight-judging competition:**
**1,197 (mean of 787 crowds) vs. 1,198 pounds (actual measurement)**

# Eg, StolenSidekick, 2006

- A woman lost a cellphone in a taxi
- A 16-year-old girl ended up having the phone
  - Refused to return the phone
- Evan Guttman, the woman's friend, sets up a blog site about the incident
  - http://stolensidekick.blogspot.com/
  - http://www.evanwashere.com/StolenSidekick/
  - Attracted a growing amount of attention → the story appeared in Digg main page → NY Times and CNN coverage → Crowds pressure on police …
- NYPD arrested the girl and re-possessed the phone

http://www.nytimes.com/2006/06/21/nyregion/21sidekick.html?_r=0

# Eg, Finding "Jim Gray", 2007

# Eg, Threadless.com

- Sells t-shirts, designed/voted by crowds
- Artists whose designs are chosen get paid

# Eg, KICKSTARTER

- Crowdfunding, started in 2009
- Project creators choose a deadline and a minimum funding goal
  - Creators only from US, UK, and Canada
- Donors pledge money to support projects, in exchange of non-monetary values
  - Eg, t-shirt, thank-u-note, dinner with creators
  - Donors can be from anywhere
- Eg, Pebble, smartwatch
  - 68K people pledged 10M

# Eg, reCAPCHA

The Norwich line steamboat train, from New-London for Boston, this morning ran off the track seven miles north of New-London.

morning

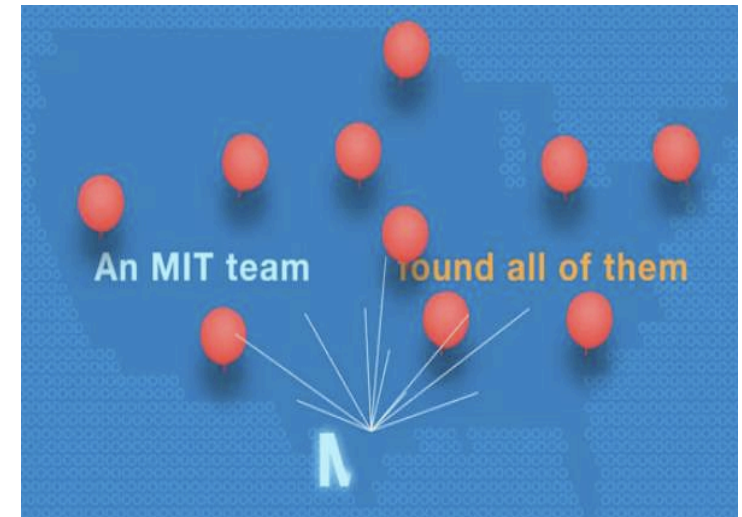morning overlooks

Type the two words:

reCAPTCHA

As of 2012

Captcha: 200M every day

ReCaptcha: 750M to date

# Eg, DARPA Challenge, 2009

- To locate 10 red balloons in arbitrary locations of US

- Winner gets $40K

- MIT team won the race with the strategy:

  - 2K per balloon to the first person, A, to send the correct coordinates

  - 1K to the person, B, who invited A

  - 0.5K to the person, C, who invited B, …



An MIT team found all of them

# Eg, Berkeley Mobile Millennium

# Eg, Who Wants to be a Millionaire?



**Asking the audience** usually works ➔ Audience members have *diverse* knowledge that can be coordinated to provide a correct answer in sum

# Eg, Who Wants to be a Millionaire?

# Eg, Game-With-A-Purpose: GWAP

- **Term coined by Luis von Ahn @ CMU**
- **Eg,**
  - ESP Game → Google Image Labeler: image recognition

  - Foldit: protein folding

  - Duolingo: language translation

# Crowdsourcing landscape Beta v2



http://www.resultsfromcrowds.com/features/crowdsourcing-landscape/

# Part 1: Crowdsourcing Basics

- Examples
- **Definitions** ⬅
- Marketplaces
- Computational Crowdsourcing
  - Preliminaries
  - Transcription
  - Sorting
- Demo

# James Surowiecki, 2004



A NEW YORK TIMES BUSINESS BESTSELLER
"As entertaining and thought-provoking as The Tipping Point by Malcolm Gladwell. . . . The Wisdom of Crowds ranges far and wide." —The Boston Globe

THE WISDOM OF CROWDS

JAMES SUROWIECKI

WITH A NEW AFTERWORD BY THE AUTHOR

"**Collective intelligence** can be brought to bear on a wide variety of problems, and complexity is no bar… conditions that are necessary for the crowd to be wise: *diversity*, *independence*, and … *decentralization*"

# Jeff Howe, WIRED, 2006

"**Crowdsourcing** represents the act of a company or institution taking a function once performed by employees and *outsourcing* it to an undefined (and generally large) network of people in the form of an open call. … The crucial prerequisite is the use of the *open call* format and the *large* network of potential laborers…"

http://www.wired.com/wired/archive/14.06/crowds.html

# "Human Computation", 2011

"**Human computation** is simply computation that is carried out by humans…
**Crowdsourcing** can be considered a method or a tool that human computation systems can use…"

By Edith Law & Luis von Ahn

# Daren Brabhan, 2013

CROWDSOURCING

DAREN C. BRABHAM

THE MIT PRESS ESSENTIAL KNOWLEDGE SERIES

"**Crowdsourcing** as an *online*, *distributed* problem-solving and production model that leverages the collective intelligence of online communities to serve specific *organizational* goals… *top-down* and *bottom-up* …"

# What is Crowdsourcing?

- Many definitions

- A few characteristics
  - **Outsourced** to **human** workers
  - **Online** and **distributed**
  - Open call & right **incentive**
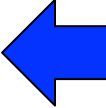  - **Diversity** and **independence**
  - Top-down & bottom-up

# What is **Computational** Crowdsourcing?

- Focus on computational aspect of crowdsourcing
  - Algorithmic aspect
  - Non-linear optimization problem

- Mainly use micro-tasks

- When to use Computational Crowdsourcing?
  1. Machine cannot do the task well
  2. Large crowds can probably do it well
  3. Task can be split to many micro-tasks

# Part 1: Crowdsourcing Basics

- Examples
- Definitions
- **Marketplaces** ⬅
- Computational Crowdsourcing
  - Preliminaries
  - Transcription
  - Sorting
- Demo

# Three Parties

- ## Requesters
  - People submit some tasks
  - Pay rewards to workers

- ## Marketplaces
  - Provide crowds with tasks

- ## Crowds
  - Workers perform tasks

**COMPANY**

**Survey**

**Submit tasks** ⬇    ⬆ **Collect answers**

**amazon** mechanicalturk™
Artificial Artificial Intelligence

**CrowdFlower**   **CloudCrowd**

**Find tasks** ⬇    ⬆ **Return answers**

# Notable Marketplaces

- Mechanical Turk
- CrowdFlower
- CloudCrowd
- Clickworker
- SamaSource

# AMT: mturk.com

Your Account | HITs | Qualifications

Introduction | Dashboard | Status | Account Settings

## Mechanical Turk is a marketplace for work.

We give businesses and developers access to an on-demand, scalable workforce.
We work with from thousands of tasks and work wh

**200,645 HITs** available. View the

**Workers**

**Requesters**

## Make Money
by working on HITs

HITs - *Human Intelligence Tasks* - are individual tasks that
you work on. Find HITs now.

