





Frontiers in Massive Data Analysis

ISBN
978-0-309-28778-4

190 pages
6 x 9
PAPERBACK (2013)

Committee on the Analysis of Massive Data; Committee on Applied and Theoretical Statistics; Board on Mathematical Sciences and Their Applications; Division on Engineering and Physical Sciences; National Research Council

 Add book to cart

 Find similar titles

 Share this PDF



Visit the National Academies Press online and register for...

- ✓ Instant access to free PDF downloads of titles from the
 - NATIONAL ACADEMY OF SCIENCES
 - NATIONAL ACADEMY OF ENGINEERING
 - INSTITUTE OF MEDICINE
 - NATIONAL RESEARCH COUNCIL
- ✓ 10% off print titles
- ✓ Custom notification of new releases in your field of interest
- ✓ Special offers and discounts

Distribution, posting, or copying of this PDF is strictly prohibited without written permission of the National Academies Press. Unless otherwise indicated, all materials in this PDF are copyrighted by the National Academy of Sciences. Request reprint permission for this book

9

Human Interaction with Data

INTRODUCTION

Until recently, data analysis was the purview of a small number of experts in a limited number of fields. In recent years, however, more and more organizations across an expanding range of domains are recognizing the importance of data analysis in meeting their objectives. Making analysis more useful to a wider range of people for a more diverse range of purposes is one of the key challenges to be addressed in the development of data-driven systems.

In many, perhaps most, scenarios today and in the near-term future, people are the ultimate consumers of the insights from data analysis. That is, the analysis of data is used to drive and/or improve human decision-making and knowledge. As such, methods for visualization and exploration of complex and vast data constitute a crucial component of an analytics infrastructure. The field of human-computer interactions has made great progress in the display and manipulation of complex information, but the increasing scale, breadth, and diversity of information provide continued challenges in the area.

People are not, however, merely consumers of data and data analysis. In many analytics usage scenarios, people are also the source (and often the subject) of the data being analyzed. Continuing improvements in network connectivity and the ubiquity of sophisticated communications and computing devices has made data collection easier, particularly as more activities are done online. Moreover, increasing network connectivity has been leveraged by a number of platforms that can allow people to participate

directly in the data analysis process. Crowdsourcing is the term used to describe the harnessing of the efforts of individual people and groups to accomplish a larger task. Crowdsourcing systems have taken on many forms, driven largely by advances in network connectivity, the development of service-oriented platforms and application programming interfaces (APIs) for accomplishing distributed work, and the emergence of user-generated content sites (sometimes referred to as “Web 2.0”) that include socially oriented and other mechanisms for filtering, vetting, and organizing content.

This chapter discusses several aspects of crowdsourcing that could contribute to extracting information from massive data. Crowdsourced data acquisition is the process of obtaining data from groups either explicitly—for example, by people deliberately contributing content to a website—or implicitly, as a side effect of computer-based or other networked activity. This has already been shown to be a powerful mechanism for tasks as varied as monitoring road traffic, identifying and locating distributed phenomena, and discovering emerging trends and events.

For the purposes of this report, a perhaps more interesting development in crowdsourcing is the involvement of people to aid directly in the analysis process. It is well known that computers and people excel at very different types of tasks. While algorithm developers continue to make progress in enabling computers to address tasks of greater complexity, there remain many types of analysis that can be more effectively done by people, even when compared to the most sophisticated computers and algorithms. Such analyses include deep language understanding and certain kinds of pattern recognition and outlier detection. Thus, there has been significant recent work and recognition of further opportunities in hybrid computer/human data analysis.

These trends are leading to the increased understanding of the role of people in all phases of the data processing lifecycle—from data collection through analysis to result consumption, and ultimately to decision making. The human dimension carries with it a new set of concerns, design constraints, and opportunities that must be addressed in the development of systems for massive data analysis. In addition, massive data calls for new approaches to data visualization, which is often used in exploratory data analysis. This chapter thus focuses on the various aspects of human interaction with data, with an emphasis on three areas: data visualization and exploration, crowdsourced data acquisition, and hybrid computer/human data analysis.

