

Connect the Dots: Supporting Intelligence Analysis with Crowdsourcing, Context Slices, and Visualization

Tianyi Li, Chris North, Kurt Luther
Department of Computer Science, Virginia Tech
{tianyili,kluther,north}@vt.edu

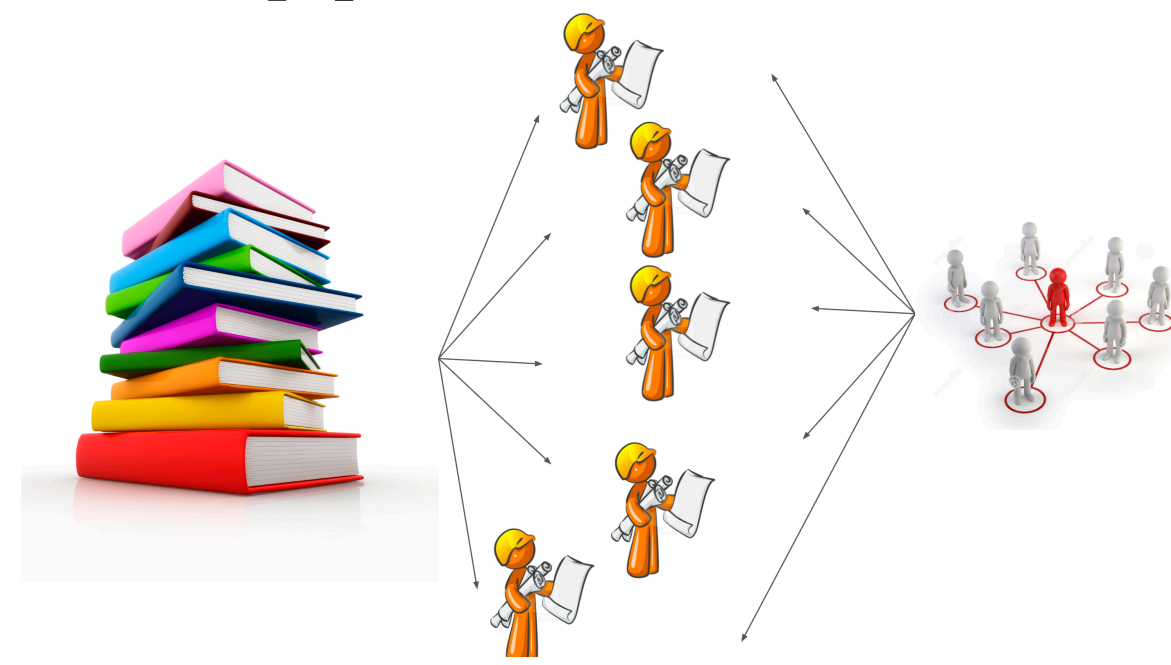
Problem Statement

Making sense of large text datasets is a difficult problem in many domains that does not scale well for individuals. Crowdsourcing presents new opportunities for large-scale sensemaking, but we must first overcome the challenge of enabling many distributed novice workers to contribute meaningfully.

