



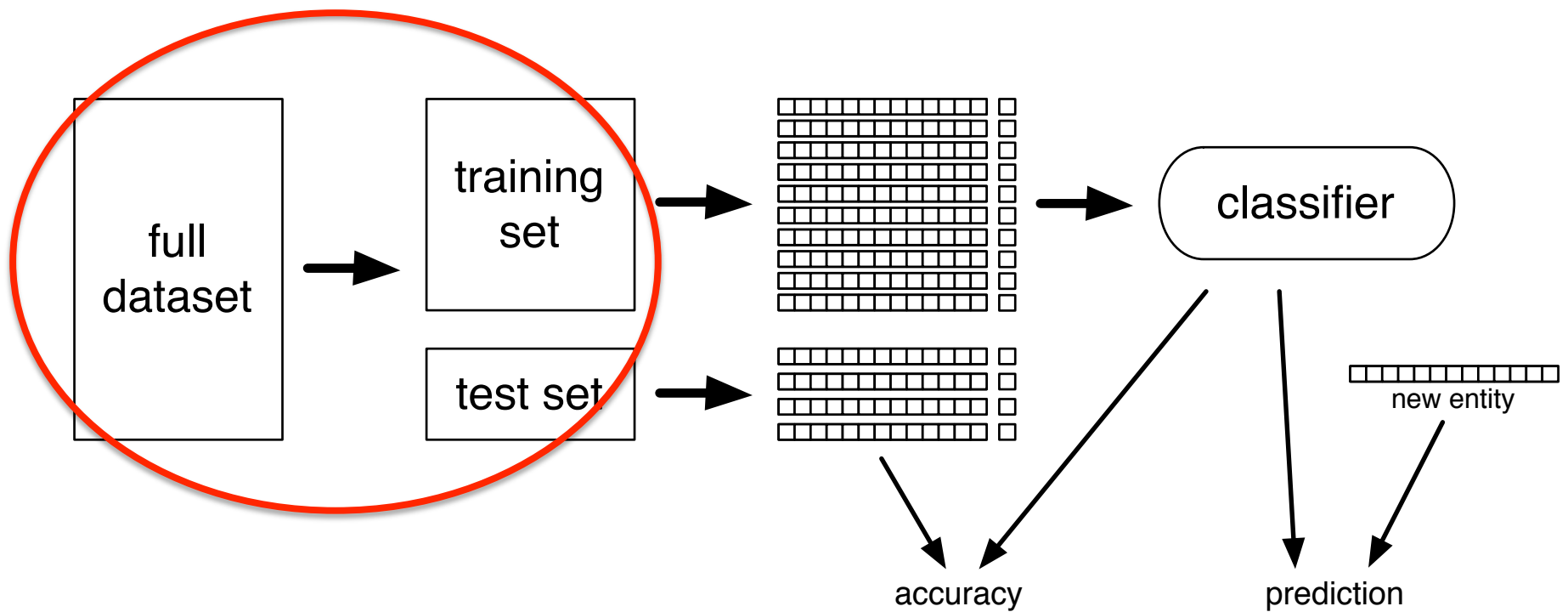
Crowdsourcing for Analytics

Tim Kraska

<kraska@cs.berkeley.edu>

Machines alone are not enough...

Classification



Machines alone are not enough...

SELECT Image
From Pictures
Where Image contains “Dog”

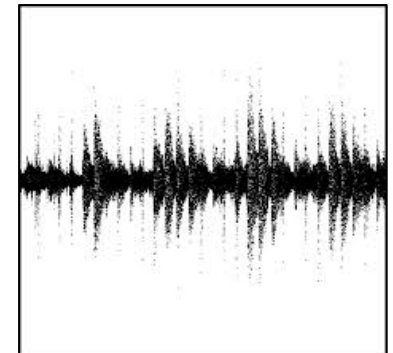


Adding People to Analytics

Data collection

America's top 10 NASDAQ companies with female CEOs

Transcription



Data cleaning

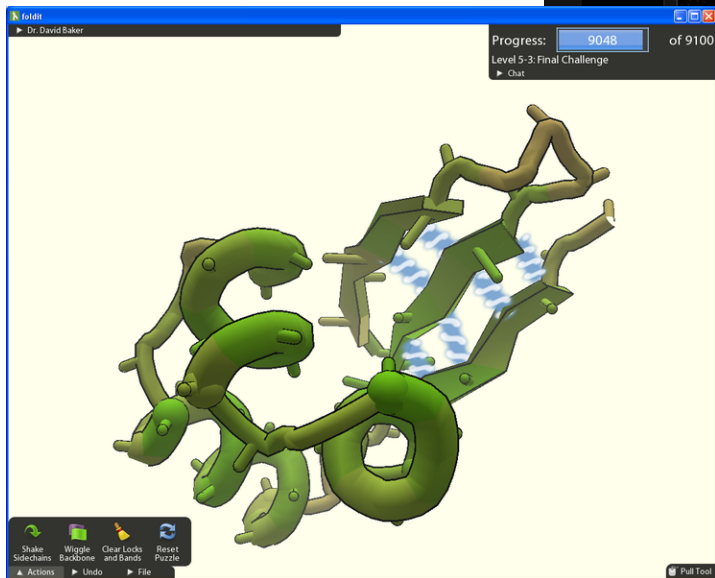
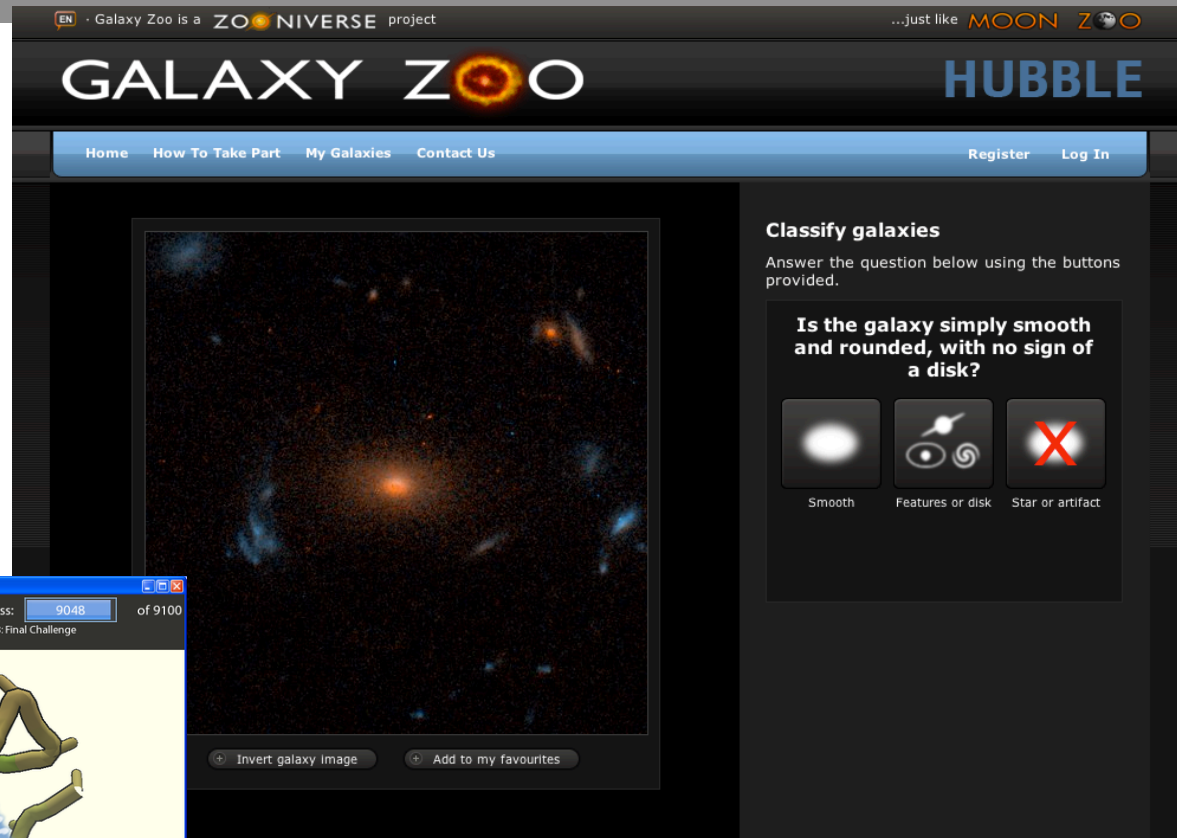