STATE OF THE ART

Data Visualization and Exploration¹

Information visualization technologies and visual analytics processes have matured rapidly in the past two decades and continued to gain commercial adoption, while the research enterprise has expanded. Successful commercial tools include some that stand alone, such as Spotfire, Tableau, Palantir, Centrifuge, i2, and Hive Group, as well as some that are embedded in other systems, such as IBM ILOG, SAS JMP, Microsoft Proclarity, Google Gapminder, and SAP Xcelsius. In addition, open-source toolkits such as R, Prefuse, ProtoVis, Piccolo, NodeXL, and Xmdv support programmers. Academic conferences and journals in this area are active, and an increasing number of graduate courses are available.

At the same time, news media and web-based blogs have focused intense public attention on interactive infographics that deal with key public events such as elections, financial developments, social media impacts, and health care/wellness. Information visualizations provide for rapid user interaction with data through rich control panels with selectors to filter data, with results displayed in multiple coordinated windows. Larger displays (10 megapixels) or display arrays (100-1,000 megapixels) enable users to operate dozens of windows at a time from a single control panel. Large displays present users with millions of markers simultaneously, allowing them to manipulate these views by dynamic query sliders within 100 ms. Rapid exploration of data sets with 10 million or more records supports hypothesis formation and testing, enabling users to gain insights about important relationships or significant outliers.

An important distinction is often made between the more established field of scientific visualization and the emerging field of information visualization. Scientific visualizations typically deal with two-dimensional and three-dimensional data about physical systems in which questions deal with position—for example, the location of highest turbulence in the airflow over aircraft wings, the location and path of intense hurricanes, or the site of blockages in arteries. In contrast, information visualizations typically deal with time series, hierarchies, networks, or multi-variate or textual data, in which questions revolve around finding relationships, clusters, gaps, outliers, and anomalies.

Information visualization problems might be typified by the following examples:

¹ The committee thanks Ben Shneiderman of the University of Maryland and Patrick Hanrahan of Stanford University for very helpful inputs to this section.

- Find the strongest correlations in daily stock market performance over 5 years for 10,000 stocks (analyze a set of time series);
- Identify duplicate directories in petabyte-scale hard drives (search through hierarchies);
- Find the most-central nodes in a social network of 500 million users (network analysis); and
- Discover related variables from 100-dimensional data sets with a billion rows (multivariate analysis).

Increasingly, text-analytics projects are searching Web-scale data sets for trending phrases (e.g., Google's *culuromics.org*), unusual combinations, or anomalous corpora that avoid certain phrases.

There are hundreds of interesting visualization techniques, including simple bar charts, line charts, treemaps, graphs, geographic maps, and textual representations such as tag clouds. Information visualization tools, however, often rely on rich interactions between multiple simultaneous visualizations. For example, a user might select a set of markers in one window, and the tool highlights related markers in all windows, so that relationships can be seen. Such a capability might allow a user to select bars on a timeline indicating earthquakes and, from that, automatically highlight markers on a map and show each earthquake in a scattergram organized by intensity versus number of fatalities, color coded by whether the quake was under water or under land areas. Similarly, the movement of any dynamic query slider immediately filters out the related markers in all views. For example, a user could filter out the quakes whose epicenter was deeper than 3 miles to study the impact of deep quakes only.

The general approach to seeking information is to overview first, zoom and filter, and then find details on demand. This simple notion conveys the basic idea of an exploratory process that has been widely applied. More recently, however, attention in the visual analytics community has shifted to process models for exploration. Such models range from simple 4-step approaches that gather information, re-represent it, develop insights, and present results, to elaborate 16-step models and domain-specific approaches for medical histories, financial transactions, or gene expression data.² The process models help guide users through steps that include data cleaning (remove errors, duplicates, missing data, etc.), filtering (select appropriate subsets), aggregation (clustering, grouping, hierarchical organization), and recording insights (marking, annotation, grouping).³

The steps in such process models have perceptual, cognitive, and domain-specific aspects that lead researchers to consider visual analytics as a

² See Thomas and Cook (2005) for an overview.

