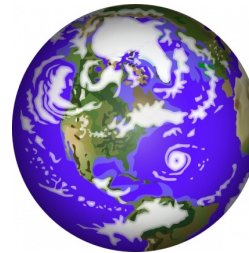# Planning to Control
# Crowd-Sourced Workflows

Daniel S. Weld
University of Washington

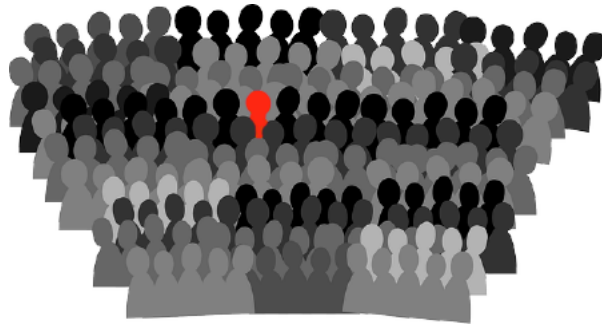**UW CrowdLab**

dub

---

## 30,000' View

- Crowdsourcing is huge & growing rapidly
  - Virtual organizations
  - Flash teams with mixed human & machine members

- Automatic organization of work
  - Reduce labor required by 30-85%

**UW CrowdLab**

dub

# Crowdsourcing

- Performing work by **soliciting effort** from many people
- **Combining the efforts** of volunteers/part-time workers (each contributing a small portion) to produce a large or significant result



# Crowdsourcing Successes



190 M reviews of 4.4 M businesses

Answers to 7.1 M prog. questions

Universal reference for anything

# Citizen Science

800,000 volunteers – Hubble images
Discovered "Hanny's Voorwerp" black-hole
"Pea galaxies"

Crowdsourced bird count & identification
Migration shift -> effect of climate change

Game to find 3D structure of proteins.
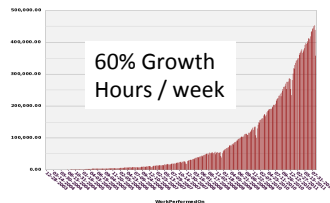Solved 15 year outstanding AIDS puzzle

5

# Labor Marketplaces
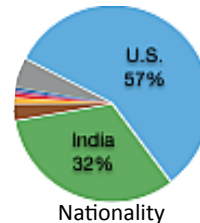Will Grow to $5B by 2018 [Staffing Industry Analysts]

- 2.7 million workers
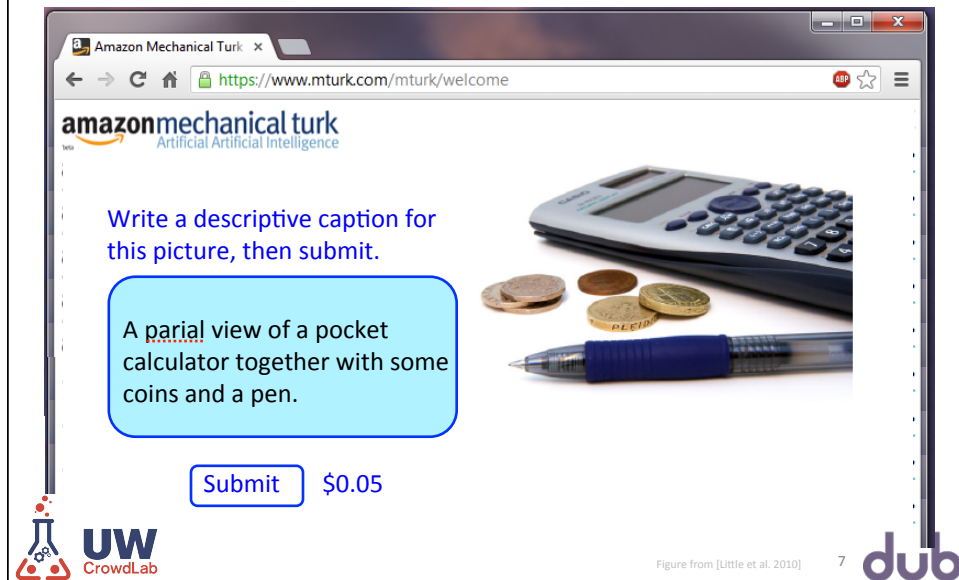- 540,000 requestors
- 35M hours worked in 2012

60% Growth
Hours / week

Charts from Panos Ipierotis' blog; phone from pixabay

U.S.
57%

India
32%

Nationality

# Example Job on Mechanical Turk



Figure from [Little et al. 2010]    7

# Big Work from Micro-Contributions

- Challenges
  - Small work units
  - Reliability & skill of individual workers vary

- Therefore
  - Use a *workflow* to aggregate results & ensure quality
  - Manage workers with (unreliable) workers

8

## *Ex:* Iterative Improvement

initial
caption

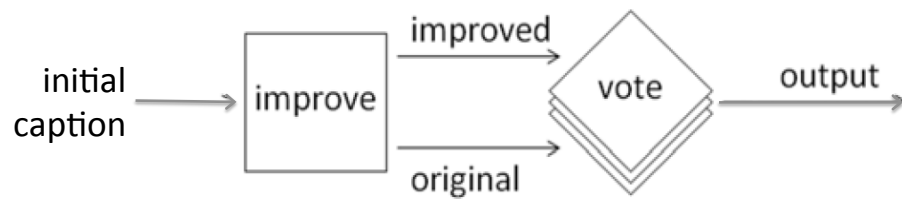[Little et al, 2010]   9

## *Ex:* Iterative Improvement

initial
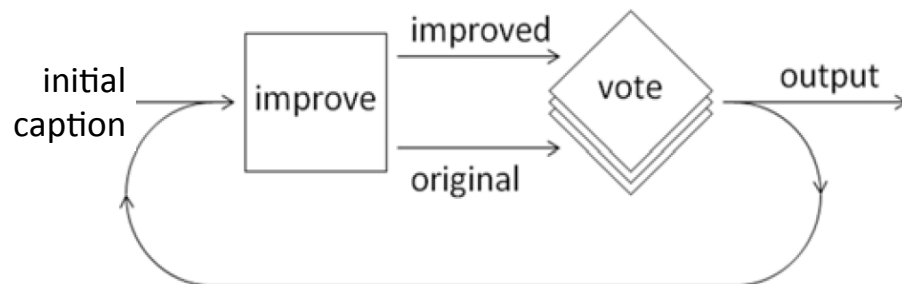caption  →  improve

[Little et al, 2010]   10

## *Ex:* Iterative Improvement



[Little et al, 2010]  11
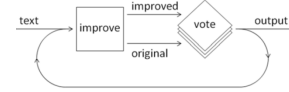
## *Ex:* Iterative Improvement