Creativity/Design/Taste



and more

Science



Knowledge Basis

Welcome to Q&A for professional and enthusiast programmers — check out the [FAQ!](#)

StackExchange v

[log in](#) | [careers](#) | [dev days](#) | [chat](#) | [meta](#) | [about](#) | [faq](#)

search



Questions

Tags

Users

Badges

Unanswered

Ask Question

Top Questions

interesting

237

featured

hot

week

month

0
votes

0
answers

1
view

[n Random rows for a given attribute - Postgres](#)

[sql](#) [postgresql](#)

44s ago [Sup3rkiddo](#) 49

1
vote

1
answer

14
views

[Branch descriptions in git, continued](#)

[git](#) [branch](#) [task-tracking](#)

48s ago [manojlds](#) 16.2k

0
votes

0
answers

1
view

[Where is hostname defined for the anchor element?](#)

[javascript](#)

56s ago [Chris Aaker](#) 868

2
votes

1
answer

12
views

[User-defined Table Variables in MySQL 5.5?](#)

[mysql](#) [stored-procedures](#) [routines](#)

1m ago [colonel_px](#) 11

0
votes

2
answers

37
views

[Closing cfpdf tag with </cfpdf> causes error](#)

[coldfusion](#) [coldfusion-8](#) [cfecclipse](#) [cfpdf](#)

1m ago [Jens Wegar](#) 81

0
votes

0
answers

6
views

[cocoa memory leak by CGAffineTranform or by view](#)

[iphone](#) [objective-c](#) [cocoa](#) [memory-leaks](#) [leak](#)

1m ago [EmptyStack](#) 9,100

Hello World!

This is a collaboratively edited question and answer site for **professional and enthusiast programmers**. It's 100% free, no registration required.

[about »](#) [faq »](#)

CAREERS 2.0

by stackoverflow

[Senior PHP Engineer](#)

[Spreetales](#)

Los Altos, CA; San Francisco, CA

[Front End Software Engineer](#)

[@Rdio](#)

[Rdio](#)

San Francisco, CA

[Senior Mobile Developer](#)


[American Public Media](#)

Oakland, CA

[Web Engineer](#)

[Monkey Inferno](#)

Structured Data



[Data](#) [Schema](#) [Apps](#) [Docs](#)

An entity graph of people, places and things, built by a community that loves open data.

Featured Data

Arts & Entertainment

Products & Services

Science & Technology

Society

Special Interests

Sports

System

Time & Space


Transportation

All

Sort by write activity

Film

80 members



43K last week


5M

641K



People


87 members



11K last week


7M

2M



TV


35 members



8K last week

8M

1M



Music

100+ members



771 last week

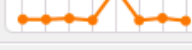
35M

10M



Business


100+ members



431 last week


2M

611K



Government

47 members



240 last week

532K

135K



Location


52 members



223 last week


10M

999K



Books

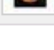
47 members



108 last week

29M

6M



**Google Refine**
An open source power tool to fix, discover, experiment, connect and customize your data. [Learn more »](#)