Background and Motivations

Challenges

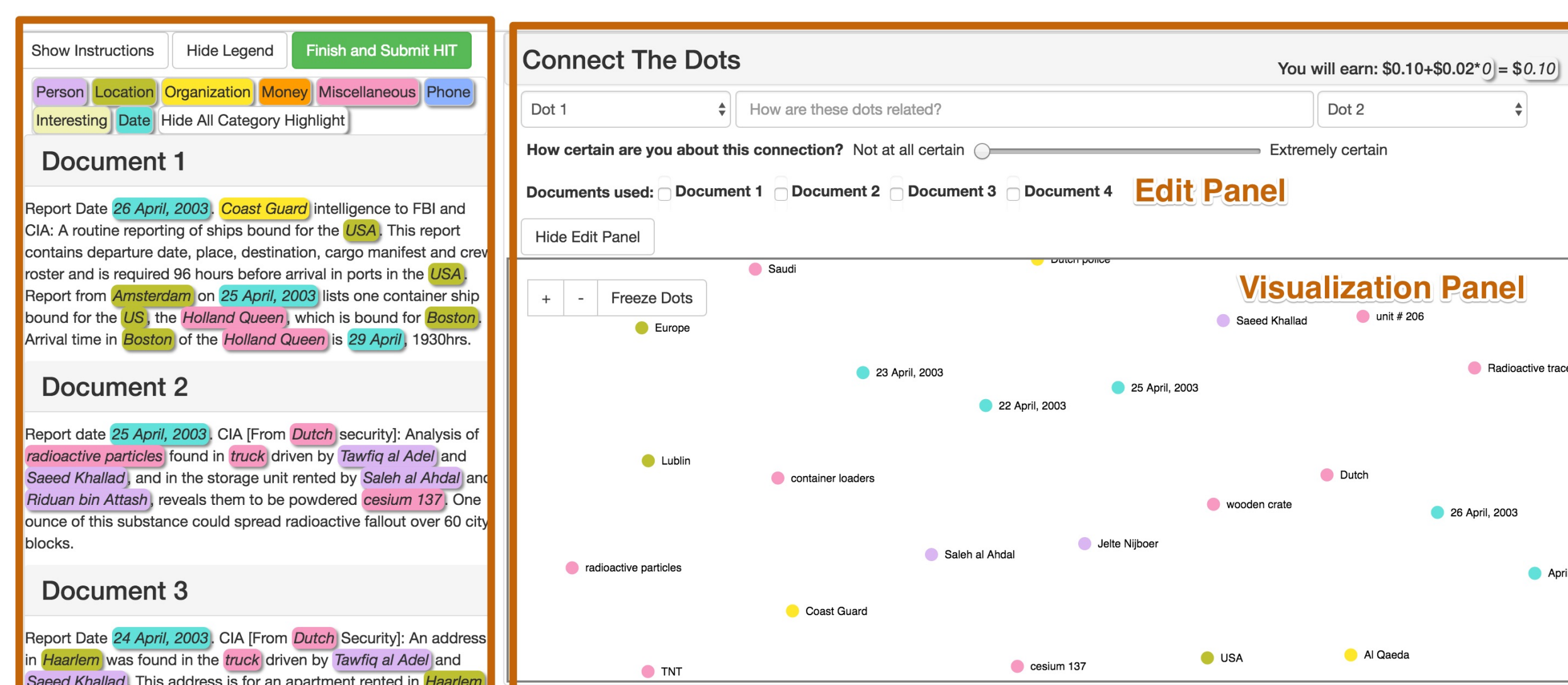
- Overwhelming amount of information [3]
- Limits of human cognition
 - Limited Time
 - Limited Attention
 - Lack of Expertise
- Coordination [4,5]
 - Communication
 - Task distribution and design
 - Expert-crowd interaction



Opportunities and supports from related work

- Crowdsourcing [2]
 - Integrate human intelligence at large scale
- Sensemaking loop [1]
 - Modularizes the process of expert's sensemaking

System Description



Document View

Connection Workspace

Document View

- Legend
- Document Texts
- Entities
- Category Highlights

Connection Workspace

- Edit panel
- Visualization panel
 - Force-directed graph
 - Zoom and drag

Experiment Design

Dataset

- *Sign of the Crescent*
 - Training materials for intelligence analysis
 - 3 coordinated terroristic actions
- In this experiment
 - 10 documents relevant to one of the actions
 - 55 context slices
 - 10 single-document slices
 - 45 double-document slices
 - all possible combinations