**As a Mechanical Turk Worker you:**

- Can work from home
- Choose your own work hours
- Get paid for doing good work

Find an interesting task → Work → Earn money

TASKS

Find HITs Now

## Get Results
from Mechanical Turk Workers

Ask workers to complete HITs - *Human Intelligence Tasks* - and
get results using Mechanical Turk. Register Now

**As a Mechanical Turk Requester you:**

- Have access to a global, on-demand, 24 x 7 workforce
- Get thousands of HITs completed in minutes
- Pay only when you're satisfied with the results

Fund your account → Load your tasks → Get results

Get Started

# AMT: Workers vs. Requesters

- Workers
  - Register w. credit account (only US workers can register as of 2013)
  - Bid to do tasks for earning money
- Requesters
  - First deposit money to account
  - Post tasks
    - Task can specify a qualification for workers
  - Gather results
  - Pay to workers if results are satisfactory

# AMT: HIT

- Tasks
  - Called **HIT** (Human Intelligence Task)
  - Micro-task
- Eg
  - Data cleaning
  - Tagging / labeling
  - Sentiment analysis
  - Categorization
  - Surveying
  - Photo moderation
  - Transcription

Translate 3 lines from English to Russian (human translation needed).
**Requester:** Sergey Vasilyev      **Reward:** $0.05 per HIT      **HITs Available:** 1      **Duration:** 15 minutes
**Qualifications Required:** HIT approval rate (%) is not less than 75

**Translate a text between the markers below from English to Russian.**

**Human translation only! Machine tranlations will be rejected.**

================ FROM HERE ================

Hello!
I am test text message to be translated from English to Russian.
If you ask me, I was born in a mind of a crazy web developer,
who tests the MTurk API to start a very promising service later.

================ TILL HERE ================

**Any notes? Advices? Emotions? (Optional)**

**Translation task**

# Micro- vs. Macro-task: Eg, oDesk



**Workers**

**Requesters**

# AMT: HIT List

# AMT: HIT Example

## Can You Find the Provided Phone Number or Street Address on this Website?

Instructions ▲

### Overview

In this task, you'll be provided a web page for a business, including its name, address, and phone number. Your goal is to answer a few questions about the business on the web page.

**IMPORTANT:** Sometimes the business will have multiple locations, and you will have to search the website for the specific business that we provide in order to verify the website.

### Step by step instructions:

- Click the link to go to the provided site.

- First, please tell us whether or not the **name** of the business on the provided website is a **close** or **identical** match to the name of the business shown at the top of the page.

- Next, please tell us whether the provided business has

- Please be sure to click the appropriate option if the site

**Wrinkles Day Spa**

| | |
|---|---|
| Phone: | +61893455333 |
| Street: | Shop 5a Stirling Central Shopping Centre, 478 Wanneroo Rd |
| City: | Westminster |
| State: | WA |
| Postalcode (Zip): | 6061 |
| Country Code: | AU |

### Click here to visit the website.

**Is the name of the business on the web page similar or identical to 'Wrinkles Day Spa'?**

- ○ Yes: the name of the business is *similar* to *Wrinkles Day Spa*
- ○ Yes: the name of the business is *nearly identical* to *Wrinkles Day Spa*
- ○ No: the name is very different from *Wrinkles Day Spa*

ℹ For the first option, the street **number** does not need to match, just the street, **Shop 5a Stirling Central Shopping Centre, 4**

# AMT: HIT Example

# Open-Source Marketplace S/W



**Build with PYBOSSA**

The only open source framework for making crowdsourcing projects

Getting Started

or View the GitHub Project

GitHub | This repository | Search | rprise

gratipay / gratipay.com

Gratitude? Gratipay! Weekly payments to people and teams you believe in. https://gratipay.com/

GitHub | This repository | Search | Explore | Features

volontariat / voluntary

Engine and Framework for open source crowdsourcing platforms like Volontari.at

# Part 1: Crowdsourcing Basics

- Examples
- Definitions
- Marketplaces
- **Computational Crowdsourcing**
  - **Preliminaries** ⬅
  - Transcription
  - Sorting
- Demo

# Three Computational Factors

- ## Latency (or execution time)
  - Worker pool size
  - Job attractiveness

- ## Monetary cost
  - Cost per question
  - # of questions (ie, HITs)
  - # of workers

- ## Quality of answers
  - Worker maliciousness
  - Worker skills
  - Task difficulty

How long do we wait for?

Latency

Quality

How much is the quality of answers satisfied?

Cost

How much $$ does we spend?

# #1: Latency

- ## Some crowdsourcing tasks finish faster than others
  - ### Eg, easier, or more rewarding tasks are popular
- ## Dependency among tasks

| This is a password-protected HIT for a particular worker. | | | | View a HIT in this group |
|---|---|---|---|---|
| **Requester:** Eric DeRosia | **HIT Expiration Date:** | Oct 7, 2016 (104 weeks 3 days) | **Reward:** | $0.00 |
| | **Time Allotted:** | 24 hours | **HITs Available:** | 1 |

| Faculty Development | | | | View a HIT in this group |
|---|---|---|---|---|
| **Requester:** Kevin Dodds | **HIT Expiration Date:** | Oct 28, 2014 (3 weeks) | **Reward:** | $0.20 |
| | **Time Allotted:** | 15 minutes | **HITs Available:** | 16 |

| Rate an online article (required screening test) | | | | View a HIT in this group |
|---|---|---|---|---|
| **Requester:** HubPages | **HIT Expiration Date:** | Oct 7, 2014 (12 hours 10 minutes) | **Reward:** | $0.15 |
| | **Time Allotted:** | 5 days | **HITs Available:** | 15 |

# #2: Cost

- Cost per question
- # of HITs

Remaining cost to pay:
$0.03 X 2075 = $62.25



| Image Keyword Verification | | | | View a HIT in this group |
|---|---|---|---|---|
| Requester: Corbis Holdings, Inc | HIT Expiration Date: | Oct 13, 2014 (6 days 19 hours) | Reward: | $0.03 |
| | Time Allotted: | 15 minutes | HITs Available: | 2075 |

| Enter information about a forum discussion thread in which a vehicle is being built, rebuilt, or restored | | | | View a HIT in this group |
|---|---|---|---|---|
| Requester: Jonathan R | HIT Expiration Date: | Oct 14, 2014 (7 days 20 hours) | Reward: | $0.20 |
| | Time Allotted: | 30 minutes | HITs Available: | 2000 |

| Transcribe up to 25 Seconds of Media to Text – Low Priority | | | | View a HIT in this group |
|---|---|---|---|---|
| Requester: Crowdsurf Support | HIT Expiration Date: | Oct 20, 2014 (2 weeks) | Reward: | $0.08 |
| | Time Allotted: | 15 minutes | HITs Available: | 1836 |

| PADs_US_consumables_20140824-Thu Sep 11 16:15:50 PDT 2014 | | | | View a HIT in this group |
|---|---|---|---|---|
| Requester: Amazon Requester Inc. | HIT Expiration Date: | Oct 11, 2014 (5 days) | Reward: | $0.00 |
| | Time Allotted: | 1 hour 46 minutes | HITs Available: | 1605 |

# #3: Quality of Answers

- Avoid spam workers
- Use workers with reputation

Store name, date, time, total, location on this receipt
**Requester:** Vishwanath Kumar          **Reward:** $0.03 per HIT          **HITs Available:** 71477          **Duration:** 60 minutes
**Qualifications Required**    Total approved HITs is greater than 1000

- Ask the same question to multiple workers to get consensus (eg, majority voting)
- Assign more number of (better-skilled) workers to more difficult questions

# Size of Comparison

- Diverse forms of questions in a HIT
- Different sizes of comparisons in a question

**Binary question**

**Accuracy**



**Which is better?**

**Which is the best?**

. . .