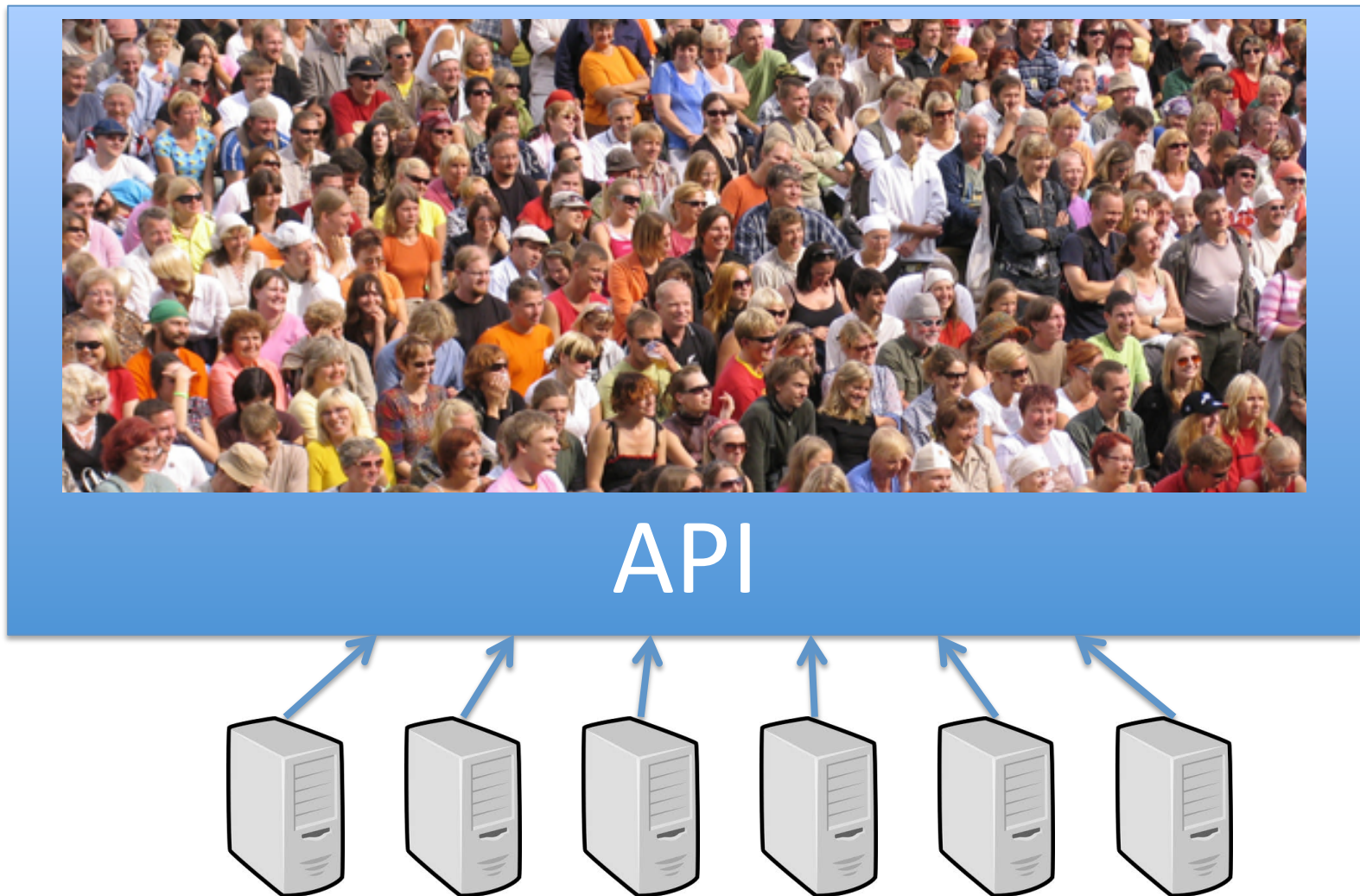
What is Freebase?
Learn what an entity graph is, what kind of information it contains, and why you should add your data!
[Learn More »](#)

Freebase for Developers

- powerful queryable API
- JavaScript-based hosting framework
- libraries for other languages

[Learn More »](#)

Crowdsourcing for Developers 101



Micro-Task CrowdSourcing



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



Get Results from Mechanical Turk Workers

Ask workers to complete HITs - *Human Intelligence Tasks* - and get results using Mechanical Turk. [Get started.](#)

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



Microtasking – Virtualized Humans

- Current leader: Amazon Mechanical Turk
- Requestors place Human Intelligence Tasks (HITs)
 - Minimum price: \$0.01
 - #of replicas (assignments), expiration, User Interface
 - API-based: “createHit()”, “getAssignments()”, “approveAssignments()”, “forceExpire()”
 - Requestors approve jobs and payment
- Workers (a.k.a. “turkers”) choose jobs, do them, get paid

Your Account

HITS

Qualifications

274,745 HITS
available now

[All HITS](#) | [HITS Available To You](#) | [HITS Assigned To You](#)

Find

HITS

containing

that pay at least \$

0.00

☐ for which you are qualified

☐ require Master Qualification

GO

All HITS

1-10 of 1157 Results

Sort by: HIT Creation Date (newest first)

GO!

[Show all details](#) | [Hide all details](#)

1 2 3 4 5 > [Next](#) >> [Last](#)

huge test

[View a HIT in this group](#)

Requester: [Mr Doe](#)

HIT Expiration Date: Aug 27, 2011 (1 day 5 hours)

Reward: \$0.01

Time Allotted: 60 seconds

HITS Available: 8424

Validate Brand/Product Information from Product Picture

[View a HIT in this group](#)

Requester: [Redwood Technologies](#)

HIT Expiration Date: Aug 27, 2011 (11 hours 59 minutes)

Reward: \$0.01

Time Allotted: 15 minutes

HITS Available: 6

Copy Brand/Product Information from Product Picture

[View a HIT in this group](#)

Requester: [Redwood Technologies](#)

HIT Expiration Date: Aug 27, 2011 (11 hours 58 minutes)

Reward: \$0.03

Time Allotted: 15 minutes

HITS Available: 8

huge job test for refactored balancer

[View a HIT in this group](#)

Requester: [Mr Doe](#)

HIT Expiration Date: Aug 27, 2011 (1 day 5 hours)

Reward: \$0.01

Time Allotted: 60 seconds

HITS Available: 12390

[All HITS](#) | [HITS Available To You](#) | [HITS Assigned To You](#)

Find

containing

that pay at least \$

☐ for which you are qualified

☐ require Master Qualification

GO

Timer: 00:00:00 of 5 minutes

Want to work on this HIT?

Want to see other HITS?

Accept HIT

Skip HIT

Total Earned: \$0.88
Total HITS Submitted: 80

CrowdDb Equal 7ae03d48-04f0-4259-8d81-492115e106ae

Requester: AMPLab

Reward: \$0.01 per HIT

HITS Available: 5

Duration: 5 minutes

Qualifications Required: None



'Are these two pictures of the same person?'

and



Yes ☐ No ☐

Please ACCEPT the hit before submitting.

Amount per Assignment:	\$0.01	Pending Review:	1
Amount to Approve Outstanding Assignments:	\$0.01 What's this?	Reviewed:	0
		Remaining:	0 Add Assignments
		Total:	1

EXPIRATION DATE

Aug 27 2011, 02:25 PM PDT [Add Time](#)

[» View HIT](#)

[Pending Review \(1\)](#) [Approved](#) [Rejected](#)

Review Submitted Assignments (showing page 1 of 1)

Select Assignments to approve or reject then click "Submit." When you approve an Assignment, the Worker is paid automatically. You will not be charged for Assignments you reject.