³ See, e.g., Perer and Shneiderman (2006).

sense-making process, which requires validation by empirical assessments. While some aspects of interface design and usage can be tested in controlled empirical user studies and expert reviews, information visualization researchers have creatively found new evaluation methods. Often, case studies of usage by actual researchers working with their own data over periods of weeks or months have been used to validate the utility of information visualization tools and visual analytics processes (Shneiderman and Plaisant, 2006).

Crowdsourced Data Acquisition

The idea of coordinating groups of people to perform computational tasks has a long history. Small groups of people were used to catalog scientific observations as early as the 1700s, and groups of hundreds of people were organized to compute and compile tables of mathematical functions in the early part of the 20th century (Grier, 2005). Recently, as the Internet has enabled large-scale organization and interaction of people, there has been a resurgence of interest in crowd-based data gathering and computation. The term “crowdsourcing” is used to describe a number of different approaches, which can be grouped into two general classes: those that leverage human activity, and those that leverage human intelligence. In the former case, data are produced and gathered, and work is performed as a by-product of individuals’ behavior on the Web or in other networked environments. In the latter case, groups of people are organized and explicitly tasked with performing a job, solving a problem, or contributing content or other work product.

The first category of crowdsourcing consists of techniques for garnering useful information generated as a by-product of human activity. Such information is sometimes referred to as “data exhaust.” For example, search companies can continuously improve their spell checking and recommendation systems using data generated as users enter misspelled search terms and then click on a differently spelled (correct) result. Many Web companies engage in similar types of activity mining, for example, to choose which content or advertising to display to specific users based on search history, access history, demographics, etc. Many other online activities, such as recommending a restaurant, “re-tweeting” an article, and so on, also provide valuable information. The implicit knowledge that is gleaned from user behaviors can be used to create a new predictive model or to augment and improve an existing one.

In these examples, users’ online activity is used to predict their intent or discern their interests. Online activity can also be used to understand events and trends in the real world. For example, there has been significant interest lately in continuously monitoring online forums and social media to

detect emerging news stories. As another example, Google researchers have demonstrated the ability to accurately detect flu outbreaks in particular geographic regions by noting patterns in search requests about flu symptoms and remedies.⁴ Importantly, they demonstrated that such methods sped up the detection of outbreaks by weeks compared to the traditional reporting methods currently used by the Centers for Disease Control and Prevention and others.

Another form of crowdsourced data acquisition is known as participatory sensing (Estrin, 2010). The convergence of sensing, communication, and computational power on mobile devices such as cellular phones creates an unprecedented opportunity for crowdsourcing data. Smartphones are increasingly integrating sensor suites (with data from the Global Positioning System, accelerometers, magnetometers, light sensors, cameras, and so on), and they are capable of processing the geolocalized data and of transmitting them. As such, participatory sensing has become a paradigm for gathering data at global scales, which can reveal patterns of humans in the built environment. Early successes have been in the area of traffic monitoring and congestion prediction,⁵ but it is possible to build many applications that integrate physical monitoring with maps. Examples of other applications include monitoring of environmental factors such as air quality, sound pollution, ground shaking (i.e., earthquake detection), and water quality and motion. Urban planning can be aided by the monitoring of vehicular as well as pedestrian traffic. Privacy concerns must be taken into account and handled carefully in some of these cases.

In all the cases described above, data are collected as a by-product of peoples' online or on-network behavior. Another class of crowdsourcing approaches more actively designs online activity with the express purpose of enticing people to provide useful data and processing. "Games with a purpose" are online games that entice users to perform useful work while playing online games (see, e.g., von Ahn and Dabbish, 2008). An early example was the ESP Game, developed at Carnegie Mellon University, in which players listed terms that describe images, simultaneously earning points in the game and labeling the images to aid in future image search queries.