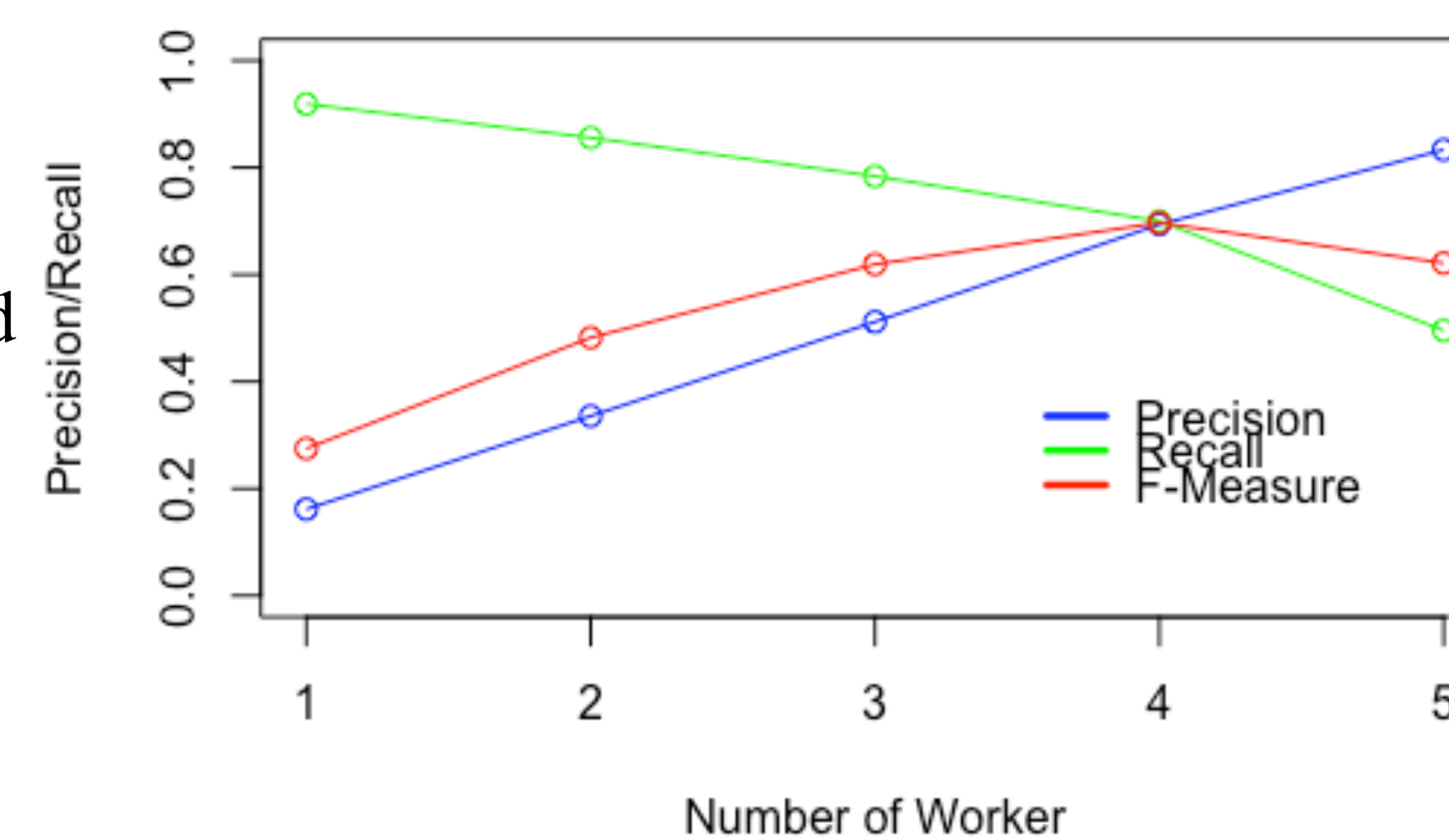
Participants

- Amazon Mechanical Turk (AMT)
- 275 crowd workers
 - 5 workers / context slice
 - Randomly assigned
 - Each worker was unique and assigned to only one HIT
- Pay
 - Workers have to make a minimum number of connections for assigned context slices
 - \$0.10+\$0.02*extra connections