**Cost Latency**

**N-ary question**

# Size of Batch

- Repetitions of questions within a HIT
- Eg, two *n*-ary questions (batch factor *b*=2)

# Response (*r*)

- # of human responses seeked for a HIT

*r* = 1

*r* = 3

**Which is better?**

**Which is better?**



**Cost, Latency**

**Accuracy**

**Smaller *r***

**Larger *r***

# Round (= Step)

- Algorithms are executed in rounds
- # of rounds ≈ latency

**Round #1**　　　　　　　　　**Round #2**

**Parallel Execution**

**Which is better?**

**Which is better?**

**Which is better?**

**Sequential Execution**

# Part 1: Crowdsourcing Basics

- Examples
- Definitions
- Marketplaces
- **Computational Crowdsourcing**
  - Preliminaries
  - **Transcription** ⬅
  - Sorting
- Demo

# Eg, Text Transcription [Miller-13]



- **Problem**: one person cannot do a good transcription
- **Key idea**: iterative improvement by many workers

**Greg Little** *et al*. "Exploring iterative and parallel human computation processes." HCOMP 2010

# Eg, Text Transcription [Miller-13]



**Handwriting Recognition Task - Mozilla Firefox**

- Please improve the transcription of this handwriting.
- People will vote whether to approve your changes.

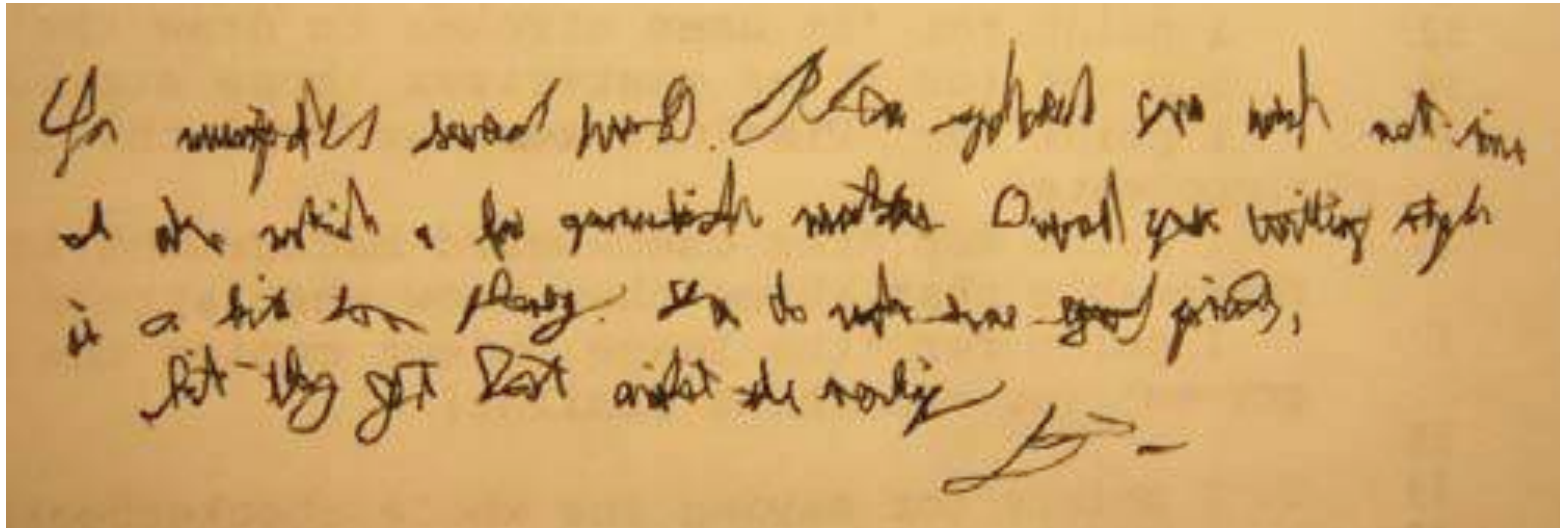You (?) (?) (?) (work). (?) (?) (?) work (not) (time). I (?) (?) a few grammatical mistakes. Overall your writing style is a bit too (phoney). You do (?) have good (points), but they got lost amidst the (writing). (signature)

**improvement $0.05**

Submit

# Eg, Text Transcription [Miller-13]



**3 votes @ $0.01**

# Eg, Text Transcription [Miller-13]

"You (misspelled) (several) (words). Please spellcheck your work next time. I also notice a few grammatical mistakes. Overall your writing style is a bit too phoney. You do make some good (points), but they got lost amidst the (writing). (signature)"
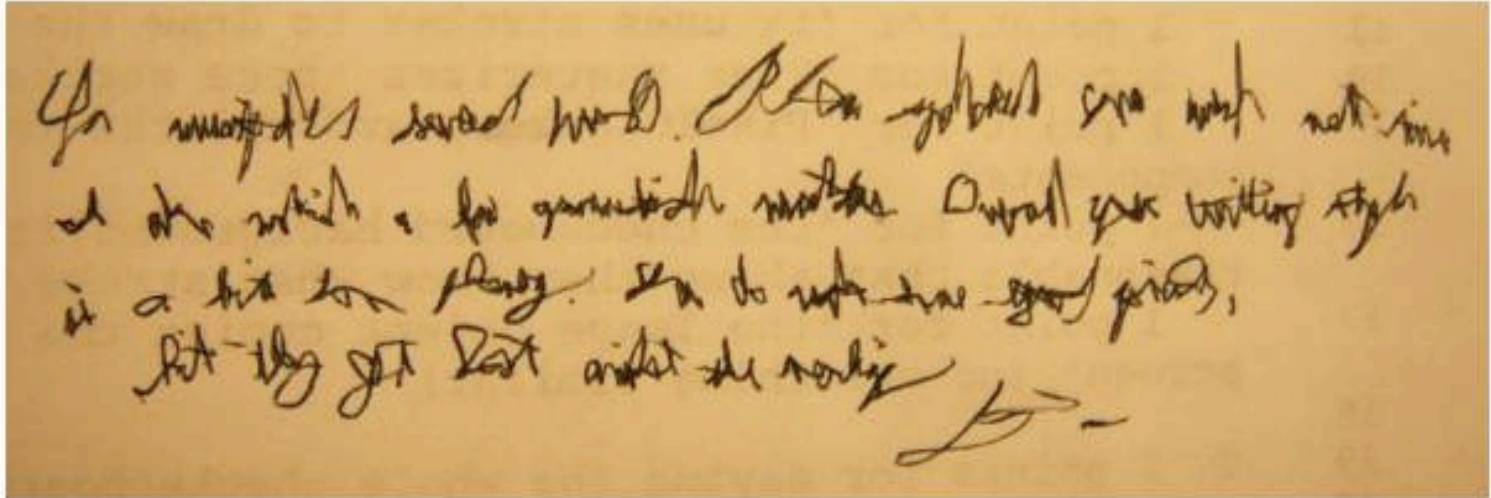
According to our ground truth, the highlighted words should be "flowery", "get", "verbiage" and "B-" respectively.

**After 9 iterations**

# Eg, Text Transcription [Miller-13]

I had intended to hit the nail, but I'm not a very good aim it seems and I ended up hitting my thumb. This is a common occurrence I know, but it doesn't make me feel any less ridiculous having done it myself. My new strategy will involve lightly tapping the nail while holding it until it is embedded into the wood enough that the wood itself is holding it straight and then I'll remove my hand and pound carefully away. We'll see how this goes.

## After 8 iterations
## with thousands of crowds

# Part 1: Crowdsourcing Basics

- Examples
- Definitions
- Marketplaces
- **Computational Crowdsourcing**
  - Preliminaries
  - Transcription
  - **Sorting**          ⬅
- Demo

# **Human-Powered Sort**

- Rank *N* items using crowdsourcing with respect to the constraint *C*

- Often *C* is subjective, fuzzy, ambiguous, and/or difficult-for-machines-to-compute

- Eg,
  - Which image is the most "representative" one of Brazil?
  - Which animal is the most "dangerous"?
  - Which actress is the most "beautiful"?