A related approach, called re-captcha,⁶ leverages human activity to augment optical character recognition (OCR). In re-captcha, users attempting to access an online resource are presented with two sets of characters to transcribe. One set of characters is known to the algorithm and is presented

⁴ E.g., Explore Flu Trends Around the World, available at <http://www.google.org/flu-trends>.

⁵ E.g., Mobile Millennium, University of California, Berkeley, Snapshot of Mobile Millennium Traffic in San Francisco and the Bay Area, available at <http://traffic.berkeley.edu/>.

⁶ The ReCAPTCHA website is available at <http://www.google.com/recaptcha>.

in a format that is difficult for machines to identify. The other set of characters presented to the user is a portion of text that an OCR algorithm was unable to recognize. The idea is that by correctly entering the first set of characters, a user verifies that he or she is not a machine, and by entering the second set of characters, the user then effectively performs an OCR task that an OCR algorithm was unable to perform.

HYBRID HUMAN/COMPUTER DATA ANALYSIS

In the crowdsourcing techniques described in the previous section, human input was obtained primarily as part of the data-collection process. More recently, a number of systems have been developed that more explicitly involve people in computational tasks. Although the fields of artificial intelligence and machine learning have made great progress in recent years in solving many problems that were long considered to require human intelligence—for example, natural language processing, language translation, chess playing, winning the television game show *Jeopardy*, and various prediction and planning tasks—there are still many tasks where human perception, and peoples' ability to disambiguate, understand context, and make subjective judgments, exceed the capabilities of even the most sophisticated computing systems. For such problems, substantial benefit can be obtained by leveraging human intelligence.

While Quinn and Bederson (2011) distinguish human computation from crowdsourcing, defining the former as replacing computers with humans and the latter as “replacing traditional human workers with members of the public,” many in both the research community and the general public do not make such a distinction. Thus, crowdsourcing is often used to refer to either type of human involvement, and that convention is followed here.

Some types of crowdsourced systems that can be used to involve people in the process of analyzing data are the following:

- *User-generated content sites.* Wikipedia is a prominent example of a user-generated content site where people create, modify, and update pages of information about a huge range of topics. More specialized sites exist for reviews and recommendations of movies, restaurants, products, and so on. In addition to creating basic content, in many of these systems users are also able to edit and curate the data, resulting in collections of data that can be useful in many analytics tasks.
- *Task platforms.* Much of the interest around crowdsourcing has been focused on an emerging set of systems known as microtask platforms. A microtask platform creates a marketplace in which requesters offer tasks and workers accept and perform the tasks.

Microtasks usually do not require any special training and typically take no longer than 1 minute to complete, although they can take longer. Typical microtasks include labeling images, cleaning and verifying data, locating missing information, and performing subjective or context-based comparisons. One of the leading platforms at present is Amazon Mechanical Turk (AMT). In AMT, workers from anywhere in the world can participate, and there are thought to be hundreds of thousands of people who perform jobs on the system.

Other task-oriented platforms have been developed or proposed to do more sophisticated work. For example, specialized platforms have been developed to crowdsource creative work such as designing logos (e.g., 99designs) or writing code (e.g., TopCoder). In addition, some groups have developed programming languages to encode more sophisticated multistep tasks, such as Turkkit (Little et al., 2010), or market-based mechanisms for organizing larger tasks (Shahaf and Horvitz, 2010). These types of platforms can be used to get human participation on a range of analytics tasks, from simple disambiguation to more sophisticated iterative processing.

- *Crowdsourced query processing.* Recently, a number of research efforts have investigated the integration of crowdsourcing with query processing as performed by relational database systems. Traditional database systems are limited in their ability to tolerate inconsistent or missing information, which has restricted the domains in which they can be applied largely to those with structured, fairly clean information. Crowdsourcing based on application programming interfaces (APIs) provides an opportunity to engage humans to help with those tasks that are not sufficiently handled by database systems today. CrowdDB (Franklin et al., 2011) and Qurk (Marcus et al., 2011) are examples of such experimental systems.
- *Question-answering systems.* Question-answering systems are another type of system for enlisting human intelligence. Many different kinds of human-powered or human-assisted sites have been developed. These include general knowledge sites where humans help answer questions (e.g., Cha Cha), general expertise-based sites, where people with expertise in particular topics answer questions on those topics (e.g., Quora), and specialized sites focused on a particular topic (e.g., StackOverflow for computer-programming-related questions).
- *Massive multi-player online games.* Another type of crowdsourcing site uses gamification to encourage people to contribute to solving a problem. Such games can be useful for simulating complex social systems, predicting events (e.g., prediction markets), or for solving

