

Crowd Computing

Rob Miller User Interface Design Group MIT CSAIL

Joint work with Greg Little, Lydia Chilton, Max Goldman, Jeff Bigham, Michael Bernstein, David Karger, Mark Ackerman, Björn Hartmann, Joel Brandt, Mason Tang, Elena Tatarchenko





motodes sweet web the golden and some with not in plang in to which you good girld,



Crowd Computing



MIT HUMAN-COMPUTER II

Coordinating a crowd (a large group of people on the web) to do micro-work (small contributions) that solves problems (that software or one user alone can't do)



Amazon Mechanical Turk



highly-available, short, cheap, programmable...
 a prototyping platform for crowd computing

amazonn	tificial Artificia	cal turk	Your A	ccount	HITs	Qualifications	90,071 HITs available now	Rob Miller	
All HITs HITs Available To You HITs Assigned To You									
Search for	HITs	🚽 contai	ning 🔤 th	at pay at	least \$ <mark>0.00</mark>	for which	you are qualifie	d 🗉 😡	
All HITs 1-10 of 1775 Results									
Sort by: HIT Creation Date (newest first) + O Show all details Hide all details 1 2 3 4 5 > Next >> Last									
Find the link and image that a forum thread is discussing View a HIT in this group									
Requester:	AEI	HIT Expira	tion Date:	Jul 12, 20	10 (59 minute	es 35 seconds) 🖪	leward:	\$0.05	
		Time Allot	ed:	8 minutes		H	IITs Available:	1	
Write a short answer to a question about writing View a HIT in this group									
Requester:	conjecture co	rporation	IT Expirati	on Date:	Jul 26, 2010	(1 week 6 days)	Reward:	\$0.04	
		1	ime Allotte	d:	60 minutes		HITs Available:	20	

Why We Need Crowd Algorithms





MIT HUMAN-COMPUTER IN

Tackling Hard Tasks with a Crowd



MIT HUMAN-COMPUT

1 work Rhan yohans swed y are troiting theh

- Problem: one person can't do a good transcription
- Key idea: iterative improvement by many workers



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

Improve-and-Vote Algorithm





MIT HUMAN-COMPUTER INTERACTION



You musight several hand Room galated some with notions at also which a few ganutische mosthe Ownall spik withing style is a tite for plang. In to refer has good good, At the get last airfut de nortige

After 9 iterations

"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.



Another Example: Blurry Text



i have been find to bet the sumi, but (or find a very gravel with it ensures, and i student up hitting, my theorem, and i contracted statements, i know it is a contracted statement, i know it is a contracted statement of the second statement is any index of the second statement is any index of the second statement is in a contract statement is any index of the second statement is in a contract statement is in a contract statement is any index of the second statement is in a contract statement is in

After 8 iterations

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 occurence 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.



Human Computation Algorithms

 Improve-and-Vote is a simple but effective crowd algorithm



• Considerations for crowd algorithm design:

Quality

- Improve-and-Vote is useful on a noisy crowd like MTurk (~30% of open-ended work is bad somehow, and ~3% of workers are spammers)
- Other crowds may have less noise

Time

• Improve-and-vote runs slowly, because it's serialized

Incentives

- For MTurk, equivalent to cost in \$
- Other crowds need other incentives



MIT HUMAN-COMPUTER

11

8

7

TurKontrol: Optimizing Improve-and-Vote

- Maintain estimates of process state ٠
 - artifact quality
 - worker ability
 - voter accuracy
- Make utility-maximizing decisions ٠
 - when to get another voter
 - when to stop iterating

Peng Dai, Mausam, Daniel S. Weld. "Artificial Intelligence for Artificial Artificial Intelligence." AAAI 2011.













Shortening A Paper to Ten Pages



Shortn

This paper presents Soylent, a word processing interface that uses crowd workers to help with proofreading, document shortening, editing and commenting tasks. Soylent is an example of a new kind of interactive user interface in which the end user has direct access to a crowd of workers for assistance with tasks that require human attention and common sense. Implementing these kinds of interfaces requires new programming patterns for interface software, since crowds behave differently than computer systems. We have introduced one important pattern, Find-Fix-Verify, which splits complex editing tasks into a series of identification, generation, and verification stages that use independent agreement and voting to produce reliable results. We evaluated Soylent with a range of editing tasks, finding and correcting 82% of grammar errors when combined with automatic checking, shortening text to approximately 85% of original length per iteration, and executing a variety of human macros successfully.

Future work falls in three categories. First are new crowd-driven features for word processing, such as readability analysis, smart find-and-replace (so that renaming "Michael" to "Michelle" also changes "he" to "she"), and figure or citation number checking. Second are new techniques for optimizing crowd-programmed algorithms to reduce wait time and cost. Finally, we believe that our research points the way toward integrating on-demand crowd work into other authoring interfaces, particularly in creative domains like image editing and programming.

