

# Heads in the Cloud

A professor and several PhD students at MIT examine the challenges and opportunities in human computation.

By Robert C. Miller, Greg Little, Michael Bernstein, Jeffrey P. Bigham, Lydia B. Chilton, Max Goldman, John J. Horton, and Rajeev Nayak

DOI: 10.1145/1869086.1869095

**C**rowd computing is quickly becoming an essential part of the technology landscape. Crowd computing encompasses the interaction among large numbers of people facilitated by software systems and networking technology. Crowds—and by “crowds,” we literally mean a mass of people—are themselves the power that fuels sites like Wikipedia, Twitter, Intrade, and even online labor markets like Amazon Mechanical Turk.

One way to think about crowd computing is as the human analogue to cloud computing. Where the cloud provides access to elastic, highly available computation, and storage resources in the network, the crowd represents access to elastic, highly-available *human* resources, such as human perception and intelligence. Crowd computing offers the strength of software with the intelligence and common sense of human beings.

## HUMAN COMPUTATION

One variant of crowd computing is human computation, which we define as using software to orchestrate a process of small contributions from a crowd to solve a problem that can't be solved by software alone.

Human computation was first popularized by Games With a Purpose (<http://gwap.com>), in which the computation is a side effect of a fun game [8]. For example, the *ESP Game* asks two players to guess words associated with an image, scoring points when their words agree, which makes the

game fun, but also generating useful labels to index the image for searching, which makes it an act of human computation.

Another human computation site is Amazon Mechanical Turk (<http://mturk.com>), a marketplace where people get paid to perform human computation. Users, or “workers,” find short tasks that are posted by “requesters” (the people who need the tasks completed) and get paid small amounts of money for completing them. CrowdFlower (<http://crowdflower.com>) is another site that pays users for computation—in not only real currency, but also virtual currencies for games like *Farmville* and *Mafia Wars*. Social networks like Facebook are also becoming platforms for human computation, motivated by social relationships rather than entertainment or monetary reward.

These platforms make it increasingly feasible to build and deploy systems that use human intelligence as an integral component. But there are at least three challenges to exploring the space

of human computation systems: 1) applications—understanding what's appropriate for human computation and what isn't; 2) programming—learning how to write software that uses human computation; and 3) systems—learning how to get good performance out of a system with humans in the loop.

## APPLICATIONS

What application areas will benefit the most from human computation? What properties do certain problems possess that make them amenable to a successful solution by a hybrid human-software system? Since the end user of such a system is also, typically, human, we can refine this question further: Why does a human end user need to request the help of a human crowd to accomplish a goal, rather than just doing it herself?

One reason is differences in capability. A group of many people has abilities and knowledge that one single end user does not, either innately or because of situational constraints. For example, VizWiz [1] helps blind us-

ers answer questions they have about things around them that they cannot see. The blind person takes a photograph with a smartphone's camera, records a spoken question (also using the phone), and then uploads the query and picture to a crowd of sighted users on the net who are better able to answer it (see **Figure 1**). For example, if a blind person grabs a can out of her cupboard but has forgotten what's inside it, she can snap a photo of the can and its label, upload it, and ask the sighted users what's in the can.

A related system, Sinch [7], draws on the crowd to provide assistance to web-enabled mobile device users who have situational disabilities, such as the limited ability to read a small screen, arthritis or hand tremors that make it difficult to click on small web page targets, and slow networks. With Sinch, the mobile users speak a question into their phone and the crowd searches the web for answers, using their more capable desktop web access, and returning web pages with the requested information highlighted.

Another reason to use a crowd is the “many eyes” principle, which has been claimed as an advantage of open-source software development (the complete phrase is “many eyes make bugs shallow”). We have exploited this principle in Soylent [2], a Microsoft Word extension that uses a crowd for proofreading, shortening, and repetitive editing. A typical run of Soylent may have dozens of people looking at each paragraph of a document, finding errors that a single writer might miss. In fact, a conference paper submitted about Soylent contained a grammatical error that was overlooked by not only Word's built-in grammar checker, but also eight authors and six reviewers. However, when we passed the paper through Soylent, the crowd caught the error.

A corollary of the many eyes principle is diversity. The fact is, a crowd comprises a wide range of ideas, opinions, and skills. For example in Soylent, the system not only identifies writing errors, but also suggests multiple ways to fix them. It can suggest text to cut to save space—a tough task even for skilled authors, who are often reluctant to make cuts. Soylent can typi-

**Figure 1: With VizWiz, blind people take photos using their mobile phones and submit them alongside a question, spoken orally into the phone, shown here above each image. A crowd of anonymous users reply, shown below, with response time given in seconds in parentheses.**



cally trim text down to 85 percent of its original length, without changing the meaning of the text or introducing errors (see **Figure 2**).