# Human-Powered Sort

```
SELECT      *
FROM        SoccerPlayers AS P
WHERE       P.WorldCupYear = '2014'
ORDER BY    CrowdOp('most-valuable')
```



. . .

# Naïve Sort

- Eg, "Which of two players is better?"
- Naïve all pair-wise comparisons takes $\binom{N}{2}$ comparisons
  - Optimal # of comparison is *O(N log N)*

# Naïve Sort

- Conflicting opinions may occur
  - o Cycle: A > B, B > C, and C > A

- If no cycle occurs
  - Naïve all pair-wise comparisons takes $\binom{N}{2}$ comparisons

- If cycle exists
  - More comparisons are required

# Sort [Marcus-VLDB11]

- N=5, S=3

# Sort [Marcus-VLDB11]

- N=5, S=3

**Sorted Result**

A

B

C

D

E



DAG

Topological Sort

A > B > C > E > D

# Part 1: Crowdsourcing Basics

- Examples
- Definitions
- Marketplaces
- Computational Crowdsourcing
  - Preliminaries
  - Transcription
  - Sorting
- **Demo**  ⬅

# Demo: Human-Powered Sorting

- From your smartphone or laptop, access the following URL or QR code:

`http://goo.gl/3tw7b5`

# Part 1 Conclusion

- Crowdsourcing ≈ Human Computation

- Academia: novel paradigm to solve the challenging problems in Computer Science

- Industry: novel entrepreneurial opportunities
  - Eg, Brazil-version Mechanical Turk?

This slide is available at

**http://goo.gl/4pNUhB**

# Reference

- **[Brabham-13]** *Crowdsourcing*, Daren Brabham, 2013
- **[Franklin-SIGMOD11]** *CrowdDB: answering queries with crowdsourcing*, Michael J. Franklin et al, SIGMOD 2011
- **[Howe-08]** *Crowdsourcing*, Jeff Howe, 2008
- **[LawAhn-11]** *Human Computation*, Edith Law and Luis von Ahn, 2011
- **[Li-HotDB12]** *Crowdsourcing: Challenges and Opportunities*, Guoliang Li, HotDB 2012
- **[Marcus-VLDB11]** *Human-powered Sorts and Joins*, Adam Marcus et al., VLDB 2011
- **[Miller-13]** *Crowd Computing and Human Computation Algorithms*, Rob Miller, 2013
- **[Shirky-08]** *Here Comes Everybody*, Clay Shirky, 2008

# Human-Powered Database Operations: Part 2

Dongwon Lee

Penn State University, USA

dongwon@psu.edu

Slide available @ **http://goo.gl/UEUEBh**

**SBBD 2014 Tutorial**

# Part 1: Crowdsourcing Basics

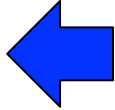- Examples
- Definitions
- Marketplaces
- Computational Crowdsourcing
  - Preliminaries
  - Transcription
  - Sorting
- Demo

# Part 2: Crowdsourced Algo. in DB

- **Preliminaries** ⬅
- Sort
- Select
- Count
- Top-1
- Top-$k$
- Join

# New Challenges

- Open-world assumption (OWA)

- Non-deterministic algorithmic behavior

- Trade-off among cost, latency, and accuracy



http://www.info.teradata.com

Database Predicate

False Propositions ? True Propositions

1094A086

0,1 p 1 q

Latency

Accuracy

Cost

# Crowdsourcing DB Projects

- CDAS @ NUS

- CrowdDB @ UC Berkeley
    & ETH Zurich

- MoDaS @ Tel Aviv U.

- Qurk @ MIT

- sCOOP @ Stanford & UCSC

# Part 2: Crowdsourced Algo. in DB

- Preliminaries
- **Sort** ⬅
- Select
- Count
- Top-1
- Top-$k$
- Join

# Sort Operation

- Rank *N* items using crowdsourcing w.r.t some criteria

- Assuming pair-wise comparison of 2 items

  - Eg, "Which of two images is better?"

- Cycle: A > B, B > C, and C > A

- If no cycle occurs

  - Naïve all pair-wise comparisons takes $\binom{N}{2}$ comparisons

- If cycle exists

  - More comparisons are required

# Sort [Marcus-VLDB11]

- Proposed 3 crowdsourced sort algorithms
- #1: <span style="color:red">Comparison-based Sort</span>
  - Workers rank $S$ items ($S \subset N$) per HIT
  - Each HIT yields $\binom{S}{2}$ pair-wise comparisons
  - Build a directed graph using all pair-wise comparisons from all workers
    - If $i > j$, then add an edge from $i$ to $j$
  - Break a cycle in the graph: "head-to-head"
    - Eg, If $i > j$ occurs 3 times and $i < j$ occurs 2 times, keep only $i > j$
  - Perform a topological sort in the DAG

# Sort [Marcus-VLDB11]

There are 2 groups of squares. We want to order the squares in each group from smallest to largest.

- Each group is surrounded by a dotted line. Only compare the squares within a group.
- Within each group, assign a number from 1 to 7 to each square, so that:
    - 1 represents the smallest square, and 7 represents the largest.
    - We do not care about the specific value of each square, only the relative order of the squares.
    - Some groups may have less than 7 squares. That is OK: use less than 7 numbers, and make sure they are ordered according to size.
    - If two squares in a group are the same size, you should assign them the same number.



**Error**

Submit

# Sort [Marcus-VLDB11]

- N=5, S=3

# Sort [Marcus-VLDB11]

- N=5, S=3

DAG

# Sort [Marcus-VLDB11]

- #2: <span style="color:red">Rating-based Sort</span>
  - *W* workers rate each item along a numerical scale
  - Compute the mean of *W* ratings of each item
  - Sort all items using their means
  - Requires *W*N* HITs: *O(N)*



**Mean rating**

| Worker | Rating |
| --- | --- |
| W1 | 4 |
| W2 | 3 |
| W3 | 4 |

. . .

| Worker | Rating |
| --- | --- |
| W1 | 1 |
| W2 | 2 |
| W3 | 1 |

1.3

3.6

. . .

8.2

# Sort [Marcus-VLDB11]

**There are 2 squares below. We want to rate squares by their size.**

- For each square, assign it a number from 1 (smallest) to 7 (largest) indicating its size.
- For perspective, here is a small number of other randomly picked squares:



smallest ○ ○ ✓ ○ ○ ○ ○ largest

1  2  3  4  5  6  7

smallest ○ ○ ○ ○ ✓ ○ ○ largest

1  2  3  4  5  6  7

Submit

# Sort [Marcus-VLDB11]

- #3: <span style="color:red">Hybrid Sort</span>
  - First, do rating-based sort $\rightarrow$ sorted list $L$
  - Second, do comparison-based sort on $S$ $(S \subset L)$

  - How to select the size of S
    - Random
    - Confidence-based
    - Sliding window

# Sort [Marcus-VLDB11]



Rank correlation btw. Comparison vs. rating

Worker agreement

# Sort [Marcus-VLDB11]

# Part II: Crowdsourced Algo. in DB

- Preliminaries
- Sort
- **Select** ⬅
- Count
- Top-1
- Top-$k$
- Join

# Select Operation

- Given *N* items, select *k* items that satisfy a predicate *P*

- ≈ Filter, Find, Screen, Search

# Select Operation

- Examples
  - **[Yan-MobiSys10]** uses crowds to search an image relevant to a query
  - **[Parameswaran-SIGMOD12]** develops human-powered filtering algorithms
  - **[Franklin-ICDE13]** efficiently enumerates items satisfying conditions via crowdsourcing
  - **[Sarma-ICDE14]** finds a bounded number of items satisfying predicates using the optimal solution by the skyline of cost and time