[Little et al, 2010]  12

# Iterative Improvement
[Little et al, 2010]

text → improve → improved / original → vote → output

### First version

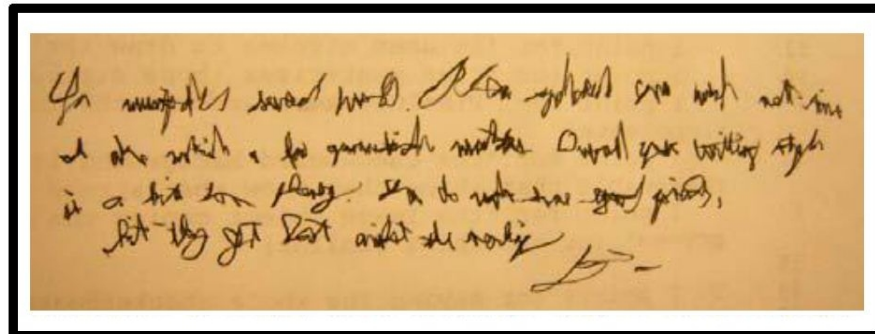A parial view of a pocket calculator together with some coins and a pen.

### After 8 iterations

A CASIO multi-function, solar powered scientific calculator.

A blue ball point pen with a blue rubber grip and the tip extended.

Six British coins; two of £1 value, three of 20p value and one of 1p value.

Seems to be a theme illustration for a brochure or document cover treating finance - probably personal finance.

UW CrowdLab

Figure from [Little et al. 2010]
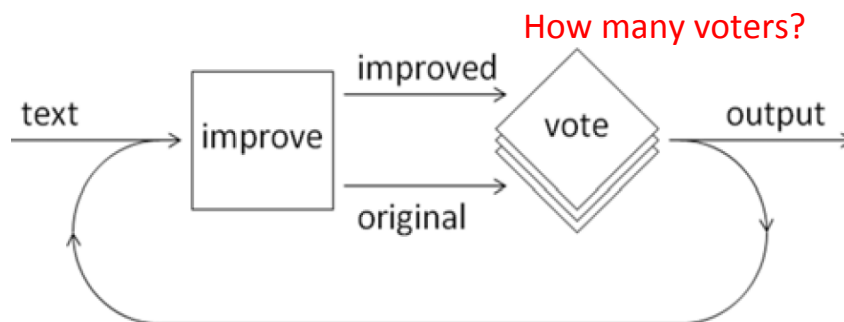
13 dub

---

UW CrowdLab

[Little et al, 2010]   14 dub

"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.

[Little et al, 2010]

# *Workflow Control Problem*



How many voters?

How many times?

Adaptive, Decision-Theoretic Control

16

4/14/15

# Outline

✓ Introduction
- **Case Study: Controlling Iterative Improvement**
- Case Study: Controlling Taxonomy Generation
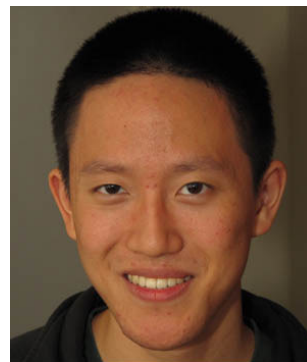- In Progress: Controlling ML Annotation

17

# TurKontrol
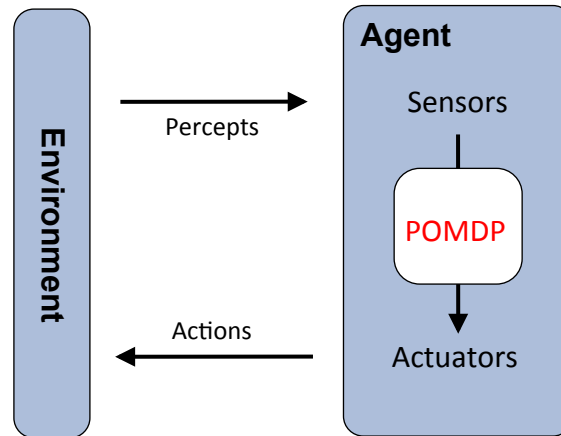### POMDP Control of Iterative Improvement

**Peng Dai**              **Chris Lin**
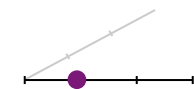
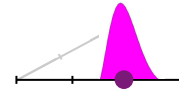Both co-advised with Mausam

# Artificial Intelligence 101

**Environment**

Percepts →

**Agent**

Sensors

POMDP

Actuators

← Actions

UW CrowdLab

19

dub

---

# Markov Decision Process

**Input:**

**World State**
s = <x, y>

**Actions**
P(s' | s, a)
Cost c

Observe: Next State s' = <x', y'>
Reward = f(s, a, s')

**Output:**

Construct **policy**, π : S→A, that chooses best action for each state
I.e., actions that **maximize expected reward – costs** over time

While **learning** action & reward probabilities
(Reinforcement learning)

UW CrowdLab

20

dub

## Partially-Observable
## Markov Decision Process

**Input:**

**Belief** State
P(s)

**Actions**
P(s' | s, a)
Cost c

Observe: Noisy Sensor = f(s')
Reward

**Output:**

Construct *policy*, π : S→A, that chooses best action for each state
I.e., actions that *maximize expected reward − costs* over time

While *learning* action & reward probabilities
(Reinforcement learning)

Figure from Dan Klein & Pieter Abbeel - UC Berkeley CS188: http://ai.berkeley.edu ]

21

---

# Solving the POMDP

**Constructing the policy**, π, to choose the best action

- Many algorithms
  - Point-based methods
  - UCT on discretized space
  - Lookahead search with beta distribution belief states

$$Q^*(s, a) = \sum_{s'} P(s' \mid s, a) \, [\, R(s, a, s') + \gamma \, \text{Max}_a \, Q^*(s, a) \,]$$