## PROGRAMMING

Prototyping a human computation system is hard if you have to entice a crowd to visit your website. Games With a Purpose handles this by making the experience fun—but not all human computation systems are fun enough to be self-motivating, particularly at the prototyping stage. Mechanical Turk is a good prototyping platform for many forms of human computation, because it offers a ready service for recruiting a crowd on demand. And the first prototypes for VizWiz and Soylent were built on Mechanical Turk.

Yet thinking about programming

**“A group of many people has abilities and knowledge that one single end-user does not... The fact is, a crowd comprises a wide range of ideas, opinions, and skills.”**

with human beings inside the system poses special problems. For example with Mechanical Turk, a request for a human to do a small task can take a few minutes and cost a few cents to get a result, which is astounding in one sense (that you can obtain human assistance so quickly and so cheaply), but is abysmally slow and expensive compared to a conventional function call.

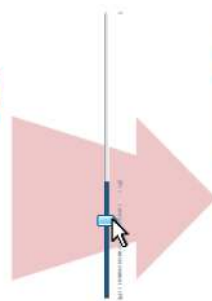
Programmers need new tools that can help them experiment with human computation in their systems. For example, our TurKit toolkit [3] integrates Mechanical Turk calls in a traditional imperative/object-oriented programming paradigm, so that programmers can write algorithms that incorporate human computation in a familiar way. TurKit does this using a novel programming model called “crash and rerun,” which is suited to long-running distributed processes where local computation (done by software) is cheap, and remote work (done by humans) is costly.

The insight of crash-and-rerun programming is that if our program crashes, it is cheap to rerun the entire program up to the place where it crashed. This is true as long as rerunning does not re-perform all the costly external operations from the previous run. The latter problem is solved by recording information in a database every time a costly operation is executed.

Costly operations are marked by a

**Figure 2: In Soylent, after the crowd has suggested words or phrases that can be edited, the end-user can shorten his or her text interactively with a slider. Red text indicates locations where cuts or rewrites have occurred.**

Automatic clustering generally helps separate different kinds of records that need to be edited differently, but it isn't perfect. Sometimes it creates more clusters than needed, because the differences in structure aren't important to the user's particular editing task. For example, if the user only needs to edit near the end of each line, then differences at the start of the line are largely irrelevant, and it isn't necessary to split 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 selection generalization would also be improved by recognizing common text structure like URLs, filenames, email addresses, dates, times, etc.



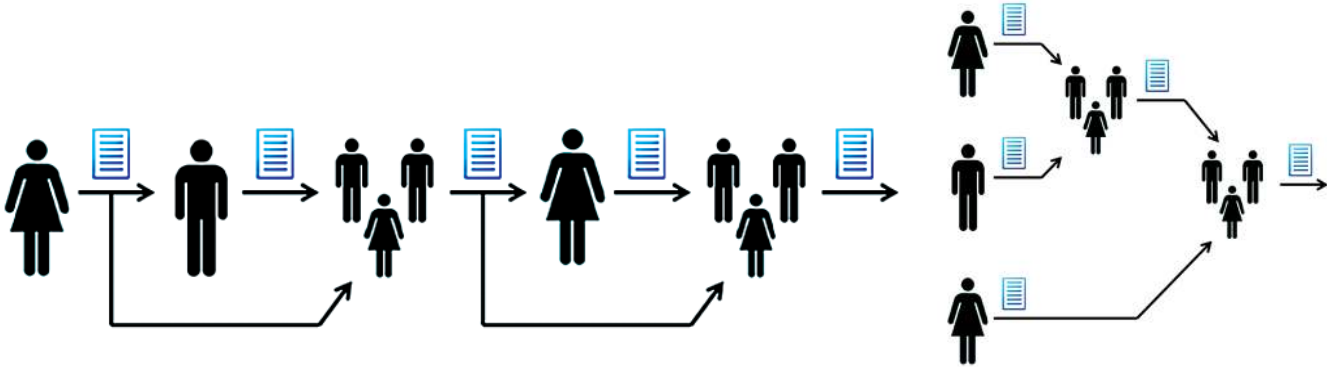
Automatic clustering generally helps separate different kinds of records that need to be edited differently, but it isn't perfect. Sometimes it creates more clusters than needed, because the differences in structure aren't relevant to a specific task. Conversely, sometimes the clustering isn't fine enough, leaving heterogeneous clusters that must be edited one line at a time. One solution to this problem would be to let the user rearrange the clustering manually using drag-and-drop edits. Clustering and selection 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 isn't perfect. Sometimes it creates more clusters than needed, because the differences in structure aren't important to the user's particular editing task. For example, if the user only needs to edit near the end of each line, then differences at the start of the line are largely irrelevant, and it isn't necessary to split 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 using drag-and-drop edits. Clustering and selection 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 isn't perfect. Sometimes it 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 let the user rearrange the clustering manually, perhaps using drag-and-drop to merge and split clusters. Clustering and selection 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 isn't perfect. Sometimes it creates more clusters than needed, because the differences in structure aren't important to the editing task. Conversely, sometimes the clustering isn't fine enough, leaving heterogeneous clusters that must be edited one line at a time. Clustering and selection generalization would also be improved by recognizing common text structure like URLs, filenames, email addresses, dates, times, etc.