# Select [Yan-MobiSys10]

- Improving mobile image search using crowdsourcing

# Select [Yan-MobiSys10]

- Ensuring accuracy with majority voting
- Given accuracy, optimize cost and latency
- Deadline as latency in mobile phones

# Select [Yan-MobiSys10]

- Goal: For a query image *Q*, find the first relevant image *I* with **min cost** before the **deadline**

# Select [Yan-MobiSys10]

- Parallel crowdsourced validation

# Select [Yan-MobiSys10]

- Sequential crowdsourced validation

# Select [Yan-MobiSys10]

- CrowdSearch: using early prediction on the delay and outcome to start the validation of next candidate early

# Select [Yan-MobiSys10]

# Select [Parameswaran-SIGMOD12]

- Novel grid-based visualization

**Same person?**

Yes ○          No ○

# Select [Parameswaran-SIGMOD12]

- Common strategies

  - Always ask X questions, return most likely answer → Triangular strategy

  - If X YES return "Pass", Y NO return "Fail", else keep asking → Rectangular strategy

  - Ask until |#YES - #NO| > X, or at most Y questions → Chopped off triangle

# Select [Parameswaran-SIGMOD12]

- What is the best strategy? Find strategy with minimum overall expected cost s.t.

  1. **Overall expected error** is less than threshold

  2. **# of questions** per item never exceeds **m**

# Part 2: Crowdsourced Algo. in DB

- Preliminaries
- Sort
- Select
- **Count** ⬅
- Top-1
- Top-$k$
- Join

# Count Operation

- Given *N* items, estimate a fraction of items *M* that satisfy a predicate *P*

- Selectivity estimation in DB → crowd-powered query optimizers

- Evaluating queries with GROUP BY + COUNT/AVG/SUM operators

- Eg, "Find photos of females with red hairs"
  - Selectivity("female") ≈ 50%
  - Selectivity("red hair") ≈ 2%
  - Better to process predicate("red hair") first

# Count Operation

- Q: "How many teens are participating in the Hong Kong demonstration?"

# Count Operation

- Using Face++, guess the age of a person



**10 - 56**        **20 - 30**        **15 - 29**

http://www.faceplusplus.com/demo-detect/

# Count [Marcus-VLDB13]

- Hypothesis: Humans can estimate the frequency of objects' properties in a <span style="color:red">batch</span> without having to explicitly label each item

- Two approaches
  - #1: Label Count
    - Sampling based
    - Have workers label samples explicitly
  - #2: Batch Count
    - Have workers estimate the frequency in a batch

# Count [Marcus-VLDB13]

- Label Count (via sampling)

There are 2 people below. Please identify the gender of each.



What is the gender of this person?
○ male ● female

What is the gender of this person?
○ male ● female

Submit

# Count [Marcus-VLDB13]

- **Batch Count**

There are 10 people below. Please provide rough estimates for how many of the people have various properties.

About how many of the 10 people are <u>male</u>? 4

About how many of the 10 people are <u>female</u>?



Submit

# Count [Marcus-VLDB13]

- Findings on accuracy

  - Images: Batch count **>** Label count

  - Texts: Batch count **<** Label count

- Further Contributions

  - Detecting spammers

  - Avoiding coordinated attacks

# Part 2: Crowdsourced Algo. in DB

- Preliminaries
- Sort
- Select
- Count
- **Top-1** ⬅
- Top-*k*
- Join

# Top-1 Operation

- Find the top-1, either MAX or MIN, among *N* items w.r.t. some criteria


- Objective
  - Avoid sorting all *N* items to find top-1

# Top-1 Operation

- Examples
  - **[Venetis-WWW12]** introduces the bubble max and tournament-based max in a parameterized framework
  - **[Guo-SIGMOD12]** studies how to find max using pair-wise questions in the tournament-like setting and how to improve accuracy by asking more questions

# Max [Venetis-WWW12]

- **Introduced two Max algorithms**
  - Bubble Max
  - Tournament Max
- **Parameterized framework**
  - $s_i$: size of sets compared at the $i$-th round
  - $r_i$: # of human responses at the $i$-th round



Which is better?

$s_i = 2$
$r_i = 3$

Which is the best?

$s_i = 3$
$r_i = 2$

# Max [Venetis-WWW12]

- Bubble Max Case #1



$s_1 = 2$
$r_1 = 3$

$s_2 = 3$
$r_2 = 3$

$s_3 = 2$
$r_3 = 5$

- $N = 5$
- $Rounds = 3$
- # of questions =
  $r_1 + r_2 + r_3 = 11$

# Max [Venetis-WWW12]

- **Bubble Max Case #2**



- $N = 5$
- $Rounds = 2$
- # of questions = $r_1 + r_2 = 8$

$s_1 = 4$
$r_1 = 3$

$s_2 = 2$
$r_2 = 5$

# Max [Venetis-WWW12]

- **Tournament Max**



$s_1 = 2$
$r_1 = 1$

$s_3 = 2$
$r_3 = 3$

$s_2 = 2$
$r_2 = 1$

$s_4 = 2$
$r_4 = 5$

- $N = 5$
- $Rounds = 3$
- # of questions

$= r_1 + r_2 + r_3 + r_4 = 10$

# Max [Venetis-WWW12]

- How to find optimal parameters?: $s_i$ and $r_i$
- Tuning Strategies (using Hill Climbing)
  - Constant $s_i$ and $r_i$
  - Constant $s_i$ and varying $r_i$
  - Varying $s_i$ and $r_i$

# Max [Venetis-WWW12]

- ## Bubble Max

  - Worst case: with $s_i=2$, O(N) comparisons needed

- ## Tournament Max

  - Worst case: with $s_i=2$, O(N) comparisons needed

- ## Bubble Max is a special case of Tournament Max

# Max [Venetis-WWW12]

# Max [Venetis-WWW12]

# Part 2: Crowdsourced Algo. in DB

- Preliminaries
- Sort
- Select
- Count
- Top-1
- **Top-*k*** ⬅
- Join

# Top-*k* Operation

- Find top-*k* items among *N* items w.r.t. some criteria

- Top-*k* <span style="color:red">list</span> vs. top-*k* <span style="color:red">set</span>

- Objective
  - Avoid sorting all *N* items to find top-*k*

# Top-*k* Operation

- Examples
  - **[Davidson-ICDT13]** investigates the variable user error model in solving top-*k* list problem
  - **[Polychronopoulous-WebDB13]** proposes tournament-based top-*k* set solution

# Top-*k* Operation

- Naïve solution is to "sort" *N* items and pick top-*k* items

- Eg, *N*=5, *k*=2, "Find two best Bali images?"
  - Ask $\binom{5}{2}$ = 10 pair-wise questions to get a total order
  - Pick top-2 images

# Top-*k*: Tournament Solution (*k* = 2)

- Phase 1: **Building a tournament tree**
  - For each comparison, only winners are promoted to the next round



**Round 3**

**Round 2**

**Total, 4 questions with 3 rounds**

**Round 1**

# Top-*k*: Tournament Solution (*k* = 2)

- Phase 2: **Updating a tournament tree**
  - **Iteratively** asking pair-wise questions from the bottom level



**Round 3**

**Round 2**

**Round 1**

# Top-*k*: Tournament Solution (*k* = 2)

- Phase 2: **Updating a tournament tree**
  - **Iteratively** asking pair-wise questions from the bottom level