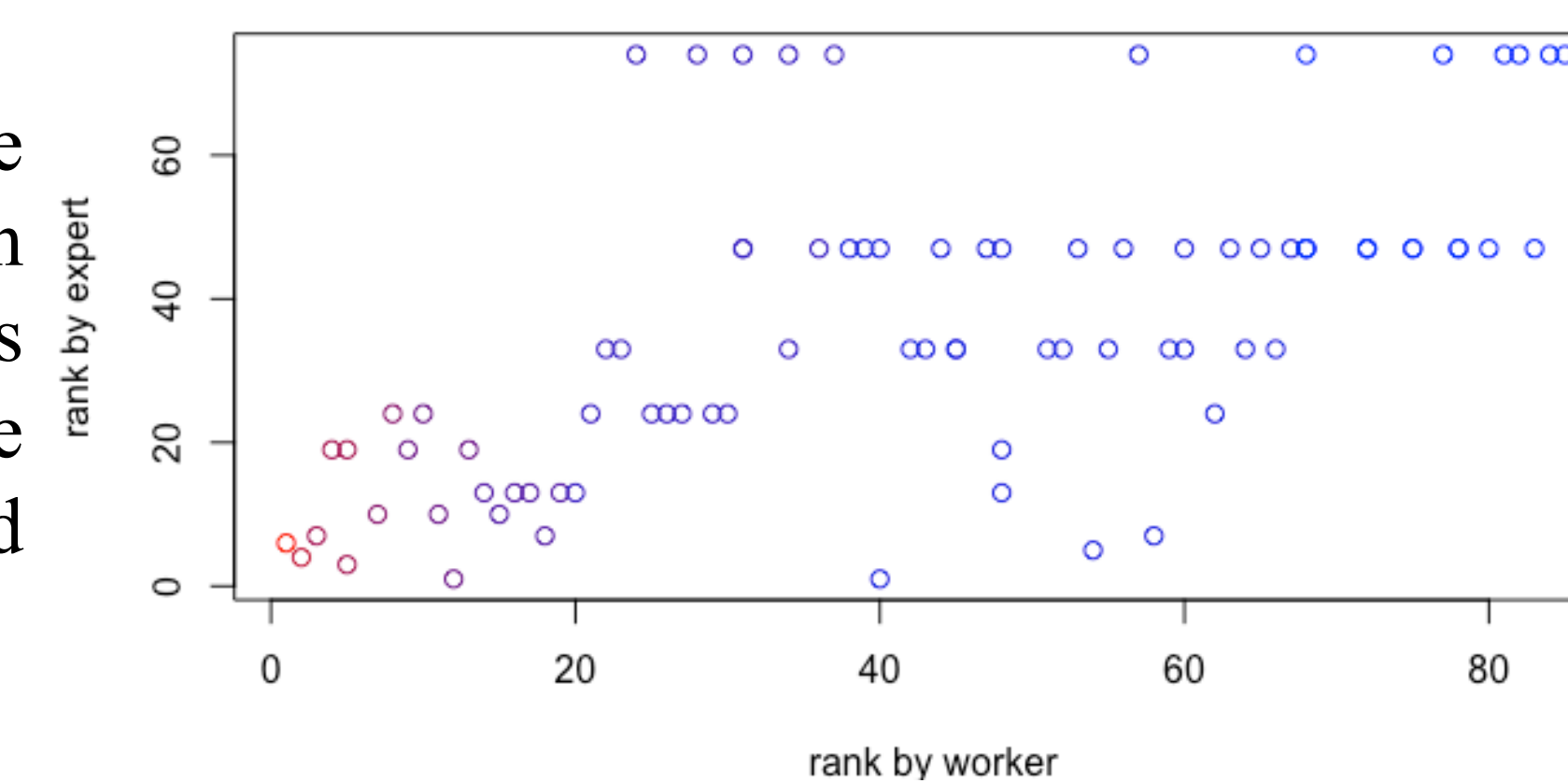
Results

The 275 crowd workers created 5992 connections from the 55 context slices in total. Removing overlapped entity pairs connected by multiple workers yields 631 pairs of entities. Average time spent is approximately 20 minutes per HIT.

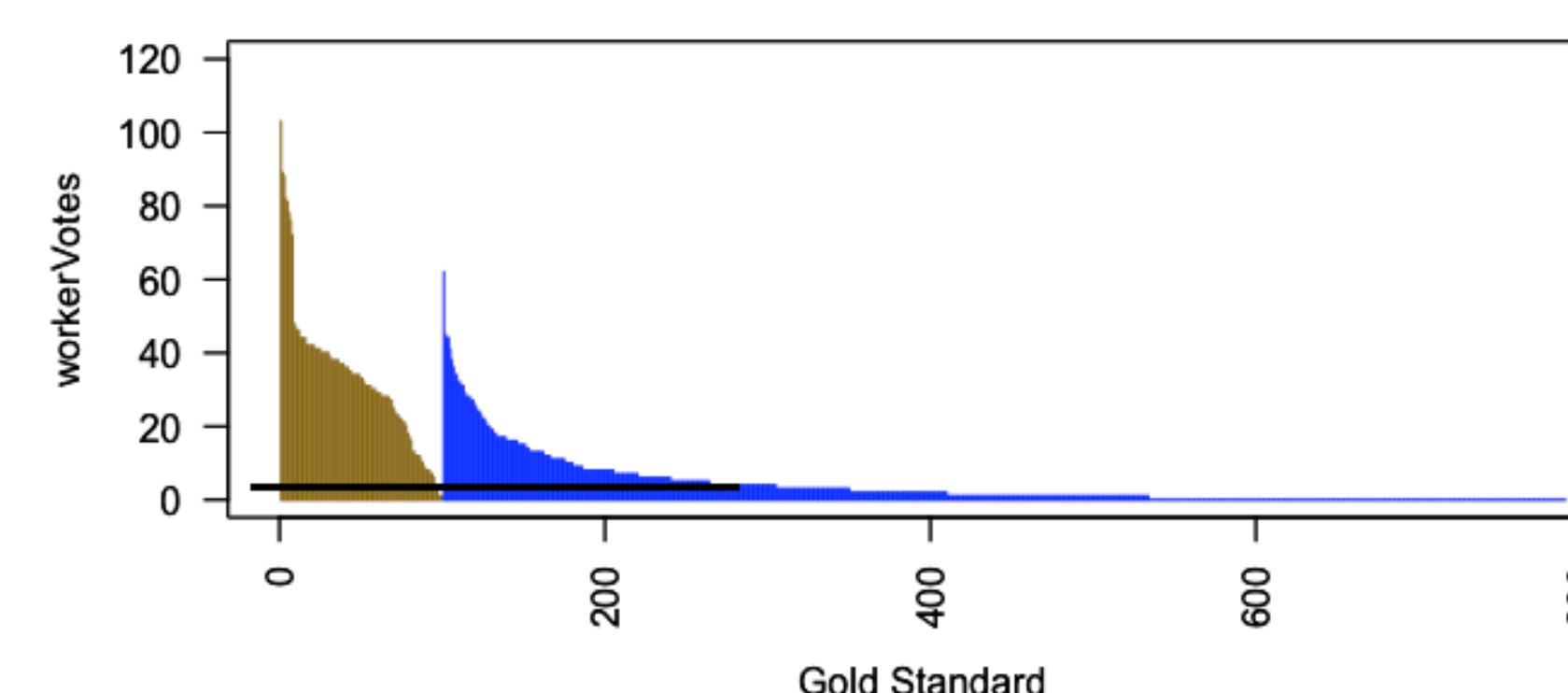
- **Rank of connectivity of entities**
Given a set of crowd-generated entity pairs \mathcal{C} (631), the overlapped entity pairs $\mathcal{O} = \mathcal{C} \cap G$.
 - The precision value $P = |\mathcal{O}|/|\mathcal{C}|$
 - The recall value $R = |\mathcal{O}|/|G|$