- Exploration / exploitation problem
  - ε-greedy
  - UCB / Multi-armed bandit

22

## From To

(Hidden)

| | | |
|---|---|---|
| **World State** | <x,y> coords | Quality $Q_1$, $Q_2 \in (0,1)$ |
| **Actions** | Move<br>Grasp | Improve caption task<br>Vote best caption |
| **Costs** | Power used | $$ paid to workers |
| **Reward** | | F(quality returned) |

Robot figure from Dan Klein & Pieter Abbeel - UC Berkeley CS188: http://ai.berkeley.edu ]

23

---

# Belief State



P   Quality$_{\alpha 1}$          P   Quality$_{\alpha 2}$

28

12

# Transition Model of Voting Action

Learned using Expectation Maximization

P    Quality$_{\alpha 1}$

P    Quality$_{\alpha 2}$

Worker votes that artifact 1 is better

P    Quality$_{\alpha 1}$

P    Quality$_{\alpha 2}$

UW CrowdLab

dub

29

# POMDP for Iterative Improvement

submit

initial artifact (α)

N

need improving ?

Y

$\alpha$   generate improve job

$\alpha'$   estimate quality of $\alpha'$

more voting ?

Y

N

$\alpha \mid \alpha'$   make ballot job

$\alpha \mid \alpha'$   update quality estimates

better of α and α'