specific types of problems. One successful example of the latter type of system is the FoldIt site,⁷ where people compete to most accurately predict the way that certain proteins will fold. FoldIt has been competitive with, and in some cases even beaten, the best algorithms for protein folding, even though many of the people participating are not experts.

- *Specialized platforms.* Some crowdsourcing systems have been developed and deployed to solve specialized types of problems. One example is Ushahidi,⁸ which provides geographic-based information and visualizations for crisis response and other applications. Another such system is Galaxy Zoo,⁹ which enables people to help identify interesting objects in astronomical images. Galaxy Zoo learns the skill sets of its participants over time and uses this knowledge to route particular images to the people who are most likely to accurately detect the phenomena in those images.
- *Collaborative analysis.* This class of systems consists of the crowdsourcing platforms that are perhaps the most directly related to data analytics at present. Such systems enable groups of people to share and discuss data and visualizations in order to detect and understand trends and anomalies. Such systems typically include a social component in which participants can directly engage each other. Examples of such systems include ManyEyes, Swivel, and Sense.us.

As can be seen from the above list, there is currently a tremendous amount of interest in and innovation around crowdsourcing in many forms. In some cases (e.g., crowdsourced query processing and collaborative analysis) crowd resources are being directly used to help make sense of data. In other cases, there is simply the potential for doing so. The next section outlines opportunities and challenges for developing hybrid human/computer analytics systems, as well as the two other areas of human interaction with data discussed above.

OPPORTUNITIES, CHALLENGES, AND DIRECTIONS

Data Visualization and Exploration

Many of the current challenges in visualization and exploration stem from scalability issues. As the volume of data to be analyzed continues to increase, it becomes increasingly difficult to provide useful visual represen-

⁷ The FoldIt website is available at <http://fold.it>.

⁸ The Ushahidi website is available at <http://ushahidi.com>.

⁹ The Galaxy Zoo website is available at <http://www.galaxyzoo.org>.

tations and interactive performance for massive data sets. These concerns are not unrelated: interactive analysis is qualitatively different from off-line approaches, particularly when exploration is required.

Aggregation strategies and visual representations are gaining importance as research topics (Shneiderman, 2008; Elmquist and Fekete, 2010). This is especially true for network visualization, in which billion-node communications or citation graphs are common and petabyte-per-day growth is a reality (Elmquist et al., 2008; Archambault et al., 2011).

In terms of performance, one would expect that the significant continuing changes in hardware architectures provide an opportunity to address the scalability issue. One appealing research direction is to support massive information visualization by way of specialized hardware. Graphics processing units (GPUs) have become low-cost and pervasive for showing three-dimensional graphics, while other emerging technologies such as data parallel computation platforms and cloud computing infrastructures must also be exploited.

A second area that requires attention is the integration of visualization with statistical methods and other analytic techniques in order to support discovery and analysis. Here, the best strategy appears to lie in combining statistical methods with information visualization (Perer and Shneiderman, 2009). Users can view initial displays of data to gain provisional insights about the distributions, identify errors or missing data, select interesting outliers or clusters, and explore high and low values. At every point they can apply statistical treatments to produce new intermediate data sets, record their insights, select groups for later analysis, or forward promising partial results to colleagues. Often users will need to combine data from several sources and apply domain knowledge to interpret the meaning of a statistical result and visual display. Although the products of an analysis may be compelling displays, the key outcome is insight about the data.