Approve	Reject	Worker ID	Result	Less	Submission Date
All · None	All · None				
<input type="checkbox"/>	<input type="checkbox"/>	A5RPKYH18EW	crowdEqual:	no	Aug 26 2011, 02:33 PM PDT
			jsessionid:	48f4e460c2ac00c3df877f1cdbe5813d	
			crowdDbCallback:	http://128.32.45.115:8082/crowdDbEqual.do? jsessionid=48f4e460c2ac00c3df877f1cdbe5813d	
			wid:		
			id:	10202458	

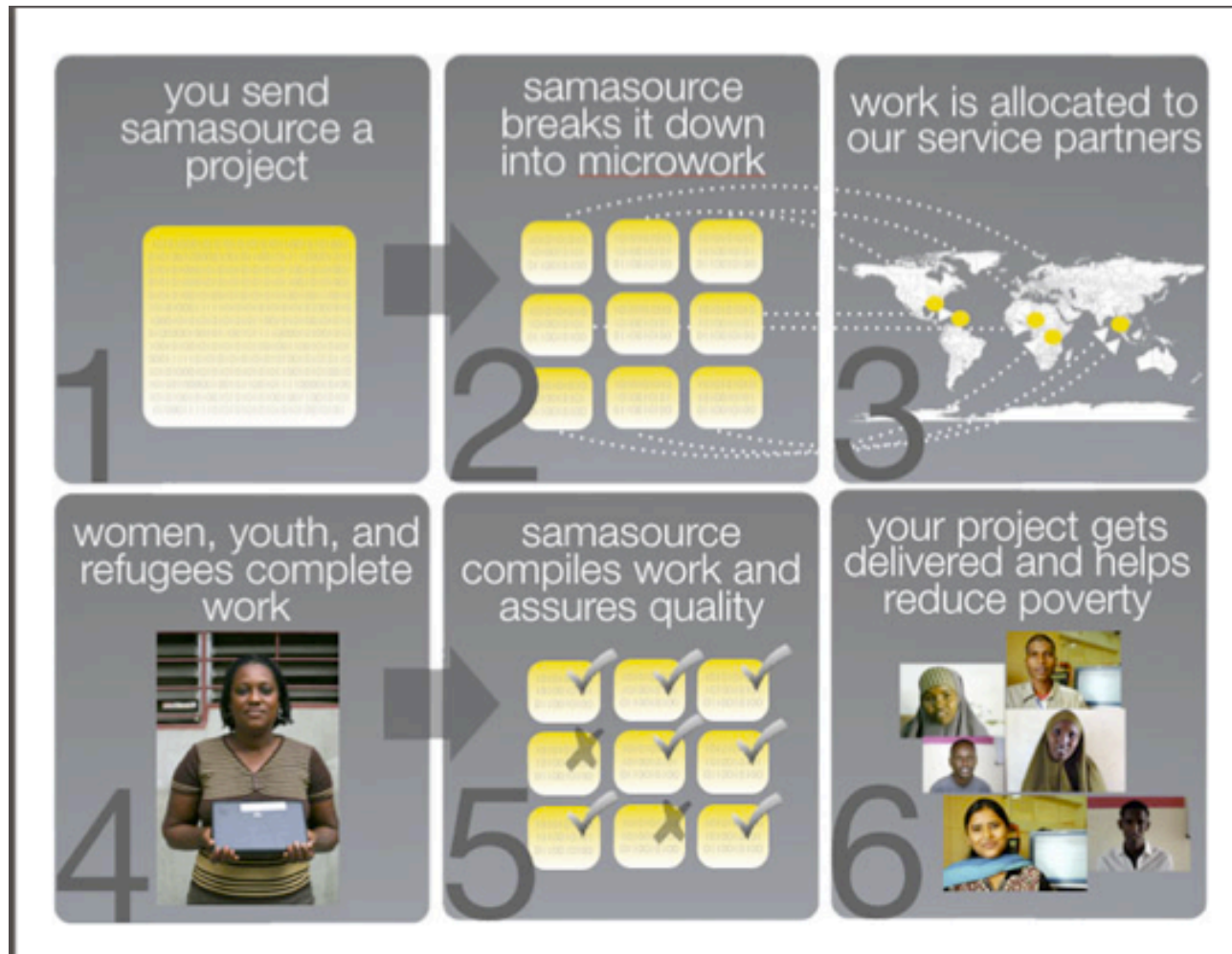
You've chosen to:

approve: none selected

reject: none selected

[Submit](#)

Samasource.org



Challenges

- Quality
- User Interface Design
- Worker motivation
- Task decomposition
- Leverage worker knowledge/capabilities
- Optimization (time/cost)
-

Challenges

- **Quality**
- Optimization (time/cost)
- User Interface Design
- Worker Motivation
- Task decomposition
- Leverage worker knowledge/capabilities
-

How Can You Trust the Crowd?



Quality Techniques

- Approval Rate / Demographic Restrictions
- Gold Sets/Honey Pots
- Redundancy
- Qualification Test
- Verification/Review
- Justification/Automatic Verification
- ...

Quality Techniques

- **Approval Rate / Demographic Restrictions**
- **Gold Sets/Honey Pots**
- **Redundancy**
- Qualification Test
- Verification/Review
- Justification/Automatic Verification
- ...

Approval Rate & Demographic Restrictions

Classify text about consumer electronics View a HIT in this group

Requester: [Buzz Evaluation](#) **HIT Expiration Date:** Sep 7, 2011 (1 week 6 days) **Reward:** \$0.02

Time Allotted: 20 minutes **HITs Available:** 3966

Description: Classify text for positive, negative, mixed or neutral tone

Keywords: [buzz](#), [classify](#), [coding](#), [tag](#), [sentiment](#), [text](#), [analysis](#), [twitter](#), [blog](#), [social](#)

Qualifications Required:
HIT approval rate (%) is not less than 95
Location is US

- + Easy to setup
- + Transparent
- Easy to defeat
- Causes a lot of trouble

Approval Rate



HIT Group » I recently did **299 HITs for this requester....** Of the 299 HITs I completed, **11 of them were rejected** without any reason being given. **Prior to this I only had 14 rejections, a .2% rejection rate. I currently have 8522 submitted HITs, with a 0.3% rejection rate** after the rejections from this requester (25 total rejections). I have attempted to contact the requester and will update if I receive a response. Until then be very wary of doing any work for this requester, as it appears that they **are rejecting about 1 in every 27 HITs being submitted.** posted by ...