This paper presents Soylent, a word processing interface that uses crowd workers to help with proofreading, document shortening, editing and commenting tasks. Soylent is an example of a new kind of interactive user interface in which the end user has direct access to a crowd of workers for assistance with tasks that require human attention and common sense. Implementing these kinds of interfaces requires new programming patterns for interface software, since crowds behave differently than computer systems. We have introduced one important pattern, Find-Fix-Verify, which splits complex editing tasks into a series of identification, generation, and verification stages that use independent agreement and voting to produce reliable results. We evaluated Soylent with a range of editing tasks, finding and correcting 82% of grammar errors when combined with automatic checking, shortening text to approximately 85% of original length per iteration, and executing a variety of human macros successfully.

--- C -- X

_ 0 _ X

sents Soylent, a word

rface that uses crowd

Shortn

00

* X

\$1.30

Future work falls in three categories. First are new crowd-driven features for word processing, such as readability analysis, smart find-and-replace_(so that renaming "Michael" to "Michelle" also changes "he" to "she"), and figure or citation number checking. Second are new techniques for optimizing crowd-programmed algorithms to reduce wait time and cost. Finally, we believe that our research points the way toward integrating on-demand crowd work into other authoring interfaces, particularly in creative domains like image editing and programming.

Amagon Mechanical Turk, olicki '10, ACM Pres for-hire, so results belong to the requester. Likewise with Clarke, J. and Lanata, M. Models for sentence compresdiate access to a pool of human expertise. Lag times in our historical precedent traditional copyeditors do not own 22. Sala, M., Partridge, K., Jacobson, L., and Begole, J. An Exploration into Activity-Informed Physical Advertission: a comparison across domains, training requirecurrent implementation are still on the order of minutes to their edits to an article. However, crowdsousced interfaces ments and evaluation measures. ACL '06, Association hours, due to worker demographics, worker availability, the will need to consider legal questions carefully ing Using PEST. Pervasive '07, Springer Berlin Heidelrelative attractiveness of our tasks, and so on. While future growth in crowdsourced work will likely shorten lag times, for Computational Linguistics (2006). 5. Cohn, T. and Lapata, M. Sentence compression beyond berg (2007). A final concern is that anonymous workers may not have this is an important avenue of future work. It may be possithe necessary domain knowledge or enough thated context to usefully contribute. We agree that some tasks, like fleshword deletion. COLLING '08, (2008). 23. Simon, I., Morris, D., and Basu, S. MySong: automatic ble to explicitly engineer for responsiveness in return for accompaniment generation for vocal melodies. Proc. 6. Cypher, A. Watek What J Do. MIT Press, Cambridge, higher monetary investment, or to keep workers on retainer ing out a related work section in an academic paper based CHU AS ACM Press (2008) MA, 1993. with distractor tasks until needed [3] on bullet points, are much more difficult to achieve on to-day's Mechanical Turk. However, a large subset of editing tasks only requires generic editing skills. We also may ef-

 Snow, R., O'Connor, B., Jucafisky, D., and Ng, A.Y. Cheap and fast—but is it good?, evaluating non-expert annotations for natural language tasks. *ACL* '03, (2005). Sorokin, A. and Forsyth, D. Utility data annotation with Amazon Mechanical Turk. CVPR V8, (2008).
 son Ahn, L. and Dathigh, L. Labeling images with a computer game. CHP V64, ACM Press (2004).



W 🔒 🤊 -

A

L

Crowdproof Sho

-

Ho

Michael Bernstein et al. "Soylent: a Word Processor with a Crowd Inside." UIST 2010. Best student paper.



*

Find-Fix-Verify Pattern



Find

"Identify at least one area that can be shortened without changing the meaning of the paragraph."

Keep patches found by at least two people

Fix

"Edit the highlighted section to shorten its length without changing the meaning of the paragraph."

the loss of a party		-	-
	 		Million.

Randomize order of suggestions



"Choose at least one rewrite that has significant style errors in it.
Choose at least one rewrite that Soylent is a prototypes...
Soylent is a prototypetest...
Soylent is a prototypetest...



Shortn Performance



• Length

15% shorter on average (10-22% overall)

• Cost

\$1.41 per paragraph
\$0.55 to Find an average of two patches
\$0.48 to Fix each patch
\$0.38 to Verify the results

• Time

Wait: median 18.5 min ($Q_1 = 8.3 \text{ min}$, $Q_3 = 41.6 \text{ min}$) Work: median 2.0 min ($Q_1 = 60 \text{ sec}$, $Q_3 = 3.6 \text{ min}$)





Find-Fix-Verify in Soylent

Both Shortn and Crowdproof use the Find-Fix-Verify pattern. We will use Shortn as an illustrative example. To provide the user with near-continuous control of paragraph length, Shortn should produce many alternative rewrites without changing the meaning of the original text or introduce⁷ grammatical errors.

⁷ Word's grammar checker, eight authors and six reviewers did not catch the error in this sentence. Crowdproof later did, and correctly suggested that "introduce" should be "introducing".