An additional requirement that arises from the interactive nature of many data analysis processes is the need for the analytics system to provide human-understandable feedback to explain analytics results and the steps taken to obtain them. For example, sometimes automated systems produce models that are difficult to understand. Currently, the understandability of the analytical processes is the biggest impediment to using such techniques in decision-making. No CEO is going to make a risky decision using a model they do not understand. Consumers also have problems when automated systems present data they do not understand. Embedding data in a semantic substrate and allowing people to ask high-level questions using interactive tools is an effective way to improve confidence in and utility of an analytics system.

A final area of opportunity is support for group-based analytics. Complex decisions are often made by groups rather than individuals. As data

become more complex, support for groups and shared expertise becomes even more important. Visual analytics researchers emphasize the social processes around information visualization, in which teams of 10 to 5,000 analysts may be working on a single problem, such as pharmaceutical drug discovery, oil/gas exploration, manufacturing process control, or intelligence analysis. These teams must coordinate their efforts over weeks or months, generate many intermediate data sets, and combine their insights to support important decisions for corporations or government agencies.

Crowdsourced Data Acquisition and Hybrid Human/Computer Data Analysis

The other two ways that people can participate in the analytics process are by helping to acquire data and by adding human intelligence where existing algorithms and systems technology cannot provide an adequate answer. These two topics are combined because they share many open issues.

One of the main research problems for crowdsourcing is the need to understand, evaluate, and improve the quality of answers obtained from people. Answers from the crowd can be subject to statistical bias, malicious or simply greedy intent (particularly when work is done for pay), or simply incorrect answers due to a lack of expertise. Such problems are exacerbated in some crowdsourcing systems where workers are more or less anonymous and, hence, not fully accountable, and in environments where monetary incentives are used, which can lead to contributors providing large numbers of random or simply incorrect answers. While many traditional statistical tests and error adjustment techniques can be brought to bear on the problem, the environment of crowdsourced work provides new challenges that must be addressed.

Another important area requiring significant work is the design of incentive mechanisms to improve the quality, cost, and timeliness of crowdsourced contributions. Incentive structures currently used include money, status, altruism, and other rewards. Also, because many crowdsourcing platforms are truly global markets, there are concerns about minimum wages, quality of work offered, and potential exploitation of workers that must be addressed.

Participatory sensing provides another set of research challenges. Large-scale sensing deployments can create massive streams of real-time data. These data can be error-prone and context sensitive. Privacy issues must also be taken into account if the sensing is being done based on monitoring individuals' activities. Finally, the sheer volume of data collected can stress even the most advanced computing platforms, particularly if data are to be maintained over long time periods.

An interesting and important problem is that of determining what types of problems are amenable to human solution as opposed to computer

solution. This question is related to the question of “AI Completeness” as described by (Shahaf and Amir, 2007). It also leads to what is likely the most important area of future work regarding crowdsourcing and analytics, namely, the design and development of hybrid human/computer systems that solve problems that are too hard for computers or people to solve alone. Designing such systems requires a deep understanding of the relative strengths and weaknesses of human and machine computation.

Given the scale of massive data, it makes sense to try to use computers wherever possible, because people are inherently slower for large number-crunching tasks and their abilities are less scalable. Thus statistical methods and machine-learning algorithms should be used when they can produce answers with sufficient confidence within time and budget constraints. People can be brought to bear to handle cases that need additional clarification or insight. Furthermore, human input can be used to train, validate, and improve models. In the longer term, the expectation is that machines and algorithms will continue to improve in terms of the scale and complexity of the tasks they can accomplish. However, it is likely also that this improvement will lead to an increase in the scope, complexity, and diversity of the analysis questions to be answered. Thus, while the roles and responsibilities of the algorithmic and human components will shift over time, the need to design and develop effective hybrid systems will remain.