fair:2 / 5 fast:4 / 5 pay:2 / 5 comm:0 / 5

Gold Sets / Honey Pots



- Gold derived from
 - Experts
 - Crowd using high quorum
- Interject trap questions
- Block users in trap and invalidate answers
- + **Often very effective**
- + **Cost efficient**
- **Not always applicable**
- **Digging gold is hard**

Defeating Honey Pots: reCAPTCHA

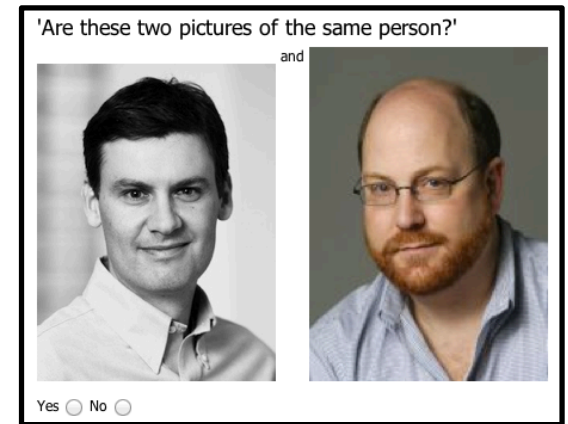
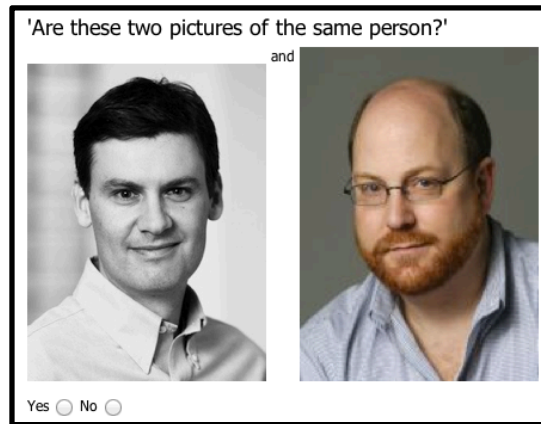
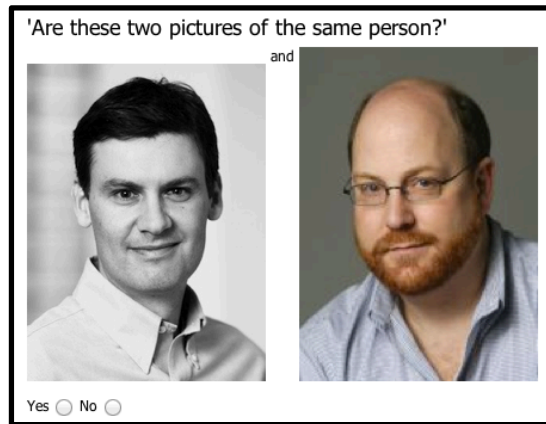


data

honey

1629 capplots | 1960-73 batizak | III, Quichot
RICHARD, Redfully | ansunjon Rupprecht | Daberny Nature
drôle Deratial | than Morgicar | Golightly, Byemore

Redundancy: Quorum Votes



majority vote

result

- + Easy to implement
- + Hard to defeat
- Increased cost
- Masks cases of ambiguity or diversity, “tail” behaviors
- Does not cover bias

Challenges

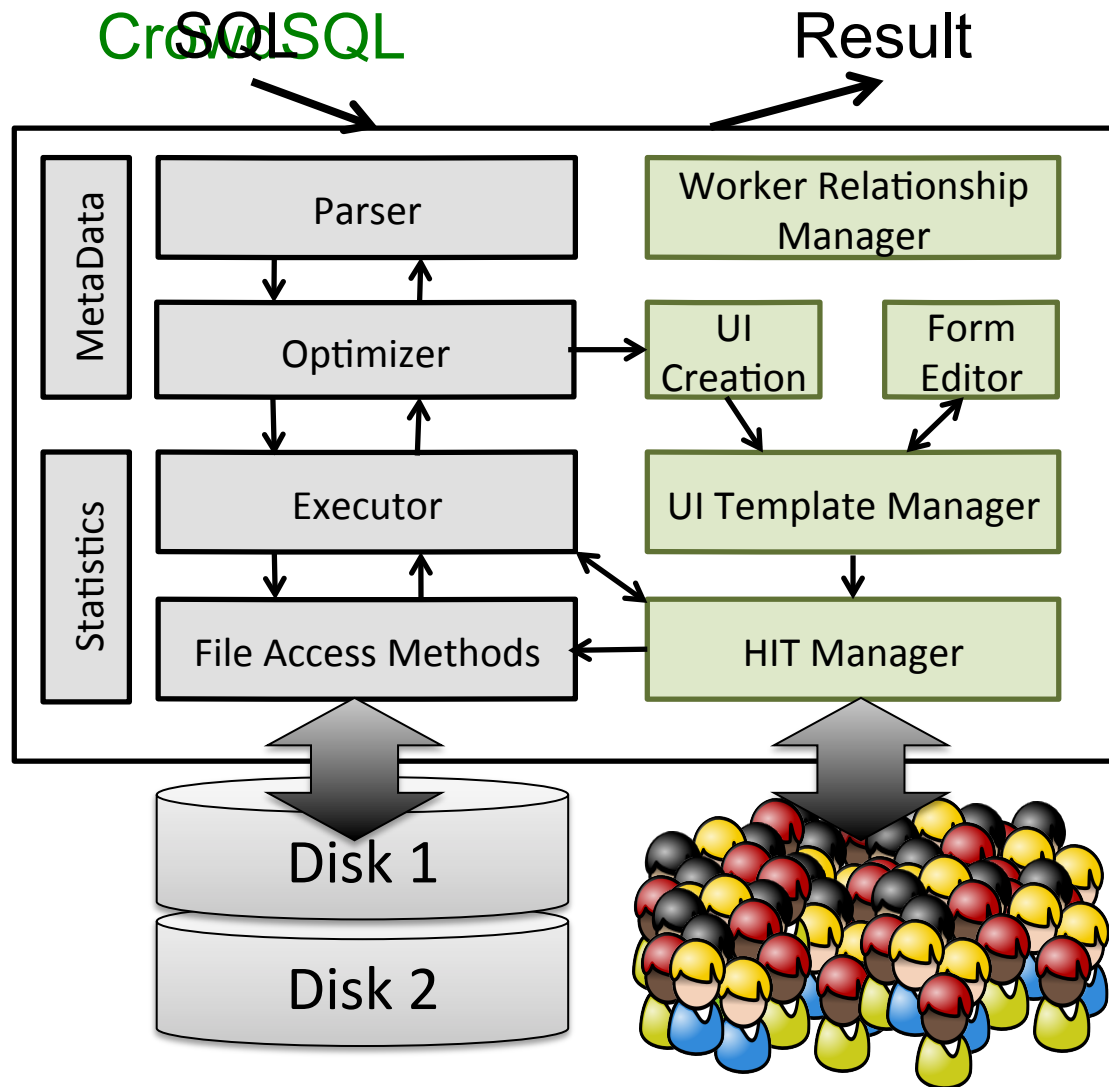
- Quality
- User Interface Design
- Worker motivation
- Task decomposition
- Leverage worker knowledge/capabilities
- Optimization (time/cost)
-



