

# eBird: A Human/Computer Learning Network to Improve Biodiversity Conservation and Research

**Steve Kelling, Jeff Gerbracht, and Daniel Fink**

Cornell Lab of Ornithology, Cornell University  
stk2@cornell.edu, jag73@cornell.edu, df17@cornell.edu

**Weng-Keen Wong and Jun Yu**

School of EECS, Oregon State University  
wong@eecs.oregonstate.edu, juyu@eecs.oregonstate.edu

**Carl Lagoze**

Information Science, Cornell University  
lagoze@cs.cornell.edu

**Theodoros Damoulas and Carla Gomes**

Department of Computer Science, Cornell University  
damoulas@cs.cornell.edu, gomes@cs.cornell.edu

## Abstract

In this paper we describe eBird, a citizen-science project that takes advantage of the human observational capacity to identify birds to species, which is then used to accurately represent patterns of bird occurrences across broad spatial and temporal extents. eBird employs artificial intelligence techniques such as machine learning to improve data quality by taking advantage of the synergies between human computation and mechanical computation. We call this a Human/Computer Learning Network, whose core is an active learning feedback loop between humans and machines that dramatically improves the quality of both, and thereby continually improves the effectiveness of the network as a whole. In this paper we explore how Human/Computer Learning Networks can leverage the contributions of a broad recruitment of human observers and processes their contributed data with Artificial Intelligence algorithms leading to a computational power that far exceeds the sum of the individual parts.

## Introduction

The transformational power of today's computing, together with information and communication technologies, are providing new opportunities to engage the public to participate in and contribute to a myriad of scientific, business and technical challenges. For example, citizen-science projects such as Galaxy Zoo, eBird, and FoldIt demonstrate the power of crowdsourcing for investigating large-scale scientific problems. These and similar projects leverage emerging techniques that integrate the speed and scalability of *mechanical computation*, using advances in Artificial Intelligence (AI), with the real intelligence of *human computation* to solve

computational problems that are beyond the scope of existing algorithms [1].

Human computational systems use the innate abilities of humans to solve certain problems that computers cannot solve [2]. Now the World Wide Web and wireless handheld devices provide the opportunity to engage large numbers of humans to solve these problems. For example, engagement can be game-based such as FoldIt, which attempts to predict the structure of a protein by taking advantage of humans' puzzle solving abilities [3]; or Galaxy Zoo, which has engaged more than 200,000 participants to classify more than 100 million galaxies [4]. Alternatively, the Web can be used to engage large numbers of participants to actively collect data and submit it to central data repositories. Projects such as eBird, engage a global network of volunteers to report bird observations that are used to generate extremely accurate estimates of species distributions [5].

Now systems are being developed that employ both human and mechanical computation to solve complex problems through active learning and feedback. These Human/Computer Learning Networks (HCLN) can leverage the contributions of broad recruitment of human observers and process their contributed data with AI algorithms for a resulting total computational power far exceeding the sum of their individual parts. This combination can be deployed in a variety of domains and holds enormous potential to solve complex computational problems.

A key factor in the power of an HCLN is the manner in which the benefits of active learning are cyclically fed back among the human participants and computational systems. We use "active learning" in both of its commonly used senses: the *machine learning* sense as a form of iterative supervised learning, and the *human* sense in which learners (our volunteers) are actively and