**Round 5**

**Round 4**

**Total, 6 questions With 5 rounds**

# Top-*k*: Tournament Solution

- This is a top-*k* **list** algorithm
- Analysis

| | k = 1 | k ≥ 2 |
|---|---|---|
| # of questions | $O(n)$ | $O(n + k \lceil \log_2 n \rceil)$ |
| # of rounds | $O(\lceil \log_2 n \rceil)$ | $O(k \lceil \log_2 n \rceil)$ |

- If there is no constraint for the number of rounds, this tournament sort based top-*k* scheme yields the optimal result

# Top-*k* [Polychronopoulous-WebDB13]

- ## Top-*k* **set** algorithm
  - Top-k items are "better" than remaining items
  - Capture NO ranking among top-*k* items

*K* items

  - Tournament-based approach

- ## Can become a Top-*k* **list** algorithm
  - Eg, Top-*k* **set** algorithm, followed by [Marcus-VLDB11] to sort *k* items

# Top-*k* [Polychronopoulous-WebDB13]

- Algorithm
  - Input: *N* items, integer *k* and *s* (ie, *s* > *k*)
  - Output: top-*k* set
  - Procedure:
    - $O \leftarrow N$ items
    - While $|O| > k$
      - Partition *O* into disjoint subsets of size *s*
      - Identify top-*k* items in each subset of size *s*: *s-rank(s)*
      - Merge all top-*k* items into *O*
    - Return *O*

- More effective when *s* and *k* are small
  - Eg, *s-rank*(20) with *k*=10 may give poor accuracy

# Top-*k* [Polychronopoulous-WebDB13]

- Eg, *N*=10, *s*=4, *k*=2

**Top-2 items**

**s-rank()**

**s-rank()**

**s-rank()**

**s-rank()**

**s-rank()**

**s-rank()**

# Top-*k* [Polychronopoulous-WebDB13]

- s-rank(s)

  // workers rank *s* items and aggregate

  - Input: *s* items, integer *k* (ie, *s* > *k*), *w* workers
  - Output: top-*k* items among *s* items
  - Procedure:
    - For each of *w* workers
      - Rank *s* items ≈ comparison-based sort [Marcus-VLDB11]
    - Merge *w* rankings of *s* items into a single ranking
      - Use median-rank aggregation [Dwork-WWW01]
    - Return top-*k* item from the merged ranking of *s* items

# Top-*k* [Polychronopoulous-WebDB13]

- Eg, s-rank(): *s=4*, *k=2*, *w=3*



**Top-2**

# Top-*k* [Polychronopoulous-WebDB13]

- Comparison to Sort [Marcus-VLDB11]

300 items,k=5,s=10, batch size 10, 20% spam



sort alg. (1303 HITs)
sort alg., 5 workers per batch(6516 HITs)
top-k alg. adaptive, (425-554 HITs)

# Top-*k* [Polychronopoulous-WebDB13]

- Comparison to Max [Venetis-WWW12]



10000 items, k=1, s=10

# Part 2: Crowdsourced Algo. in DB

- Preliminaries
- Sort
- Select
- Count
- Top-1
- Top-$k$
- **Join** ⬅

# Join Operation

- Identify matching records or entities within or across tables

  - ≈ similarity join, entity resolution (ER), record linkage, de-duplication, …

  - Beyond the exact matching

- [Chaudhuri-ICDE06] similarity join

  - $R$ JOIN$_p$ $S$, where p=$sim(R.A, S.A) > t$

  - $sim()$ can be implemented as UDFs in SQL

  - Often, the evaluation is expensive

    - DB applies UDF-based join predicate after Cartesian product of R and S

# Join Operation

- Examples
  - **[Marcus-VLDB11]** proposes 3 types of joins
  - **[Wang-VLDB12]** generates near-optimal cluster-based HIT design to reduce join cost
  - **[Wang-SIGMOD13]** reduces join cost further by exploiting transitivity among items
  - **[Whang-VLDB13]** selects right questions to ask to crowds to improve join accuracy
  - **[Gokhale-SIGMOD14]** proposes the hands-off crowdsourcing for join workflow

# Join [Marcus-VLDB11]

- To join tables *R* and *S*
- #1: <span style="color:red">Simple Join</span>
  - Pair-wise comparison HIT
  - $|R||S|$ HITs needed
- #2: <span style="color:red">Naïve Batching Join</span>
  - Repetition of #1 with a batch factor *b*
  - $|R||S|/b$ HITs needed
- #3: <span style="color:red">Smart Batching Join</span>
  - Show *r* and *s* images from *R* and *S*
  - Workers pair them up
  - $|R||S|/rs$ HITs needed

# Join [Marcus-VLDB11]

Is the same celebrity in the image on the left and the image on the right?

**#1 Simple Join**

Yes  No

# Join [Marcus-VLDB11]

# Join [Marcus-VLDB11]



#3 Smart Batching Join

# Join [Marcus-VLDB11]



MV: Majority Voting
QA: Quality Adjustment

# Join [Marcus-VLDB11]



Last 50% of wait time is spent completing the last 5% of tasks

# Join [Wang-VLDB12]

- [Marcus-VLDB11] proposed two batch joins
  - More efficient smart batch join still generates $|R||S|/rs$ # of HITs
  - Eg, (10,000 X 10,000) / (20 x 20) = 250,000 HITs → Still too many !

- [Wang-VLDB12] contributes CrowdER:
  1. A hybrid human-machine join
     - #1 machine-join prunes obvious non-matches
     - #2 human-join examines likely matching cases
       - Eg, candidate pairs with high similarity scores
  2. Algorithm to generate min # of HITs for step #2

# Join [Wang-VLDB12]

- Hybrid idea: generate candidate pairs using existing similarity measures (eg, Jaccard)



| ID | Product Name | Price |
|----|--------------|-------|
| $r_1$ | iPad Two 16GB WiFi White | $490 |
| $r_2$ | iPad 2nd generation 16GB WiFi White | $469 |
| $r_3$ | iPhone 4th generation White 16GB | $545 |
| $r_4$ | Apple iPhone 4 16GB White | $520 |
| $r_5$ | Apple iPhone 3rd generation Black 16GB | $375 |
| $r_6$ | iPhone 4 32GB White | $599 |
| $r_7$ | Apple iPad2 16GB WiFi White | $499 |
| $r_8$ | Apple iPod shuffle 2GB Blue | $49 |
| $r_9$ | Apple iPod shuffle USB Cable | $19 |

$(r_1, r_2, 0.57)$
$(r_4, r_6, 0.50)$
$(r_1, r_7, 0.43)$
$(r_3, r_4, 0.43)$
$(r_4, r_7, 0.43)$
$(r_8, r_9, 0.43)$
$(r_2, r_3, 0.38)$
$(r_2, r_7, 0.38)$
$(r_3, r_5, 0.38)$
$(r_4, r_5, 0.38)$
$(r_3, r_6, 0.29)$
$(r_1, r_3, 0.25)$
...