VizWiz: Helping the Blind See







Jeffrey Bigham *et al.* "VizWiz: Nearly Real-time Answers to Visual Questions." UIST 2010. Best paper award.



Helping the Blind See









Deployment



10,000+ downloads, dozens of questions/day





Adrenaline: A Crowd-Powered Camera







Mechanical Turk



Ten second video Time to final picture:



Michael Bernstein, Joel Brandt, et al. "Crowds in Two Seconds: Enabling Realtime Crowd-Powered Interfaces." UIST 2011.



Step 1: Get the Crowd Fast



- Retainer model
 - Recruit crowd in advance, and pay them to wait a few minutes



• A worker on retainer costs \$0.30 per hour



Step 2: Get the Crowd to Work Fast



- Rapid refinement pattern
 - Recognize potential agreement early, and use it to reduce the search space







Algorithm	Histogram of Executio	n Times Ν=72 12.6 sec, σ=2.2 sec
Rapid Refinement		22¢ 16.3 sec, σ=9.8 sec
First Answer		22¢ -5.3 sec, σ=14.0 sec
All Answers +Vote	•• .•-•	53¢



Caesar: Crowd-Driven Code Review



- Students in MIT software engineering classes write lots of code
 - roughly 10kloc in problem sets and projects
- Automatic grading is necessary but not sufficient



- we need human readers, and we want line-by-line feedback



Caesar: Crowd-Driven Code Review



- Chop up student programs into chunks
- Assign the chunks for review by a mixed crowd



• Results from 2 semesters in MIT 6.005 software engineering

13 problem sets, 2200 submissions

21,500 comments 5% alums 8% staff 87% students

16.2% upvoted 0.7% downvoted

9.6 comments per submission comments come back in < 3 days



Kinds of Comments









- alum: avoid abbreviating variable names... 'hi' would be especially confusing to someone who isn't familiar with english
- student: Idk about this one though. I've seen to and hi at various places. I'd say this is fine. Though, the integer N should be lower case, just to conform with the java naming convention.
- alum: When you're in industry and working on code with people who aren't necessarily familiar with English... it's a problem to use phonetic abbreviations. There is really no good reason to abbreviate variable names especially since IDEs autocomplete.
- student: For that, yes, I agree. Sadly, some IDEs don't autocomplete very well. *cough*
 - code author: These were the variables that were given to us















Deployable Wizard of Oz





- Wizard of Oz
 - tried-and-true prototyping technique when we don't know how to write the software part of a system
- Crowd computing enables Wizard of Oz systems that are useful and deployable
 - collect data from real use
 - use AI for performance/cost improvement



30

The User Plays a Key Role



- The end-user knows the goal and the context of the work
- The end-user provides coherence to a creative product
- The end-user has final responsibility for the result

Automatic clustering generally helps separate different kinds of records that need to be edited differently, but it storn perfect. Sometimes it creates more clusters than needed, because the differences in structure aren't migrature to be add near the end of each ine, then differences at the start of the line are targely irrelevant, and it isn't necessary to spit based on those differences. Conversely, sometimes the clustering isn't fine enough, leaving heterogeneous clusters that must be edited one line at a time. One solution to this problem would be to let the user rearrange the clustering manually, perhaps using drag-and-drop to merge and split clusters. Clustering and selecting generalization would also be improved by recognizing common text structure like URLs, filenames, email addresses, dates, times, etc. Automatic clustering generally helps separate different kinds of records that need to be edited differently, but it sn't perfect. Sometimes il creates more clusters than needed, because the differences in structure aren't relevant to a specific task. Conversely, sometimes the clustering isn't fine enough, leaving heterogeneous clusters that must be edited one line at a time. One solution to this problem would be to leit the user. Clustering this selection generalization would also be improved by recognizing common text structure like URLs, tilenames, email addresses, dates, times, etc.

Crowd



MIT HUMAN-COMPUTER INTERACTION

Many Crowds, Many Incentives





MIT HUMAN-COMPUTER INTERACT

Hints for Crowd System Design



- Divide work into small chunks
 - for parallelism and fault-tolerance
- Expect noise
 - either design a good algorithm, or refine the crowd

Crowd







Crowds Are Human!



- Crowd's abilities should complement the user's
 - diversity ("many eyes")
 - different competence

Find-Fix-Verify in Soylent

Both Shortn and Crowdproof use the Find-Fix-Verify pattern. We will use Shortn as an illustrative example. To provide the user with near-continuous control of paragraph length, Shortn should produce many alternative rewrites without changing the meaning of the original text or introduce⁷ grammatical errors.





Crowd

Software



Ongoing Work





- Open questions in crowd computing
 - Scalability
 - Designing incentives for different kinds of crowds
 - Mobile crowds
 - Expert crowds
 - Task routing
 - Transitioning from crowd to AI



Conclusion





Thanks to support from NSF, Quanta Computer, Xerox ... and a crowd of more than 20,000 turkers