UW CrowdLab

dub

35

4/14/15

# POMDP for Iterative Improvement



# POMDP for Iterative Improvement



14

# POMDP for Iterative Improvement



# POMDP for Iterative Improvement

POMDP for Iterative Improvement



POMDP for Iterative Improvement
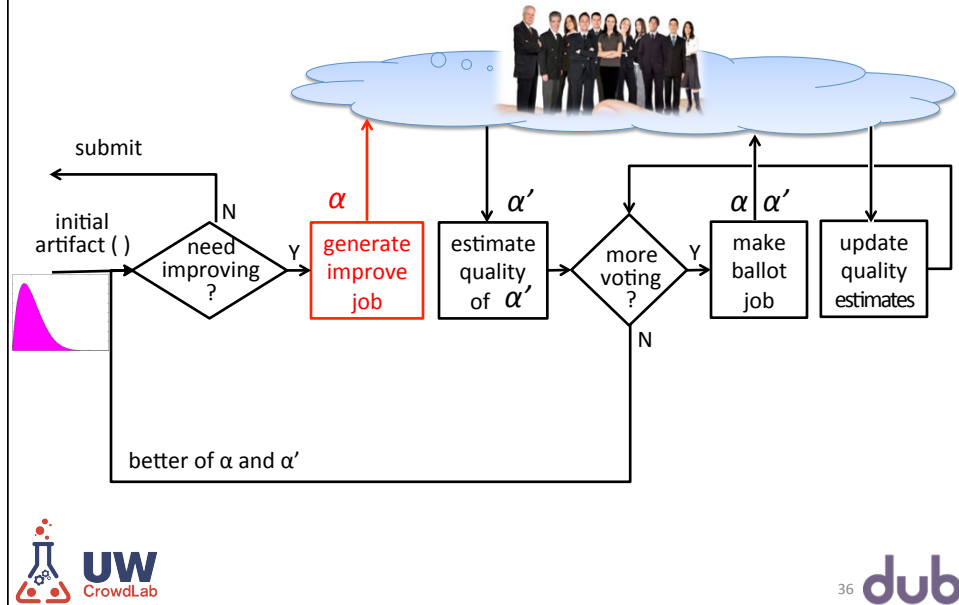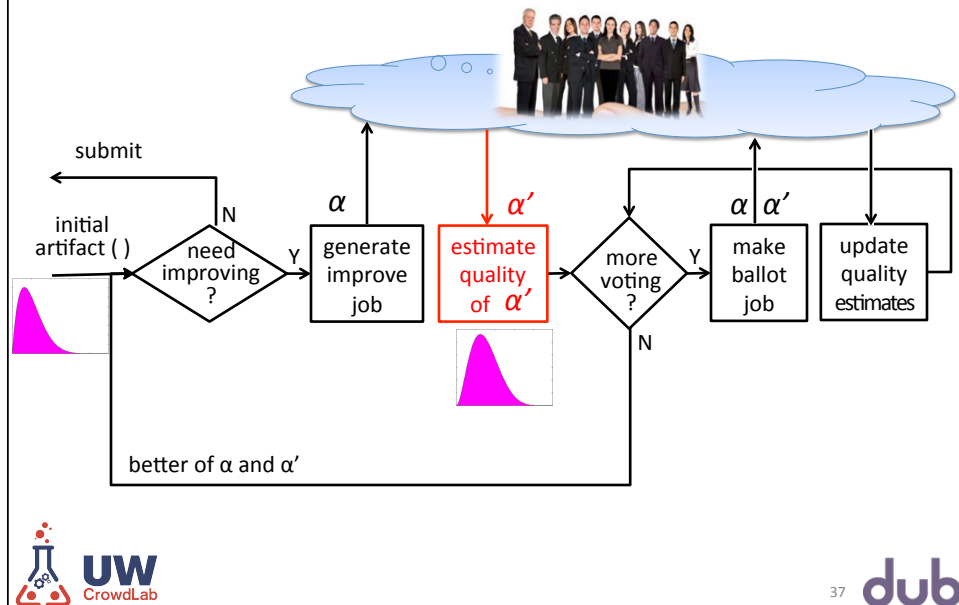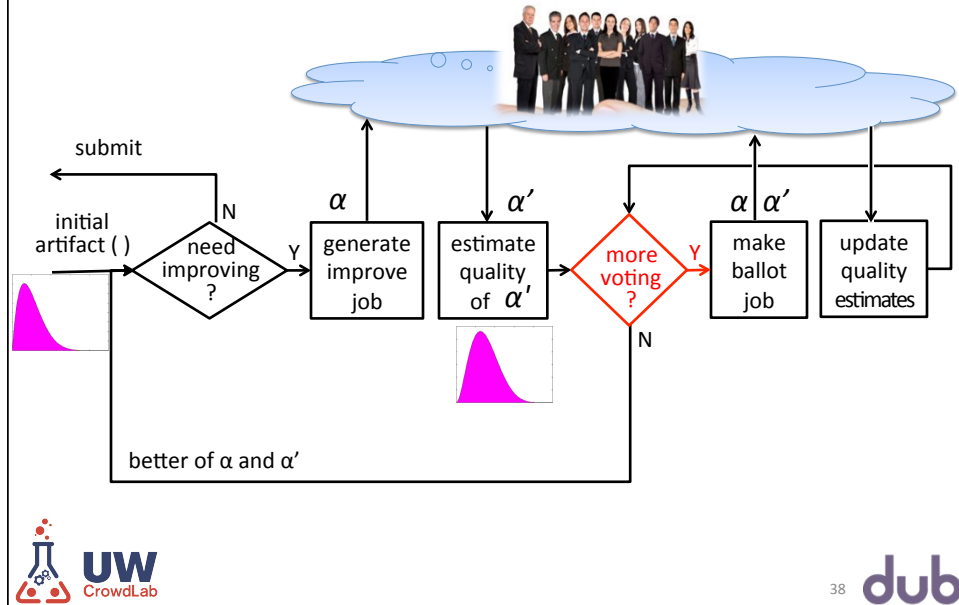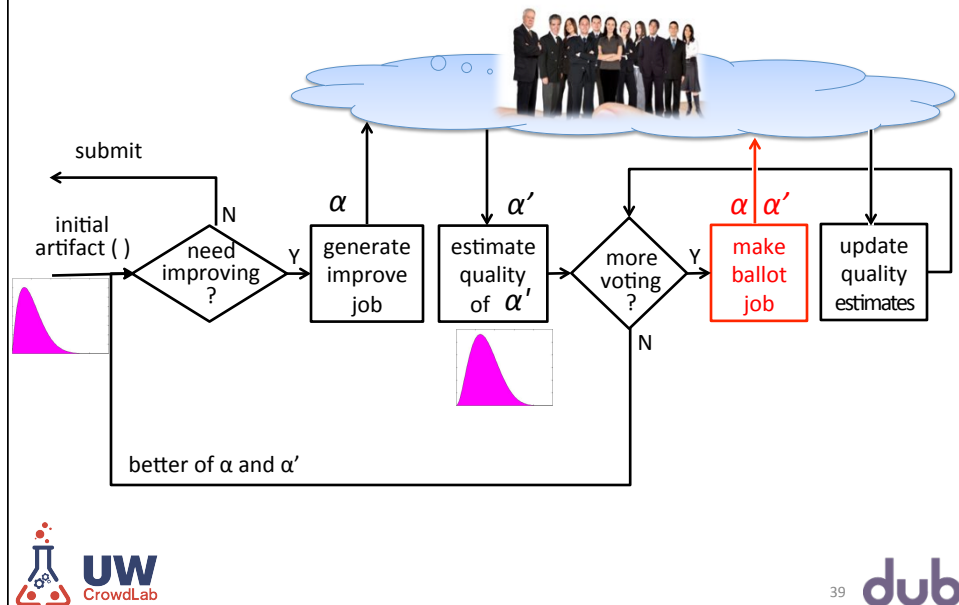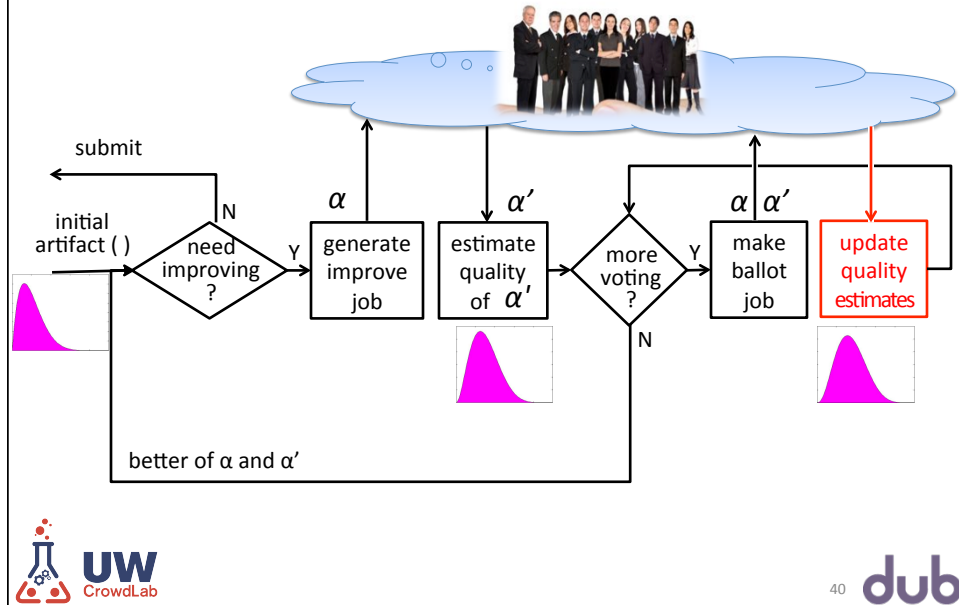
# POMDP for Iterative Improvement



# POMDP for Iterative Improvement

# POMDP for Iterative Improvement
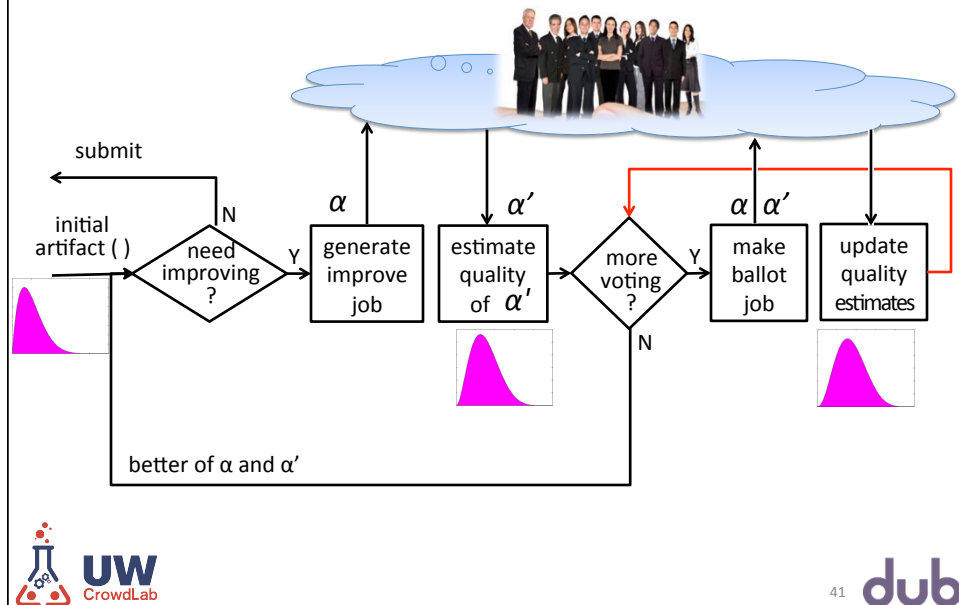
submit

initial
artifact ( )

N
need
improving
?