- **Rank of connectivity of entities**
The entity with the highest degree in both the crowd-generated graph and the gold standard graph is ranked #1 and the entity with the fewest (if not 0) links is ranked with larger numbers.



- **Rank of votes of entity pairs**
The gold region (left) of the bar plot shows the node pairs from gold standard connections. The blue region (right) shows the remaining ones. Y-axis shows the number of workers connecting them.



Discussion

Crowd coverage vs. gold standard and thresholds

- The crowd-generated connections combined covers nearly all gold standard connections.
 - 92%, or 102 of 111
- Eight of the missing connections from the gold standard were made from more than two documents.
 - Thus, crowd workers were not given enough context in this experiment to connect them;
 - Testing larger context slices might be helpful to uncover these.
- The remaining missed connection represented supporting evidence for other less obvious connections in the solution, which were successfully found by crowd workers.
 - Experts are very likely to be able to recover that connection from the connections found by the crowd.

Prioritized and guided search focus

- Ranking nodes according to their degrees in the graph of crowd-generated connections yielded similar results to those from the gold standard connections.
- The crowd showed great potential in finding important entities and make connections with them.
- This ranking can serve as a starting point in expert analysts' sensemaking process to help guide and refine their search of the solution space and prioritize entities within the same contexts.

Acknowledgements

This research is funded by NSF IIS #1527453. I want to thank my advisors Dr. Chris North and Dr. Kurt Luther for their great support. I also want to thank my lab mate Dr. Maoyuan Sun and Yali Bian, and my undergraduate research assistants Edward McEnrue, Jazmine Zurita and Chris Lai.

References

- 1) Pirolli, Peter, and Stuart Card. "The sensemaking process and leverage points for analyst technology as identified through cognitive task analysis." *Proceedings of international conference on intelligence analysis*. Vol. 5. 2005.
- 2) Paul André, Aniket Kittur, and Steven P. Dow. 2014. Crowd Synthesis: Extracting Categories and Clusters from Complex Data. In *Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work & Social Computing* (CSCW '14), 989–998. <http://doi.org/10.1145/2531602.2531653>
- 3) James J. Thomas and Kristin A. Cook (eds.). 2005. *Illuminating the Path: The Research and Development Agenda for Visual Analytics*. National Visualization and Analytics Ctr
- 4) Haoqi Zhang, Edith Law, Rob Miller, Krzysztof Gajos, David Parkes, and Eric Horvitz. 2012. Human Computation Tasks with Global Constraints. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (CHI '12), 217–226. <http://doi.org/10.1145/2207676.2207708>
- 5) Narges Mahyar and Melanie Tory. 2014. Supporting Communication and Coordination in Collaborative Sensemaking. *IEEE transactions on visualization and computer graphics* 20, 12: 1633–1642. <http://doi.org/10.1109/TVCG.2014.2346573>