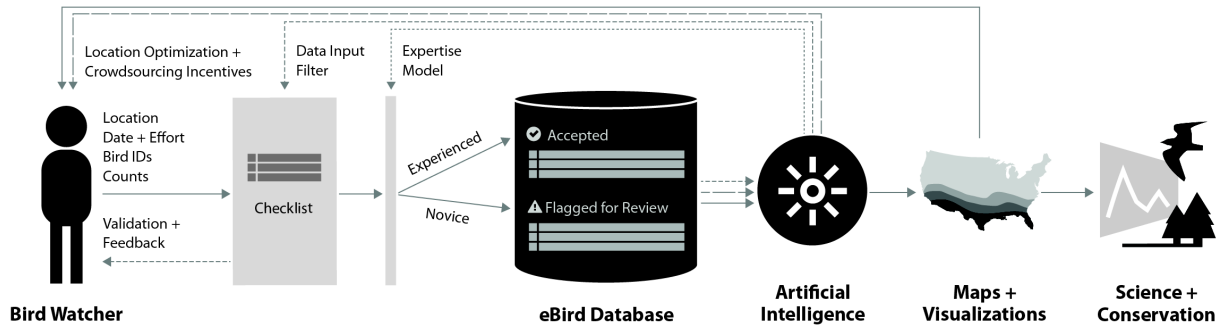


Figure 1. An HCLN example. Human observers and AI processes synergistically improve the overall quality of the entire system. Additionally, AI is used to generate analyses. These analyses also improve as the quantity and quality of the incoming data improves.

dynamically guided to new levels of expertise. The role of active learning in a HCLN is illustrated in Figure 1. In our example, broad networks of volunteers act as intelligent and trainable sensors to gather observations. AI processes dramatically improve the quality of the observational data that volunteers provide by filtering inputs based on aggregated historical data and observers' expertise. By guiding observers with immediate feedback on observation accuracy AI processes contribute to advancing observer expertise. Simultaneously, as observer data quality improves, the training data on which the AI processes make their decisions also improves. This feedback loop increases the accuracy of the analysis, which enhances the general utility of the data for scientific purposes.

A successful HCLN must be able to address the 4 following challenges. First, a task must be identified that human computational systems can complete but mechanical computational systems cannot [1]. Second, the task must be sufficiently straightforward and incentivized to maximize participation [6]. Third, the complimentary abilities of both humans and machines must be clearly identified so that they can be leveraged to increase the accuracy and efficiency of the network [7]. Finally novel methods for extracting biological insights from the noisy and complex data provided by multiple human computers must be employed [8]. In this paper we use our experience with eBird as a model to address these 4 HCLN challenges.

### Challenge 1: Species Identification

Few mechanical computational systems have been developed to classify organisms to the species level. Those that do exist typically can only identify a single or small group of species, and cannot classify a multitude of organisms. Only human observers can reliably identify organisms to the species level [9], and are capable of classifying hundreds of species. This is because identifying a species is a complex task that relies on a combination of factors. First, observers must be able to process impressions of shape, size, and behavior under variable observation conditions. As this process continues, the

observer must combine these impressions with a mental list of species most likely to occur at that specific location and date until the species is correctly identified.

eBird (<http://ebird.org>) [5] is a citizen science project that engages a global network of bird watchers to identify birds to species and report their observations to a centralized database. Anyone can submit observations of birds to eBird via the web or wireless handheld devices (e.g. iPhone and Android). To date more than 91,000 individuals have volunteered more than 4 million hours and collected over 100 million bird observations; arguably the largest biodiversity data collection project in existence. These amassed observations provide researchers, scientists, students, educators, and amateur naturalists with data about bird distribution and abundance across varying spatio-temporal extents. Dynamic and interactive maps, graphs and other visualizations are available on the eBird website, and all data are readily accessible through the Avian Knowledge Network [10]. Since 2006 eBird data have been the used in more than 60 peer-reviewed publications and reports, from highlighting the importance of public lands in conservation [11], to studies of evolution [12], climate change [13] and biogeography [14].

### Challenge 2: Maximizing Participation

eBird uses crowdsourcing techniques to engage a large numbers of people to perform tasks that automated sensors and computers cannot readily accomplish [15]. This is accomplished through the development of straightforward rules for participation and incentives for contributing. Initially, the incentive for participating in eBird was to help scientists study birds, which led to very disappointing participation in eBird. Recognizing this, project managers changed the emphasis of the project from having birders help scientists, to tools that appealed to the birding community. New features were developed for eBird that allowed participants to: (i) keep track of their bird records; (ii) sort their personal bird lists by date and region; (iii) share their lists with others; and (iv) visualize their observations on maps and graphs. By providing these

record-keeping, exploration, and visualization facilities as a direct reward for participation eBird participation has grown exponentially (Figure 2). eBird appeals to the competitiveness of participants, and through the further development of eBird more interactive and varied tools allowed participants to determine their relative status compared to other participants (e.g. numbers of species seen) and geographical regions (e.g. checklists submitted per state and province). Thus, by changing the emphasis of eBird to one that supports the needs and desires of the birding community, growth in eBird has been exponential (Figure 2). For example, more data were gathered in May 2012, than during the first 3 years of the project.