Y

$\alpha$ generate improve job

$\alpha'$ estimate quality of $\alpha'$

more voting ?

Y

N

$\alpha \mid \alpha'$ make ballot job

update quality estimates

better of α and α'

UW CrowdLab

48

dub

# POMDP for Iterative Improvement

submit α

initial
artifact ( )

N
need
improving
?

Y

$\alpha$ generate improve job

$\alpha'$ estimate quality of $\alpha'$

more voting ?

Y

N

$\alpha \mid \alpha'$ make ballot job

update quality estimates

better of α and α'

UW CrowdLab

49

dub

# Comparison

Quality

POMDP  Hand Coded

0.8
0.75
0.7
0.65
0.6
0.55
0.5

40 images, same average cost

Controlling quality: *POMDP 30% less labor*

[Dai, Mausam & W, AAAI-11]
[Dai *et al.* AIJ 2013]

50

# Allocation of Human Labor

Average # of Ballots

POMDP (TurKontrol)
Hand Coded

7
6
5
4
3
2
1
0

1   2   3   4   5   6   7   8
**Iteration Number**

51

# Human Labor Redirected

# Lessons So Far

- Reduced labor costs
- Improved quality

- POMDP planning
  - Update belief states about uncertain world
  - Model sensing actions
- Expectation maximization & prob inference

# Outline

✓ Introduction

✓ Case Study: Controlling Iterative Improvement

- **Case Study: Controlling Taxonomy Generation**

- In Progress: Controlling ML Annotation

54

---

# Cascade
**Crowdsourcing Taxonomy Creation**

**Lydia Chilton**
Co-advised with James Landay

# Image Data Sets

# Q&A Site Responses

# Crowdsourcing Taxonomy Generation
## *Is Hard !*

- Good taxonomy requires a global perspective

- But workers see only a tiny fraction of data…?

**UW** CrowdLab          58    dub

---

# Iterative Improvement?

**Task: Add Tips to the Hierarchy of Travel Advice**

```
* Packing/packing
** Clothing

* food

* flying
** carry-on luggage
** in flight meals/airline food
** airport check in/Check-In
** On Board Entertainment
** flying on a budget
** customs/Immigration

* luggage

* Communication/communication
** long distance calls

* Insurance

* Security Control/security

* Accomodation/lodging
** Hostels/the benefits of hostels

* Budget Travel/Thrifty Travel Tips
** flying on a budget

* personal valuables
* San Fransisco International
* tours
* Travel Etiquette
```

`<` Previous 4     **Tips that still need categorization**     `>` Next 4

#100
"-Luggage with a lifetime guarantee is worth the slight premium in price. Briggs and Riley make a very sturdy bag that's strong enough you can sit on it during a long pre-boarding wait, and with zippers that rarely break. And when they do - in 5 or 10 year"

#101
"-If you're tall or otherwise picky about airplane seats, use seatguru.com to understand the seat layout of your flight. Seatguru will warn you about equipment boxes under the seat in front of you, cold seats, or seats with a lot of bathroom traffic."

#102

#103

`Submit`

**Problems**
1. The growing hierarchy becomes overwhelming
2. Workers confused

**Lesson**: Decompose the task into smaller steps

**UW** CrowdLab          59    dub

25

# Initial Approach 2:
# Category Comparison



**Problem**
Without context it's hard to judge relationships:
- **TSA liquids** *vs.* **removing liquids**
- **Packing** *vs.* **what to bring**

**Lesson**: Don't compare abstractions

UW CrowdLab

60

dub

---

# Cascade Overview

[Chilton *et al.*, CHI-13]

Use the *crowd* to:

1. Generate category names
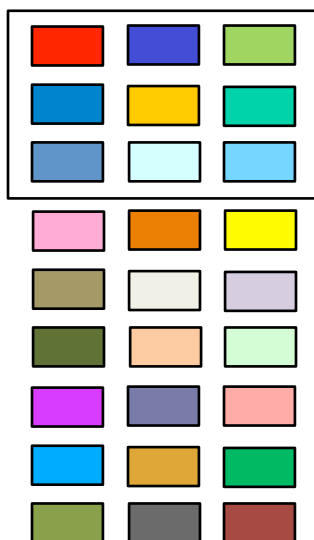2. Select the best categories
3. Place the data into the best categories

Use *machines* to:

4. Infer global structure of categories

UW CrowdLab

61

dub

Example Input: 100 Random Colors
Step 0. Sample Data
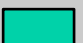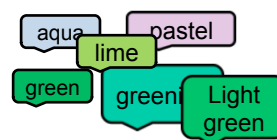


Step 1. Generate Categories

**For each color** **Task** **Crowd responses**

What category do you suggest for this color?

greenish

aqua  pastel  lime  green  greeni  Light green

**This generates an initial set of category names.**

# Step 2. Select Best Categories

**For each color**

**Task**

What is the best category for this color?

| Category | Best? |
|----------|-------|
| Aqua | ☐ |
| Greenish | ☑ |
| Lime | ☐ |
| Pastel | ☐ |

**Crowd responses**

| Category | Votes |
|----------|-------|
| Aqua | 1/5 |
| Greenish | 4/5 |
| Lime | 0/5 |
| Pastel | 0/5 |

**An early filter for spam and vague categories**

UW CrowdLab

64

dub

---

# Step 3. Label Data

**For each color and category**

| Categories |
|------------|
| Green |
| Greenish |
| Yellow |
| Pink |

**Task**

What categories does this belong to?

| Category | Fits | Doesn't Fit |
|----------|------|-------------|
| Green | ☑ | ☐ |
| Greenish | ☑ | ☐ |
| Yellow | ☐ | ☑ |
| Pink | ☐ | ☑ |