0.3

$(r_1, r_2)$ ⊙YES ○NO
$(r_4, r_6)$ ○YES ⊙NO

$(r_1, r_7)$ ⊙YES ○NO
$(r_3, r_4)$ ⊙YES ○NO

$(r_4, r_7)$ ○YES ⊙NO
$(r_8, r_9)$ ○YES ⊙NO

$(r_2, r_3)$ ○YES ⊙NO
$(r_2, r_7)$ ⊙YES ○NO

$(r_3, r_5)$ ○YES ⊙NO
$(r_4, r_5)$ ○YES ⊙NO

$(r_1, r_2)$
$(r_1, r_7)$
$(r_3, r_4)$
$(r_2, r_7)$

(a) Remove the pairs whose likelihood < 0.3

(b) Generate HITs to verify the pairs of records

c) Output matching pairs

**Main Issue: HIT Generation Problem**

# Join [Wang-VLDB12]

**Pair**-based HIT Generation
≈ Naïve Batching in
[Marcus-VLDB11]

**Cluster**-based HIT Generation
≈ Smart Batching in
[Marcus-VLDB11]

# Join [Wang-VLDB12]

- ## HIT Generation Problem
  - Input: pairs of records $P$, # of records in HIT $k$
  - Output: <span style="color:red">minimum</span> # of HITs s.t.
    1. All HITs have at most $k$ records
    2. Each pair $(p_i, p_j) \in P$ must be in at least one HIT

1. ## Pair-based HIT Generation
   - Trivial: $P/k$ # of HITs s.t. each HIT contains $k$ pairs in $P$

2. ## Cluster-based HIT Generation
   - <span style="color:red">NP-hard</span> problem → approximation solution

# Join [Wang-VLDB12]

| ID | Product Name | Price |
|----|--------------|-------|
| $r_1$ | iPad Two 16GB WiFi White | $490 |
| $r_2$ | iPad 2nd generation 16GB WiFi White | $469 |
| $r_3$ | iPhone 4th generation White 16GB | $545 |
| $r_4$ | Apple iPhone 4 16GB White | $520 |
| $r_5$ | Apple iPhone 3rd generation Black 16GB | $375 |
| $r_6$ | iPhone 4 32GB White | $599 |
| $r_7$ | Apple iPad2 16GB WiFi White | $499 |
| $r_8$ | Apple iPod shuffle 2GB Blue | $49 |
| $r_9$ | Apple iPod shuffle USB Cable | $19 |

$(r_1, r_2, 0.57)$
$(r_4, r_6, 0.50)$
$(r_1, r_7, 0.43)$
$(r_3, r_4, 0.43)$
$(r_4, r_7, 0.43)$
$(r_8, r_9, 0.43)$
$(r_2, r_3, 0.38)$
$(r_2, r_7, 0.38)$
$(r_3, r_5, 0.38)$
$(r_4, r_5, 0.38)$

$k = 4$

**Cluster-based HIT #1**

$r_1, r_2, r_3, r_7$

**Cluster-based HIT #2**

$r_3, r_4, r_5, r_6$

**Cluster-based HIT #3**

$r_4, r_7, r_8, r_9$

**This is the minimal # of cluster-based HITs satisfying previous two conditions**

# Join [Wang-VLDB12]

- ## Two-tiered Greedy Algorithm

  - Build a graph $G$ from pairs of records in $P$

  - CC $\leftarrow$ connected components in $G$

    - LCC: large CC with more than $k$ nodes

    - SCC: small CC with no more than $k$ nodes

  - Step 1: Partition LCC into SCCs

  - Step 2: Pack SCCs into HITs with $k$ nodes

    - Integer programming based

# Join [Wang-VLDB12]

- Eg, Generate cluster-based HITs ($k = 4$)
  1. Partition the LCC into 3 SCCs
     - $\{r_1, r_2, r_3, r_7\}$, $\{r_3, r_4, r_5, r_6\}$, $\{r_4, r_7\}$
  2. Pack SCCs into HITs
     - A single HIT per $\{r_1, r_2, r_3, r_7\}$ and $\{r_3, r_4, r_5, r_6\}$
     - Pack $\{r_4, r_7\}$ and $\{r_8, r_9\}$ into a HIT

# Join [Wang-VLDB12]

- ## Step 1: Partition
  - Input: LCC, $k$        Output: SCCs
  - $r_{max}$ ← node in LCC with the max degree
  - scc ← $\{r_{max}\}$
  - conn ← nodes in LCC directly connected to $r_{max}$
  - while $|scc| < k$ and $|conn| > 0$
    - $r_{new}$ ← node in conn with max indegree (# of edges to scc) and min outdegree (# of edges to non-scc) if tie
    - move $r_{new}$ from conn to scc
    - update conn using new scc
  - add scc into SCC

# Join [Wang-VLDB12]

# Join [Wang-VLDB12]



(a) Initialize scc={r₄}

(b) conn = {r₃, r₅, r₆, r₇}
Add r₆ into scc

(c) conn = {r₃, r₅, r₇}
Add r₅ into scc

(d) conn = {r₃, r₇}
Add r₃ into scc

(e) Output scc

(f) Output other scc

# Join [Wang-VLDB12]

# Join [Wang-SIGMOD13]

- Use the same hybrid machine-human framework as [Wang-VLDB12]
- Aim to reduce # of HITs further

- Exploit transitivity among records



http://blogs.oc.edu/ece/transitivity/

# Join [Wang-SIGMOD13]

- ## Positive transitive relation
    - ### If a=b, and b=c, then a=c

    > iPad $2^{nd}$ Gen = iPad Two
    >
    > iPad Two = iPad 2

    ➡ iPad $2^{nd}$ Gen = iPad 2

- ## Negative transitive relation
    - ### If a = b, b ≠ c, then a ≠ c

    > iPad $2^{nd}$ Gen = iPad Two
    >
    > iPad Two ≠ iPad 3

    ➡ iPad $2^{nd}$ Gen ≠ iPad 3

# Join [Wang-SIGMOD13]

- Three transitive relations
  - If there exists a path from o to o' which only consists of **matching pairs**, then (o, o') can be deduced as a **matching pair**
  - If there exists a path from o to o' which only contains **a single non-matching pair**, then (o, o') can be deduced as a **non-matching pair**
  - If any path from o to o' contains **more than one non-matching pairs**, (o, o') **cannot** be deduced.

# Join [Wang-SIGMOD13]



$(o_3, o_5) \rightarrow$ match

$(o_5, o_7) \rightarrow$ non-match

$(o_1, o_7) \rightarrow$ ?

# Join [Wang-SIGMOD13]

- Given a pair $(o_i, o_j)$, to check the transitivity
  - Enumerate path from $o_i$ to $o_j$ → exponential !
  - Count # of non-matching pairs in each path
- Solution: Build a cluster graph
  - Merge matching pairs to a cluster
  - Add inter-cluster edge for non-matching pairs



$(o_5, o_6) →$ ?

$(o_1, o_5) →$ ?

# Join [Wang-SIGMOD13]

- ## Problem Definition:
  - Given a set of pairs that need to be labeled, **minimize the # of pairs** requested to crowd workers based on **transitive relations**

| ID | Object |
|---|---|
| $o_1$ | iPhone 2nd Gen |
| $o_2$ | iPhone Two |
| $o_3$ | iPhone 2 |
| $o_4$ | iPad Two |
| $o_5$ | iPad 2 |
| $o_6$ | iPad 3rd Gen |

| ID | Object Pairs | Likelihood |
|---|---|---|
| $p_1$ | $(o_2, o_3)$ | 0.85 |
| $p_2$ | $(o_1, o_2)$ | 0.75 |
| $p_3$ | $(o_1, o_6)$ | 0.72 |
| $p_4$ | $(o_1, o_3)$ | 0.65 |
| $p_5$ | $(o_4, o_5)$ | 0.55 |
| $p_6$ | $(o_4, o_6)$ | 0.48 |
| $p_7$ | $(o_2, o_4)$ | 0.45 |
| $p_8$ | $(o_5, o_6)$ | 0.42 |

**?**

# Join [Wang-SIGMOD13]

- Labeling order matters !



$(o_1, o_2), (o_1, o_6), (o_2, o_6)$

vs.