**Figure 3: Some human computation processes are iterative (left), involving a succession of interleaved improvement steps (by one person) and voting steps (by several people). Other processes are parallel (right), in which individuals generate original content, and voters simply choose among the alternatives.**



new primitive called "once," meaning they should only be executed once over all reruns of a program. Subsequent runs of the program check the database before performing operations marked with "once" to see if they have already been executed. This model makes it much easier to code algorithms involving human computation. For example, a TurkKit program can sort a list of images using human preference judgments by calling the human computation in the sort algorithm's comparison function, and wrapping those calls in "once" to make them persistent.

Another programming challenge is the development of algorithms and

design patterns that handle the idiosyncrasies of human beings. Humans are not programmable machines, and they don't always follow instructions, unintentionally or otherwise. Sometimes this should be embraced and supported, to harness the creativity and diversity of the crowd. Other times, it simply produces noisy, erroneous, or useless results.

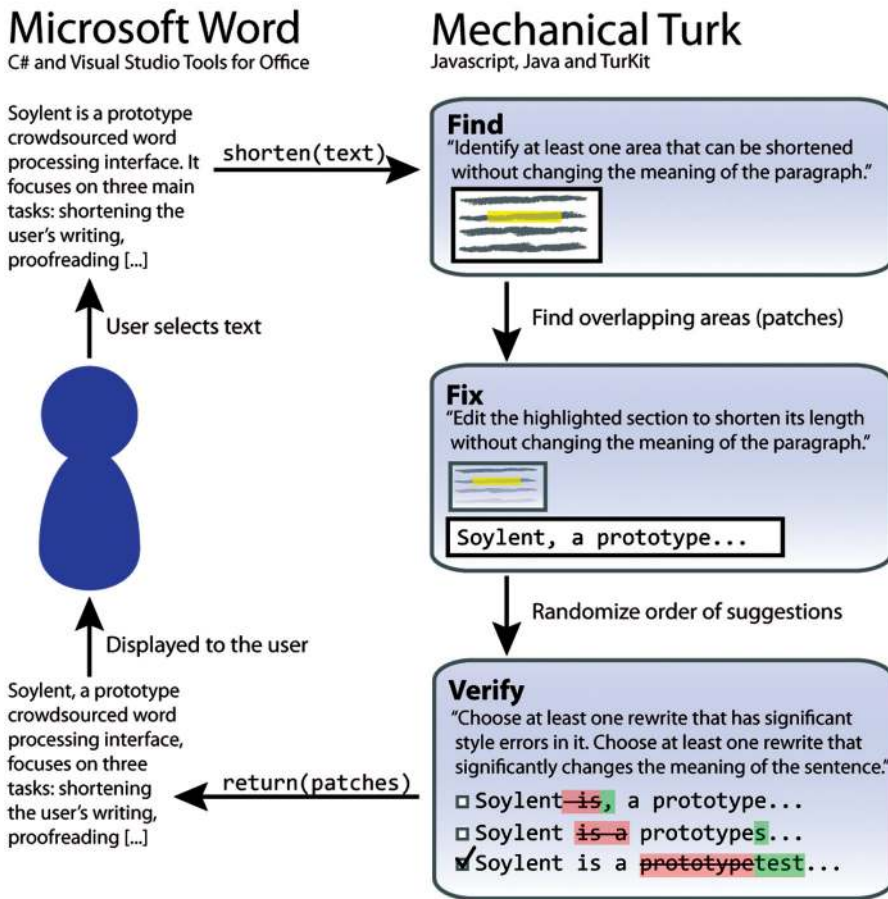
For example, we have studied alternative algorithms for content creation [4]. Iterative processes are similar to Wikipedia or open-source software development. People build on existing content created by others, with voting or independent review ensuring

that the process stays on track. Parallel processes are often seen in design contests, like Threadless.com, where people generate content independently, and then the best is chosen through a vote. See Figure 3.