**Crowd responses**

| Category | Votes |
|----------|-------|
| Green | 4/5 |
| Greenish | 5/5 |
| Yellow | 1/5 |
| Pink | 0/5 |

**This determines category membership.**

UW CrowdLab

65

dub

# Step 4. Global Structure Inference



**Determine parent/child relations; eliminate duplicates.**

66

# Finally, … Recurse



**May lead to new tags & recomputing taxonomy**

67

29

## Evaluation

**Quality**

inter-annotator agreement

- Cascade vs. 4 Experts
- Experts vs. 3 other experts

**Decision-Theoretic Control!**

**Cost**

- Cascade
- Expert

69

---

## *Deluge*
## (Decision-Theoretic Control of Cascade)

## Jonathan Bragg
Co-advised with Mausam

70

30

# Why is Cascade Expensive?

# POMDP Model Agent Belief State

*World state* = taxonomy & labels applying to item

# POMDP Model Agent Belief State

***Belief state*** includes … distribution over taxonomies
label probabilities for item



*Etc.*

7%            5%            …

Learn & refine taxonomy during execution
Too complex for off-the-shelf POMDP solver

74

# Decision Cycle for New Item

[Bragg, Mausam & W HCOMP-13]

**Agent**

Sensors ← Percepts

**Probabilistic inference to update**
• Posterior probabilities
• Co-occurrence model for labels
• Worker accuracy

?



Independent        Joint (naive Bayes)        Joint (MRF / CRF)

Actuators → Actions    Ask about label with max VOI

76

# Performance of Decision-Theoretic Model



# Performance of Decision-Theoretic Model

## Now crowd is cheaper than experts!

Reaches same performance as Cascade
with only *13%* as many voting jobs

# Lessons So Far

- Decision-theoretic planning
  - Probabilistic inference
  - Expectation maximization

- Reduced labor & improved quality
  - Iterative Improvement
  - Taxonomy Generation
  - *???*

79

---

# Outline

✓ Introduction
✓ Case Study: Controlling Iterative Improvement
✓ Case Study: Controlling Taxonomy Generation
- In Progress: Controlling ML Annotation

80

4/14/15

# Information Omnivore Project

- **Large Scale Information Extraction**



- **Train via 2 kinds of Weak Supervision**
  - **Align Corpus to Background Knowledge Base**
    [Wu & W CIKM-07; … Koch *et al.* EMNLP-14]
  - **Identify & Extract Events from Newswire**
    [Zhang & W EMNLP-13; Zhang, Soderland & W TACL-15]

NewsSpike

81

---

# Information Omnivore Project

- **Augment with Crowdsourced Annotations**
  - Eg:
    "Calling himself Guccifer, Marcel-Lehel Lazar rampaged through the email accounts of rich and powerful Americans…"

    AliasOf(p, p)
  - **For improved machine learning performance**

- **Train via Semi-Distant Supervision**
  - **Align Corpus to Background Knowledge Base**
    [Wu & W CIKM-07; … Koch *et al.* EMNLP-14]
  - **Identify & Extract Events from Newswire**
    [Zhang & W EMNLP-13; Zhang, Soderland & W TACL-14]

NewsSpike

82

35

# Observation

- Vast proportion of micro-task crowdsourcing… is used to create training data for ML classifiers
  - Chris Caliston-Burch (UPenn) $250,000 on MTurk
  - LDC: 44 FT employees *just* creating NLP training data
  - Google, MSFT – internal CS: each larger than MTurk

- Common approach
  - Get two humans to annotate
  - If they agree, … done
  - Else recruit a third to arbitrate
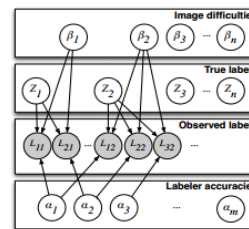
*"2/3 Relabeling"*
*j/k Relabeling*

UW CrowdLab

83 dub

---

$$U(p(\beta_r|y_{i'},l')) = \left\| E_{p(\beta_r)}(\beta_r) - E_{p(\beta_r|y_{i'},l')}(\beta_r) \right\|_2 \quad (12)$$

$$\approx \left\| E\left( \frac{1}{S-1} \sum_{s=2}^{S} Z_r^{s-1^\top} \left[ (\gamma|\gamma^{s-1}, Z^{s-1}) - (\gamma_{(i',l')}|\gamma^{s-1}, Z^{s-1}) \right] \right) \right\|_2. \quad (13)$$

$$
\begin{aligned}
Q(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= E\left[ \ln p(\mathbf{l}, \mathbf{z}|\boldsymbol{\alpha}, \boldsymbol{\beta}) \right] \\
&= E\left[ \ln \prod_j \left( p(z_j) \prod_i p(l_{ij}|z_j, \alpha_i, \beta_j) \right) \right] \\
&\text{since } l_{ij} \text{ are cond. indep. given } \mathbf{z}, \boldsymbol{\alpha}, \boldsymbol{\beta} \\
&= \sum_j E\left[ \ln p(z_j) \right] + \sum_{ij} E\left[ \ln p(l_{ij}|z_j, \alpha_i, \beta_j) \right]
\end{aligned}
$$

$$p(z|L, \theta) = \int p(z, q|L, \theta)dq = \prod_{j \in [M]} \int_0^1 p(q_j|\theta) q_j^{c_j} (1-q_j)^{\gamma_j - c_j} dq_j \overset{def}{=} \prod_{j \in [M]} \psi_j(z_{\mathcal{N}_j}), \quad (4)$$