Use Cases

- Data collection:
 - How do my prices compare to the prices of my competitors
 - Finding job candidates (who is graduating from HPI next year)
 - Find green-tech companies in the Bay Area
 - ...
- Data cleaning
 - Verifying customer addresses
 - Duplicate elimination
 - ...
- Extending data
 - Labeling (spam/not_spam)
 - ...

CrowdDB



CrowdSQL

DDL Extensions:

Crowdsourced columns

```
CREATE TABLE company (  
  name STRING PRIMARY KEY,  
  hq_address CROWD STRING);
```

Crowdsourced tables

```
CREATE CROWD TABLE department (  
  university STRING,  
  department STRING,  
  phone_no STRING)  
PRIMARY KEY (university, department);
```

DML Extensions:

CrowdEqual:

```
SELECT *  
FROM companies  
WHERE Name ~ "Big Blue"
```

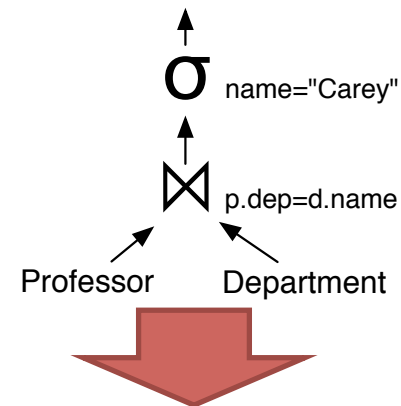
CROWDORDER operators (currently UDFs):

```
SELECT p FROM picture  
WHERE subject =  
  "Golden Gate Bridge"  
ORDER BY CROWDORDER(p, "which  
pic shows better %subject");
```

Optimization: Quality

CROWD TABLE professor(name, e-mail)
CROWD TABLE department(name, phone-nb)

SELECT *
FROM professor p, department d
WHERE d.name = p.dep
AND p.name = "Michael J. Carey"



MTJoin (Professor)
p.name = "carey"

Please fill out the missing professor data

Name: Carey

Department name: CS

E-Mail:

Submit

(Department first)
Inefficient

MTJoin (Dep)
p.dep = d.name

Please fill out the missing department data

Department Name: CS

Phone:

Submit

Please fill out the missing professor data

Name: Carey

E-Mail:

Department:

Submit

(Professor first)
≈10% Error-Rate

MTProbe (Professor, Dep)
name=Carey

Please fill out the missing professor data

Name: Carey

E-Mail:

Department:

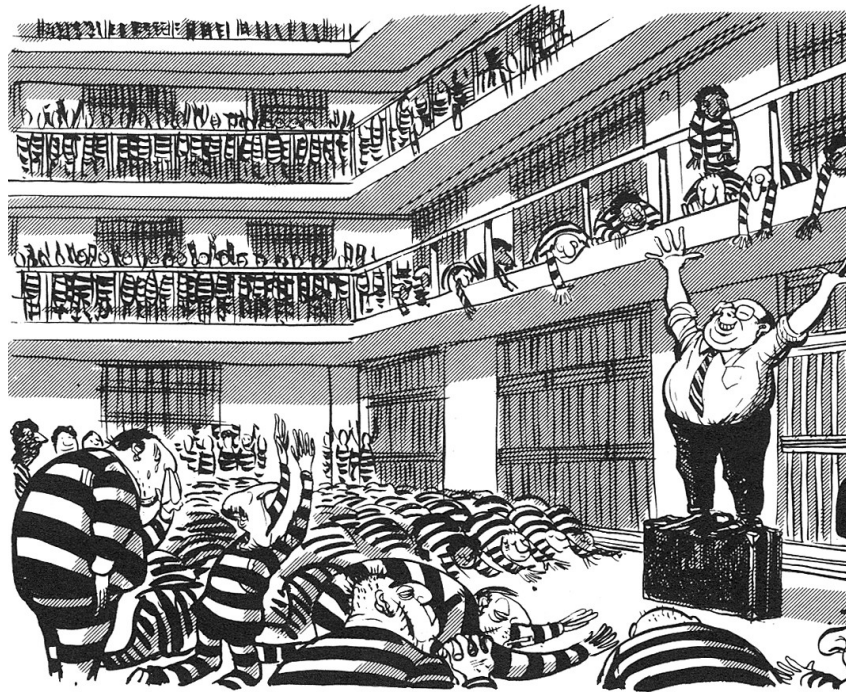
Department Phone:

Submit

(De-normalized Probe)
≈80% Error-Rate

A Bigger(?) Underlying Issue

Closed-World



Open-World



What Does This Query Mean?

```
SELECT COUNT(*) FROM IceCreamFlavors
```

Enter a flavor of ice cream

Requester: trush **Reward:** \$0.01 per HIT **HITs Available:** 51 **Duration:** 10 minutes

Qualifications Required: HIT approval rate (%) is greater than 90

Enter a flavor of ice cream

In the textbox below, please enter a flavor of ice cream

answer:

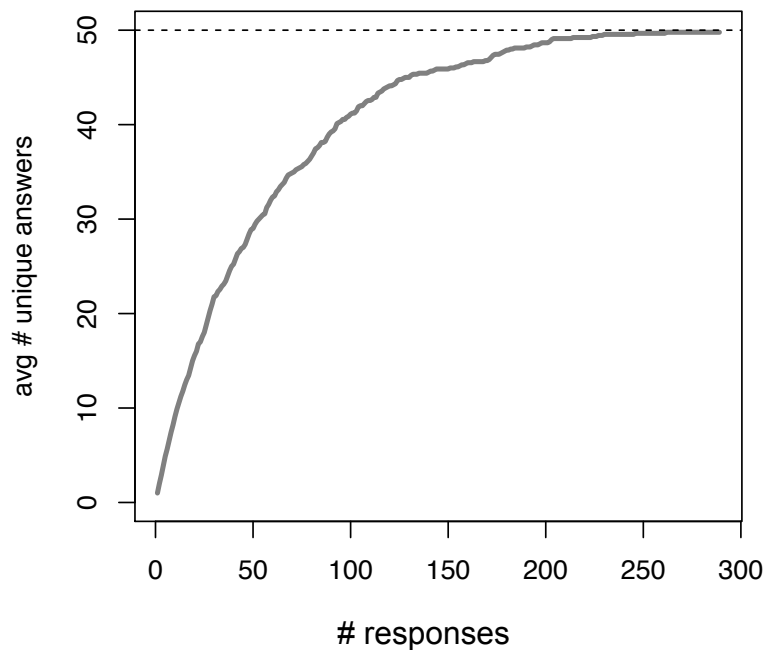
Trushkowsky *et al.* Getting it All From the Crowd, (in preparation) on arxiv