In experiments involving various kinds of work, such as handwriting transcription, image description, and brainstorming, our results show that iterative processes generally produce higher than average quality than parallel processes. However, in the case of brainstorming, workers riff on good ideas that they see to create more good ones, but the very best ideas seem to come from people working alone. And



**Figure 4: The find-fix-verify algorithm in Soylent identifies patches in need of editing, suggests fixes to the patches, and votes on those fixes.**



with transcription tasks, it turns out that showing workers the guesses of other workers often leads them astray, especially if the guesses are self-consistent but wrong.

Crowd workers exhibit high variance in the amount of effort they invest in a task. Some are lazy turkers, who do as little work as necessary to get paid, while others are eager beavers, who go above and beyond the requirements, either to be helpful or to signal that they aren't lazy turkers, but in counterproductive ways. We need new design patterns for algorithms involving human computation that recognize and control this behavior.

For example, Soylent uses a find-fix-verify pattern to improve the quality of proofreading and document shortening (Figure 4). In this pattern, some workers find problems, other workers fix them, and still other workers verify the fixes. But questions remain. What

other algorithms and design patterns are useful? How should algorithms involving human computation be evaluated and compared from a theoretical point of view?

**SYSTEMS PROBLEMS**

Moving from prototyping to actual deployment requires facing questions about how to obtain a reliable and well-performing source of human computation for the system. How can we recruit a crowd to help, and motivate it to continue to help over time, while optimizing for cost, latency, bandwidth, quality, churn, and other parameters?

For paid crowds, these questions intersect with labor economics. Some of our recent work has found that workers in human computation markets like Mechanical Turk behave in unusual ways. For example, instead of seeking work that provides a target wage, they often seek a target earning amount,

and simply work until they reach their target, consistent with game-playing behavior [5].

Another difference in these markets is the overwhelming importance of searchability. Workers' ability to find tasks they want to do is strongly affected by the kind of interface the market offers. Mechanical Turk, for example, typically displays a list of thousands of available tasks, divided into hundreds of result pages, with few effective tools for searching or filtering this list. We have found that most workers simply choose a particular sort order and work their way through the list. They most often sort by newest task, or most tasks available, and surprisingly not by price. The speed of completion of a task is strongly affected by its ability to be found, which may not be strongly related to the monetary reward it offers [6].

We can also think about human computation in computer systems terms, such as cost, latency, and parallelism. Services like VizWiz and Sinch need to return answers quickly, and to support that, we have developed an approach (and accompanying implementation) called quikTurkit that provides a layer of abstraction on top of Mechanical Turk to intelligently recruit multiple workers before they're needed.

In a field deployment of VizWiz, users had to wait a little longer than two minutes on average to get their first answer. Wait times decreased sharply when questions and photos were easy for workers to understand. Answers were returned at an average cost per question of only \$0.07 for 3.3 answers. Given that other visual-assistance tools for the blind can cost upwards of \$1,000 (the equivalent of nearly 15,000 uses of VizWiz), we believe that human computation embedded in an inexpensive software system can be not only more effective but also competitive with, or even cheaper than, existing pure software solutions. When set to maintain a steady pool of workers (at a cost of less than \$5 per hour), quikTurkit can obtain answers in less than 30 seconds.

Beyond monetary compensation, many other reasons entice people to participate in human computation, in-

cluding altruism, entertainment, and friendship. How do those motivations influence system performance? And how should the systems be designed to encourage some motivations, and perhaps discourage others?

After demonstrating that VizWiz was feasible using paid strangers on Mechanical Turk, we also ported it to Facebook, so that a blind user's sighted friends can help. We are currently studying how people (at least in this context) choose to trade off the strengths and weaknesses of each service. Mechanical Turk is fast but costs money. Facebook is free, and the user's friends might be more motivated to answer, or even more capable since they know more about the person. On the other hand, the user might be less willing to ask certain personal questions to his or her friends, rather than asking an anonymous Mechanical Turk crowd.

## PEOPLE VS SYSTEMS

The gap between what software can do and what people can do is shrinking, but a gap of some sort will exist for a long time. Automatic techniques need to be able to fallback to people when necessary to fill in the gaps, enabling interactions that automatic techniques alone can't yet support and helping us design for the future.