[Dawid *et al* 79, Whitehill *et al* 09, Welinder *et al* 10, Raykar *et al* 10, Karger *et al* 11, Kajino *et al* 12, Baba *et al* 13, Liu *et al* 12, etc, etc…]
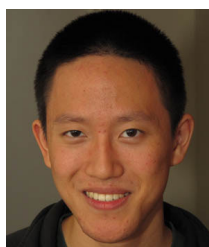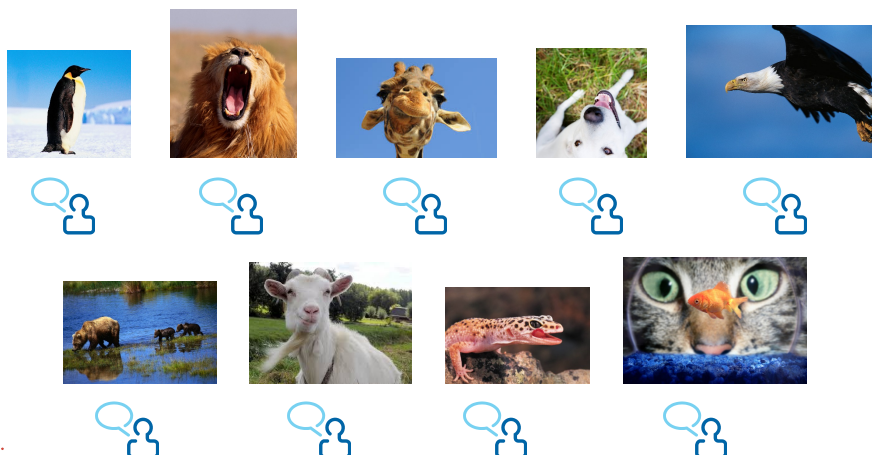
UW CrowdLab

84 dub

[Lin, Mausam & W HCOMP-14]

How should one best spend a fixed annotation budget… *when training an ML classifier?*

85

---

# Unilabel?
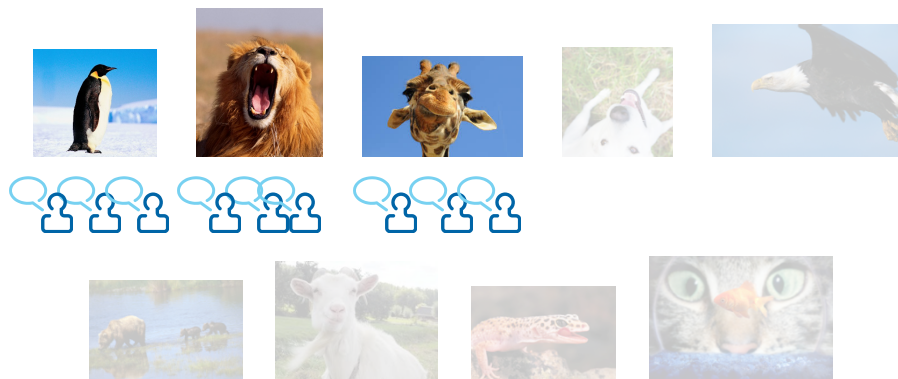
9 examples with labels that are 75% accurate?

# 2/3 Relabel?

3 examples each with 3 labels – consensus 84% accurate?



# Or Even?

1 example with 9 labels – consensus 98% accurate?

## Existing Data Sets?

| Dataset | # Features | # Examples |
|---|---|---|
| (a) Breast Cancer | 9 | 699 |
| (b) Bank Note Authentication | 4 | 1372 |
| (c) Seismic Bumps | 18 | 2584 |
| (d) EEG Eye State | 14 | 14980 |
| (e) Sonar | 60 | 208 |
| (f) Breast Cancer Diagnostic | 30 | 569 |
| (g) Hill-Valley | 100 | 606 |
| (h) Hill-Valley with Noise | 100 | 606 |
| (i) Internet Ads | 1558 | 2359 |
| (j) Gisette | 5000 | 6000 |
| (k) Farm Ads | 54877 | 4143 |
| (l) Spambase | 57 | 4601 |

## Factors that Affect Relabeling Efficacy

Inductive Bias of Classifier

"Strong" → limited expressiveness

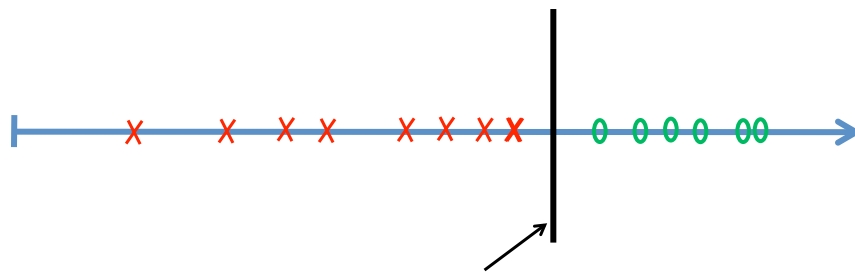"Weak" → can learn many different concepts

Worker Accuracy

Budget

---

# If Data was Clean



True Concept: 65 and older -> "Senior Citizen"

92

## With Noisy Annotation

True Concept: 65 and older -> "Senior Citizen"

93

## (low expressiveness)
## Strong Inductive Bias Classifier

True Concept: 65 and older -> "Senior Citizen"

# Overfitting to Noise

(high expressiveness)

Weak Inductive Bias Classifier

True Concept: 65 and older -> "Senior Citizen"

95

# Conjecture

- Relabeling more important for classifiers with weak inductive bias

  (e.g., in domains with myriad features)

96

# Experiments on Synthetic Data



Weaker Inductive Bias

Relabeling Better

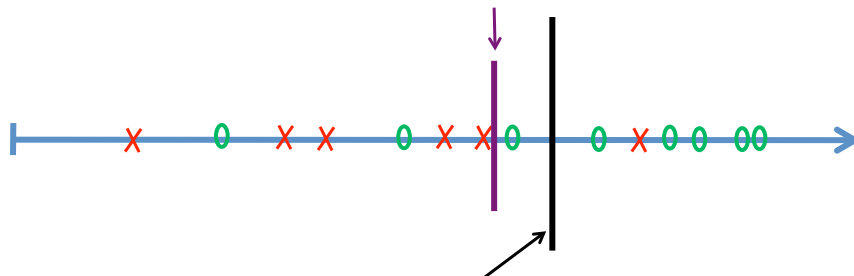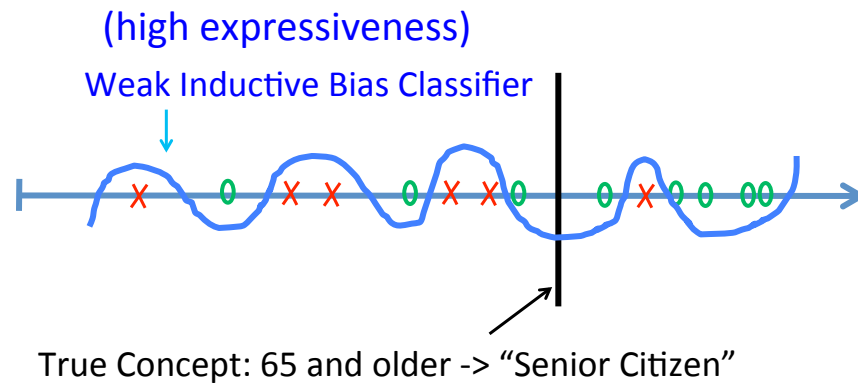Relabeling Accuracy / Unilabeling Accuracy