$(o_1, o_6), (o_2, o_6), (o_1, o_2)$

➔ Given a set of pairs to label, how to order them affects the # of pairs to deduce using the transitivity

# Join [Wang-SIGMOD13]

- Theorem: Optimal labeling order

  $$w = <p_1, \ldots, p_{i-1}, p_i, p_{i+1}, \ldots, p_n>$$

  $$w' = <p_1, \ldots, p_{i-1}, p_{i+1}, p_i, \ldots, p_n>$$

  - If $p_i$ is a matching pair and $p_{i+1}$ is a non-matching pair, then $C(w) \leq C(w')$

    - $C(w)$: # of crowdsourced pairs required for $w$

- That is, always better to first label a matching pair and then a non-matching pair

- In reality, optimal label order cannot be achieved

# Join [Wang-SIGMOD13]

- Expected optimal labeling order

  - $w = <p_1, p_2, ..., p_n>$

  - $C(w)$ = # of crowdsourced pairs required for $w$

  $$\mathrm{E}\big[\mathcal{C}(\omega)\big] = \sum_{i=1}^{n} \mathbb{P}(p_i = \text{crowdsourced})$$

  - $P(p_i = \text{crowdsourced})$

    o Enumerate all possible labels of $<p_1, p_2, ..., p_{i-1}>$, and for each possibility, derive whether $p_i$ is crowdsourced or not

    o Sum of the probability of each possibility that whether $p_i$ is crowdsourced

# Join [Wang-SIGMOD13]

- ## Expected optimal labeling order

  - $w_1 = <p_1, p_2, p_3>$

  - $E[C(w_1)] = 1 + 1 + 0.05 = $ **2.05**

    - $P_1$: $P(P_1 = crowdsourced) = 1$

    - $P_2$: $P(P_2 = crowdsourced) = 1$

    - $P_3$: $P(P_3 = crowdsourced) = P(both\ P_1\ and\ P_2\ are\ non\text{-}matching) = (1\text{-}0.9)(1\text{-}0.5) = 0.05$

| Probability of matching | |
|---|---|
| $P_1$ | 0.9 |
| $P_2$ | 0.5 |
| $P_3$ | 0.1 |



| Expected value | |
|---|---|
| $w_1 = <p_1, p_2, p_3>$ | **2.05** |
| $w_2 = <p_1, p_3, p_2>$ | 2.09 |
| $w_3 = <p_2, p_3, p_1>$ | 2.45 |
| $w_4 = <p_2, p_1, p_3>$ | **2.05** |
| … | … |

# Join [Wang-SIGMOD13]

- Theorem: Expected optimal labeling order

  - Label the pairs in **the decreasing order of the probability** that they are a matching pair

  - Eg, $p_1$, $p_2$, $p_3$, $p_4$, $p_5$, $p_6$, $p_7$, $p_8$



$$E\big[\mathcal{C}(\omega)\big] = \sum_{i=1}^{n} \mathbb{P}(p_i = \text{crowdsourced})$$

| ID | Object Pairs | Likelihood | High |
|----|--------------|------------|------|
| $p_1$ | $(o_2, o_3)$ | 0.85 | |
| $p_2$ | $(o_1, o_2)$ | 0.75 | |
| $p_3$ | $(o_1, o_6)$ | 0.72 | |
| $p_4$ | $(o_1, o_3)$ | 0.65 | |
| $p_5$ | $(o_4, o_5)$ | 0.55 | |
| $p_6$ | $(o_4, o_6)$ | 0.48 | |
| $p_7$ | $(o_2, o_4)$ | 0.45 | |
| $p_8$ | $(o_5, o_6)$ | 0.42 | |

# Join [Wang-SIGMOD13]

- Two data sets
    - Paper: 997 (author, title, venue, date, and pages)
    - Product: 1081 product (abt.com), 1092 product (buy.com)



(a) Paper

(b) Product

# Join [Wang-SIGMOD13]

- Transitivity



(a) Paper

(b) Product

# Machine vs. Human

- Human-Powered Crowdsourcing → "**Human-in-the-loop**" Crowdsourcing

  - Should use machine to process majority of big data

  - Should use human to process a small fraction of challenging cases in big data

- How to split tasks and combine results for machines and human automatically is an open issue



http://www.theoddblog.us/2014/
02/21/damienwaltershumanloop/

# Conclusion

- New opportunities
  - Open-world assumption
  - Non-deterministic algorithmic behavior
  - Trade-off among cost, latency, and accuracy
- Crowdsourcing for Big Data?

This slide is available at

**http://goo.gl/UEUEBh**

# Reference

- **[Brabham-13]** *Crowdsourcing*, Daren Brabham, 2013
- **[Cao-VLDB12]** *Whom to Ask? Jury Selection for Decision Making Tasks on Microblog Services*, Caleb Chen Cao et al., VLDB 2012
- **[Chaudhuri-ICDE06]** A Primitive Operator for Similarity Join in Data Cleaning, Surajit Chaudhuri et al., ICDE 2006
- **[Davidson-ICDT13]** *Using the crowd for top-k and group-by queries*, Susan Davidson et al., ICDT 2013
- **[Dwork-WWW01]** *Rank Aggregation Methods for the Web*, Cynthia Dwork et al., WWW 2001
- **[Franklin-SIGMOD11]** *CrowdDB: answering queries with crowdsourcing*, Michael J. Franklin et al, SIGMOD 2011
- **[Franklin-ICDE13]** *Crowdsourced enumeration queries*, Michael J. Franklin et al, ICDE 2013
- **[Gokhale-SIGMOD14]** *Corleone: Hands-Off Crowdsourcing for Entity Matching*, Chaitanya Gokhale et al., SIGMOD 2014
- **[Guo-SIGMOD12]** *So who won?: dynamic max discovery with the crowd*, Stephen Guo et al., SIGMOD 2012
- **[Howe-08]** *Crowdsourcing*, Jeff Howe, 2008

# Reference

- **[LawAhn-11]** *Human Computation*, Edith Law and Luis von Ahn, 2011
- **[Li-HotDB12]** *Crowdsourcing: Challenges and Opportunities*, Guoliang Li, HotDB 2012
- **[Liu-VLDB12]** *CDAS: A Crowdsourcing Data Analytics System,* Xuan Liu et al., VLDB 2012
- **[Marcus-VLDB11]** *Human-powered Sorts and Joins*, Adam Marcus et al., VLDB 2011
- **[Marcus-VLDB12]** *Counting with the Crowd*, Adam Marcus et al., VLDB 2012
- **[Miller-13]** *Crowd Computing and Human Computation Algorithms*, Rob Miller, 2013
- **[Parameswaran-SIGMOD12]** *CrowdScreen: Algorithms for Filtering Data with Humans*, Aditya Parameswaran et al., SIGMOD 2012
- **[Polychronopoulous-WebDB13]** *Human-Powered Top-k Lists*, Vassilis Polychronopoulous et al., WebDB 2013
- **[Sarma-ICDE14]** Crowd-Powered Find Algorithms, Anish Das Sarma et al., ICDE 2014
- **[Shirky-08]** *Here Comes Everybody*, Clay Shirky, 2008

# Reference

- **[Surowiecki-04]** *The Wisdom of Crowds*, James Surowiecki, 2004
- **[Venetis-WWW12]** *Max Algorithms in Crowdsourcing Environments*, Petros Venetis et al., WWW 2012
- **[Wang-VLDB12]** *CrowdER: Crowdsourcing Entity Resolution*, Jiannan Wang et al., VLDB 2012
- **[Wang-SIGMOD13]** *Leveraging Transitive Relations for Crowdsourced Joins*, Jiannan Wang et al., SIGMOD 2013
- **[Whang-VLDB13]** *Question Selection for Crowd Entity Resolution*, Steven Whang et al., VLDB 2013
- **[Yan-MobiSys10]** *CrowdSearch: exploiting crowds for accurate real-time image search on mobile phones*, T. Yan et al., MobiSys 2010