So-called *Wizard of Oz* prototyping (wherein a human is hiding behind the curtain, so to speak) is a venerable technique in human-computer interaction and artificial intelligence that makes an intelligent system (or even a not-so-intelligent one) appear to work even though a software backend isn't ready yet. With platforms like Mechanical Turk and Facebook that make human computation practical, we are now at the point where *Wizard of Oz* is not just for prototyping anymore. We can build useful systems with human power inside, and actually deploy them to real users. These systems will stretch the limits of what software can do, and allow us to find out whether the ideas even work and how people would use them.

In addition, we can collect data from actual system use, like VizWiz queries and photos, that might eventually help to replace some or all of the human

power with artificial intelligence. From this perspective, AI would speed up performance and reduce labor costs. But human computation made the system possible in the first place.

## Acknowledgements

Ideas and work in this article come from many students and collaborators, including Mark Ackerman, David Crowell, Bjoern Hartmann, David Karger, Marc Grimson, and Katrina Panovich. This work was supported in part by Quanta Computer, NSF, and Xerox.

## Biographies

Robert C. Miller is an associate professor of computer science at Massachusetts Institute of Technology. He grew up in rural Louisiana and plays the accordion, although extremely poorly.

Danny "Greg" Little is a PhD student in computer science at MIT. He is probably asleep right now.

Michael Bernstein is a PhD student in computer science at MIT and a features editor for this magazine.

Jeffrey P. Bigham is an assistant professor of computer science at University of Rochester, and he wants to be mayor of your living room. Don't let him in.

Lydia B. Chilton is a PhD student in computer science at University of Washington. When not studying human computation, she is a consultant to the United Federation of Planets and makes the occasional journey on the USS Enterprise with her old pals Kirk, Spock, and McCoy.

Max Goldman is a PhD student in computer science at MIT. When he gives a talk, his performance is so lively and engaging that it distracts the audience from the actual research results, but everybody goes away happy.

John J. Horton is a PhD student in public policy at Harvard University. He is having trouble thinking of a public policy angle for his human computation research and needs to defend his dissertation soon. This is a growing problem that, to date, Turkers have been unable to solve.

Rajeev Nayak is a graduate student in computer science at MIT. He bats .400, shoots .575 from the field, sings a *cappella*, and watches *Glee* religiously.

## References

1. Bigham, J.P., Jayant, C., Ji, H., Little, G., Miller, A., Miller, R.C., Miller, R., Whyte, B., White, S., Yeh, T. VizWiz: Nearly Real-time Answers to Visual Questions. *UIST 2010*.
2. Bernstein, M., Little, G., Miller, R.C., Hartmann, B., Ackerman, M.S., Karger, D.R., Crowell, D., Panovich, K. Soylent: A Word Processor with a Crowd Inside. *UIST 2010*.
3. Little, G., Chilton, L., Goldman, M., Miller, R.C. TurKit: Human Computation Algorithms on Mechanical Turk. *UIST 2010*.
4. Little, G., Chilton, L., Goldman, M., Miller, R.C. Exploring Iterative and Parallel Human Computation Processes. *HCOMP 2010*, to appear.
5. Horton, J.J. and Chilton, L. The Labor Economics of Paid Crowdsourcing. *EC 2010*.
6. Chilton, L., Horton, J.J., Miller, R.C., Azenkot, S. Task Search in a Human Computation Market. *HCOMP 2010*, to appear.
7. Nayak, R., et al. Sinch: Searching Intelligently on a Mobile Device. *CHI 2011*, in submission.
8. von Ahn, L., and Dabbish, L. Designing Games with a Purpose. *CACM*, 51, 8, August 2008.

## ACRONYMS

**AMT** Amazon Mechanical Turk: a web service owned by Amazon that facilitates crowdsourcing

**CAPTCHA** complete automated public Turing test to tell computers and humans apart; it's a contrived acronym intentionally redolent of the word "capture," used to describe a test issued on web forms to protect against automated responses

**GWAP** Game with a Purpose: a computer game that layers a recreational challenge on top of a problem that demands human intelligence for efficient solution, e.g., protein folding

**HCIR** human-computer information retrieval

**HIT** human intelligence task: a task that an AMT requester is willing to pay to have accomplished by AMT providers. More generally, a task that may be best completed via crowdsourcing

**HuGS** human-guided search: A research project investigating a strategy for search and optimization problems that incorporates human intuition and insight

**reCAPTCHA** a kind of "CAPTCHA" (see above) service that helps to digitize books, newspapers and old radio shows. Used to detect whether a user is a human or a computer [bot]

**TF-IDF** term frequency-inverse document frequency: A weight formula used in data mining to determine the importance of a particular term to a document in a corpus of text