Estimating Completeness

`SELECT COUNT(*) FROM US States`

US States using Mechanical Turk

Unique items over time

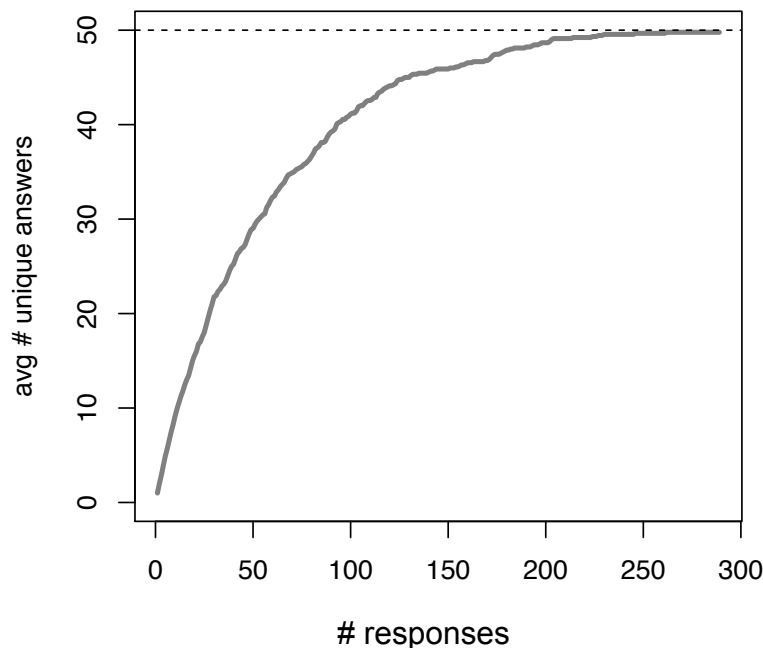


Estimating Completeness

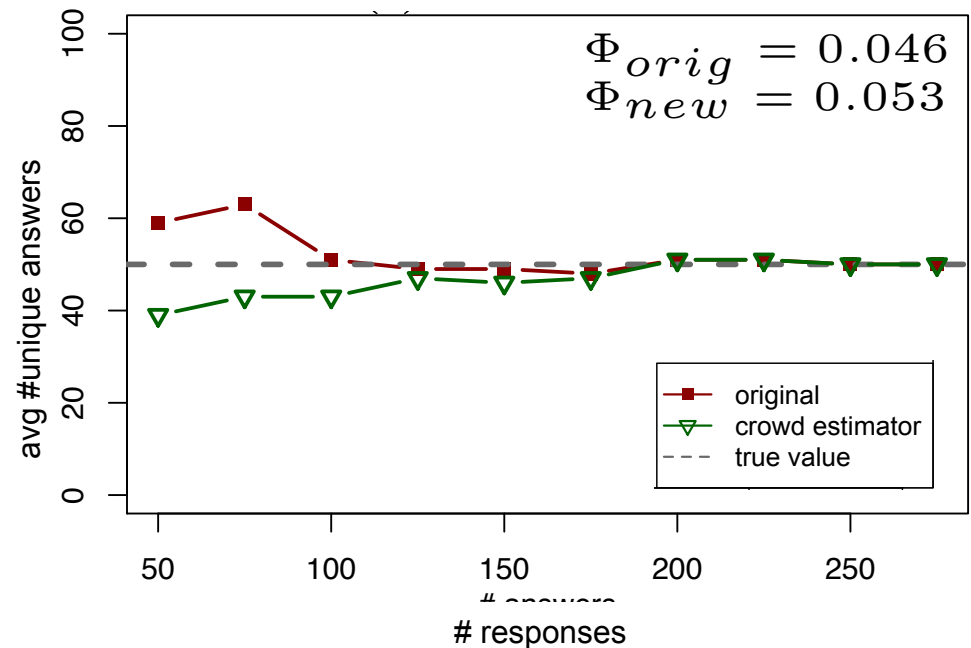
SELECT COUNT(*) FROM *US States*

US States using Mechanical Turk

Unique items over time



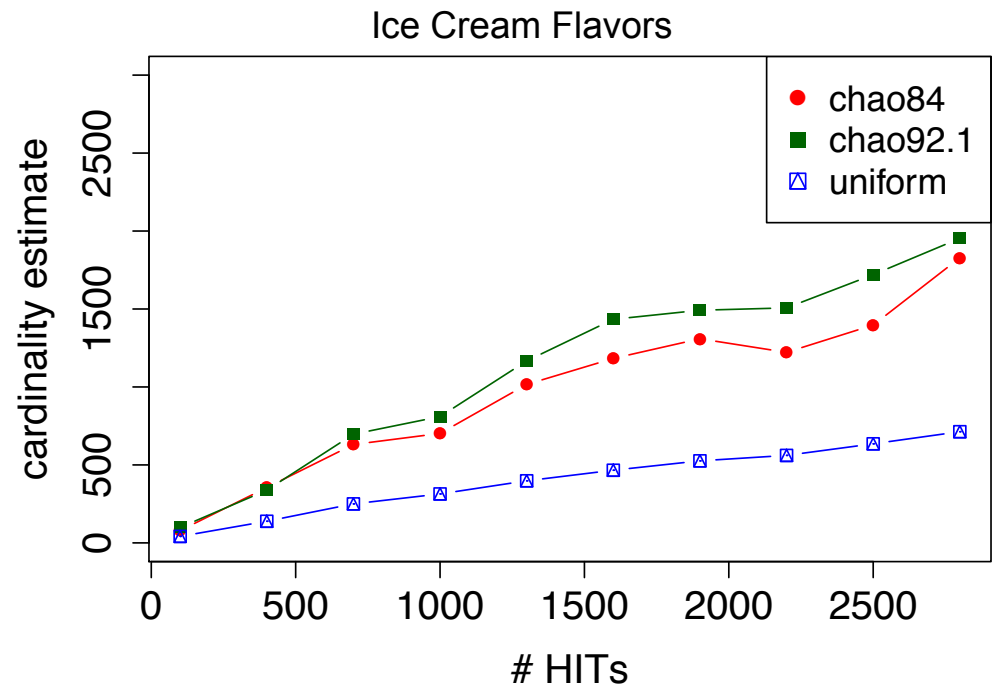
Crowd set-size estimation



Estimating Completeness

SELECT COUNT(*) FROM *IceCreamFlavors*

- Ice Cream Flavors
 - Estimators don't converge
 - Very highly skewed (CV = 5.8)
 - Detect that # HITs insufficient (beginning of curve)



Few, short lists of ice cream flavors
(e.g. “alumni swirl, apple cobbler crunch,
arboretum breeze,...” from Penn State
Creamery)

Pay-As-You-Go

- *“I don’t believe it is usually possible to estimate the number of species... but only an appropriate lower bound for that number. This is because there is nearly always a good chance that there are a very large number of extremely rare species”*
– Good, 1953
- So instead, can ask: “What’s the benefit of m additional HITs?”

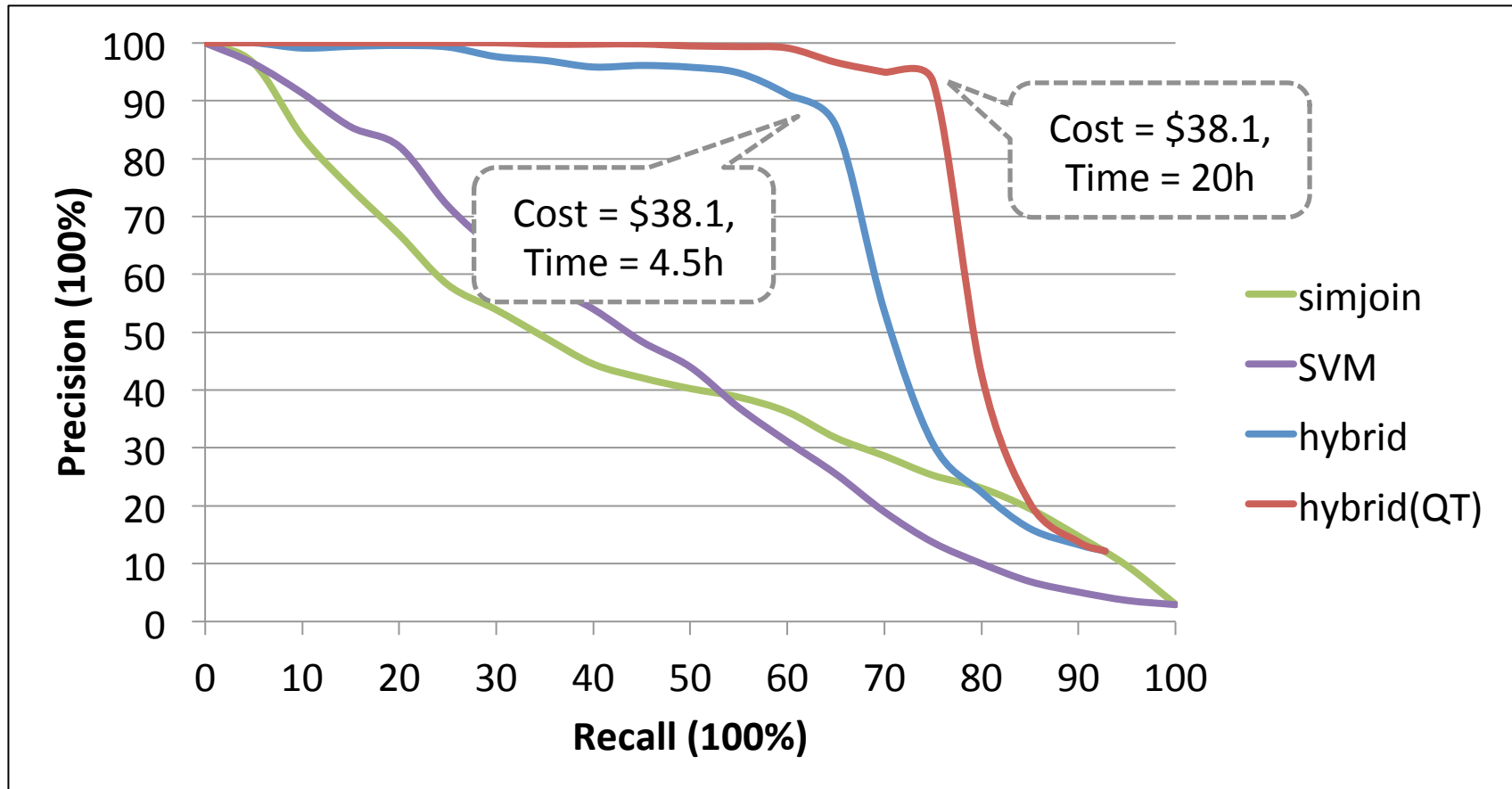
Ice Cream after 1500 HITs

m	Actual	Shen	Spline
10	1	1.79	1.62
50	7	8.91	8.22
200	39	35.4	32.9

Entity Resolution

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

Hybrid Entity Resolution



J. Wang *et al.* CrowdER: Crowdsourcing Entity Resolution, PVLDB 2012

Human-Tolerant Computing

Adding People into the Analytics Lifecycle:

- Inconsistent answer quality
- Incentives
- Latency & Variance
- Open vs. Closed world
- Hybrid Human/Machine Design



Approaches:

- Statistical methods for error and bias
- Quality-conscious Interface design
- Cost (time, quality)-based optimization

Summary

The AMPLab looks into integrating Algorithms, Machines and People for big data analytics

- Crowdsourcing can help with Big Data analytics where machines are not enough
- CrowdDB is a first hybrid Crowd/Cloud data management system following this vision
- Full tutorial: *Crowdsourcing Applications and Platforms: A Data Management Perspective. VLDB, 2011*
- Try it out at mturk.com

Tim Kraska
kraska@cs.berkeley.edu