REFERENCES

- Archambault, D., T. Munzner, and D. Auber. 2011. Tugging graphs faster: Efficiently modifying path-preserving hierarchies for browsing paths. *IEEE Transactions on Visualization and Computer Graphics* 17(3):276-289.
- Elmqvist, N., and J.-D. Fekete. 2010. Hierarchical aggregation for information visualization: Overview, techniques, and design guidelines. *IEEE Transactions on Visualization and Computer Graphics* 16(3):439-454.
- Elmqvist, N., D. Thanh-Nghi, H. Goodell, N. Henry, and J.-D. Fekete. 2008. ZAME: Interactive Large-Scale Graph Visualization. Pp. 215-222 in *Proceedings of the IEEE Pacific Visualization Symposium 2008* (PacificVIS '08), doi:10.1109/PACIFICVIS.2008.4475479.
- Estrin, D. 2010. Participatory sensing: Applications and architecture. *IEEE Internet Computing* 14(1):12-42.
- Franklin, M.J., D. Kossmann, T. Kraska, S. Ramesh, and R. Xin. 2011. CrowdDB: Answering queries with crowdsourcing. Pp. 61-72 in *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2011*. Association for Computing Machinery, New York, N.Y.
- Grier, D. 2005. *When Computers Were Human*. Princeton University Press, Princeton, N.J.
- Little, G., L.B. Chilton, M. Goldman, and R.C. Miller. 2010. TurKit: Human computation algorithms on Mechanical Turk. Pp. 57-66 in *Proceedings of the 23rd Annual ACM Symposium on User Interface Software and Technology*. Association for Computing Machinery, New York, N.Y.

- Marcus, A., E. Wu, S. Madden, and R.C. Miller. 2011. Crowdsourced databases: Query processing with people. Pp. 211-214 in *Proceedings of the 2011 Conference on Innovative Data Systems Research (CIDR)*. Available at <http://www.cidrdb.org/cidr2011/program.html>, pp. 211-214.
- Perer, A., and B. Shneiderman. 2006. Balancing systematic and flexible exploration of social networks. *IEEE Transactions on Visualization and Computer Graphics* 12(5):693-700.
- Perer, A., and B. Shneiderman. 2009. The importance of integrating statistics and visualization: Long-term case studies supporting exploratory data analysis of social networks. *IEEE Computer Graphics and Applications* 29(3):39-51.
- Quinn, A.J., and B.B. Bederson. 2011. Human computation: A survey and taxonomy of a growing field. Pp. 1403-1412 in *Proceedings of the 2011 SIGCHI Conference on Human Factors in Computing*. Association for Computing Machinery, New York, N.Y.
- Shahaf, D., and E. Amir. 2007. Towards a theory of AI completeness. Presented at COMMONSENSE 2007: 8th International Symposium on Logical Formalizations of Commonsense Reasoning. In *Logical Formalizations of Commonsense Reasoning: Papers from the AAI Spring Symposium*. AAI Technical Report SS-07-05. AAI Press, Menlo Park, Calif. Available at <http://www.ucl.ac.uk/commonsense07/papers/>.
- Shahaf, D., and E. Horvitz. 2010. Generalized task markets for human and machine computation. Pp. 986-993 in *Proceedings of the 24th AAI Conference on Artificial Intelligence*. Association for the Advancement of Artificial Intelligence, Palo Alto, Calif.
- Shneiderman, B. 2008. Extreme visualization: Squeezing a billion records into a million pixels. Pp. 3-12 in *Proceedings of the ACM SIGMOD 2008 International Conference on the Management of Data*. Association for Computing Machinery, New York, N.Y. Available at <http://portal.acm.org/citation.cfm?doid=1376616.1376618>.
- Shneiderman, B., and C. Plaisant. 2006. Strategies for evaluating information visualization tools: Multi-dimensional in-depth long-term case studies. Pp. 1-7 in *Proceedings of the 2006 Advanced Visual Interfaces Workshop on Beyond Time and Errors: Novel Evaluation Methods for Information Visualization*. Association for Computing Machinery, New York, N.Y.
- Thomas, J.J., and K.A. Cook, eds. 2005. *Illuminating the Path: Research and Development Agenda for Visual Analytics*. IEEE Press. Available at <http://nvac.pnl.gov/agenda.stm>.
- von Ahn, L., and L. Dabbish. 2008. Games with a purpose. *Communications of the ACM* 51(8):58-67.