Number of Features (≈VC Dimension)

2/3 Relabeling
3/5 Relabeling
4/7 Relabeling

98

# Revisiting the Real Data



< 100 features    > 100 features

Relabeling Accuracy / Unilabeling Accuracy

b  a  d  c  f  l  e  g  h  i  j  k

Dataset

2/3 Relabeling
3/5 Relabeling
4/7 Relabeling

# Factors that Affect Relabeling Efficacy

Inductive Bias of Classifier

Worker Accuracy

# Accuracy of **Training Data**

4/14/15

# Factors that Affect Relabeling Efficacy

**Inductive Bias of Classifier**

**Worker Accuracy**



Results on simulated Gaussian data, fixed dimensionality = 50

---

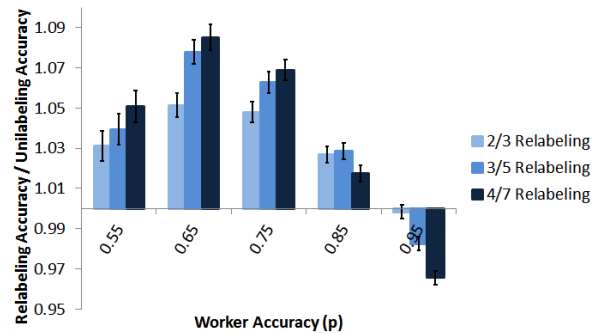# Factors that Affect Relabeling Efficacy

**Inductive Bias of Classifier**

**Worker Accuracy**

**Budget**

## Future Work

**Relax Assumptions**

**Complete Decision-Theoretic Control**

103

45

# Outline

✓ Introduction
✓ Case Study: Controlling Iterative Improvement
✓ Case Study: Controlling Taxonomy Generation
- Future Challenges

104

# Other Challenges

- Usually assume **workers choose** job to perform

amazonmechanical turk
Artificial Artificial Intelligence
beta

- What if employer can **assign** jobs to best workers?
  - Google internal crowdsourcing
    - Street-view/maps, knowledge graph, search relevance
    - Task routing (expert / novice) in citizen science

4/14/15

# Matching Jobs to Workers

- Set of jobs, each with difficulty
- Set of workers, each with
  - Skill
  - Capacity (bound on # jobs)
  - Independent errors (conditioned on difficulty)
- Minimize overall error *wrt* fixed budget

- Knapsack?
  - "Pack" jobs with workers



106

# Unknown Difficulty ∨ **Skill**

Jonathan Bragg    Andrey Kolobov

- If skill levels are known…
  - Assigning unknown problem is like MAB "arm"
  - Once find hard problem (workers disagree), add expert
- If difficulty is known…
  - Assigning unknown worker is like MAB "arm"

- Exploration / Exploitation Tradeoff
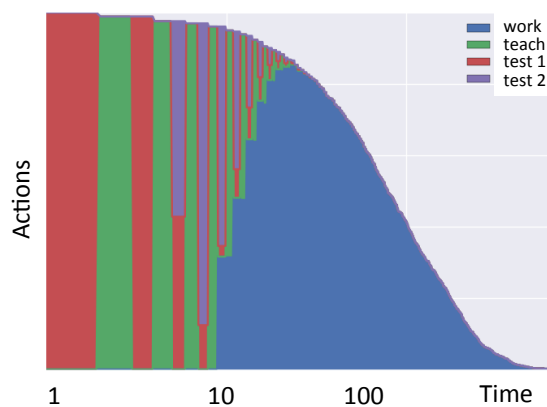  - [Bragg, Kolobov, Mausam & W  HCOMP-14]

107

47

# Additional Challenges

- Balancing worker desires w/ central needs
  - Frenzy [Chilton et al. CHI-14]

---

# Additional Challenges

- Balancing worker desires w/ central needs
- Optimizing for time
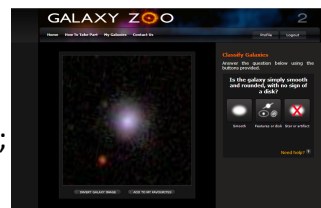- Balancing **work**, **teaching** & **testing**



Jonathan Bragg

109

# Additional Challenges

- Balancing worker desires w/ central needs
- Optimizing for time
- Interleaving work, education & testing
- Workers improving job instructions
- Aggregation when majority is wrong
  - Bayesian truth serum
  - MicroTalk – focused argumentation

---

# Related Work

- DT Crowdsourcing / Active Learning with Noise
  - GalaxyZoo – [Kamar & Horvitz 2012]
  - BBMC – [Wauthier & Jordan 2011]
  - ITS – Poppovic & Brunskill
  - [Sheng et al. 2008, Donmez et al. 2009;
  - Etc.
- Crowdsourcing Global Structure
  - Mobi – [Zhang *et al.* 2012]
  - Context Trees - [Verroios & Bernstein 2014]
- Information Omnivore
  - Never-Ending Language Learning – [Carlson et al 2012]
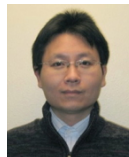  - [Angeli & Manning 2014, Pershina *et al.* 2014]

# Conclusion

- **Crowdsourcing is huge & growing rapidly**
  - Specialized communities, citizen science & labor mkts

- **Decision theoretic planning – large potential**
  - Reduce required labor by 30-85%
  - Sequential decision making is crucial
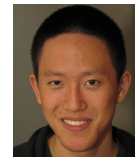  - Must model uncertainty & noisy sensors

- **Many open questions...**

---

| | | | | | |
|---|---|---|---|---|---|
| Jonathan Bragg | Lydia Chilton | Peng Dai | Shih-Wen Huang | James Landay | Chris Lin |

Thanks

| | | | |
|---|---|---|---|
| Angli Liu | Andrey Kolobov | Mausam | Stephen Soderland |

# Extra Slides