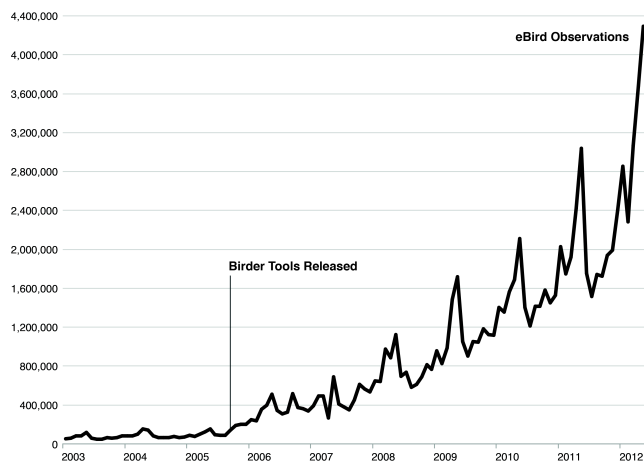


Figure 2. The number of observations submitted monthly to eBird since its inception in late 2003. Tools to better engage the bird watching community were released in mid-2005. Note the annual peaks of submission, which occur each May, when spring migration is at its peak and birders are most active.

An additional key component of eBird's success has been the implementation of a sound data management strategy, which reduces the risk of data loss and allows for efficient use and re-use of the data. All eBird data contain the following information: observer identification, location, visit, and what was collected. These data form the core observational data model [16] and provide the opportunity for integration, visualization, experimentation and analysis. For example, eBird collects the name and contact information for every observer, which allows each observation to be attributed to a specific person. Location data such as the site name the coordinates where the observations were made and the geographic area represented by the location are stored with every visit to that location. Information about a specific visit consists of data and time of visit, amount of effort expended, such as distance traveled, time spent and area covered, and whether

or not all species observed were reported. Species observations consist of a checklist of birds observed and how many individuals of each species were counted.

### Challenge 3: Identifying the Synergies Between Humans and Machines

While eBird has been successful in engaging a large community of volunteers to contribute large quantities of observations of birds, there are many challenges to using eBird data for analysis. First, observers are bound to misidentify birds, which is the largest source of error in the eBird database. Second, there is much variability in a participant's ability to identify birds, with some eBird contributors being experts in bird identification, while others are novices. Third, participation in eBird is not uniformly distributed in space; most eBird observations occur in regions where human population densities are fairly high. Improving eBird data quality is a constant and major effort. This is because as data quality improves the accuracy in estimating patterns of bird occurrence also improves. In this section we describe how the implementation of HCLN processes allow us address these 3 data quality issues.

#### How can we efficiently filter erroneous data before data enter the database?

Data quality is a major issue for eBird, particularly as it pertains to an observer's ability to correctly identify birds to the species level. While eBird has motivated tens of thousands of volunteers to collect large amounts of data at relatively little cost, the misidentification of birds is a major concern. Since its inception eBird has employed a data validation system that relied heavily on a network of volunteers who were experts in the patterns of bird occurrence. However, beginning in 2010, the sheer volume of data being gathered had overwhelmed the volunteer network of record reviewers. Initially records had been filtered regionally (i.e., country, state, county), and temporally at the monthly scale. The basic filter mechanism assigned a specific region with a value for a given month, which corresponded to an expert's opinion for a maximum allowable count for a given region. If a submission exceeded the maximum allowable amount, it was "flagged" for review by one of more than 450 volunteer reviewers. Reviewers contacted those individuals who submitted flagged records to obtain additional information, such as field notes or photographs, in order to confirm unusual reports. In 2010, 4% or 720,000 of the 23 million records submitted to eBird were reviewed. This number put a severe strain on the volunteer network, with many reviewers complaining they were overwhelmed from the sheer volume of records to review.

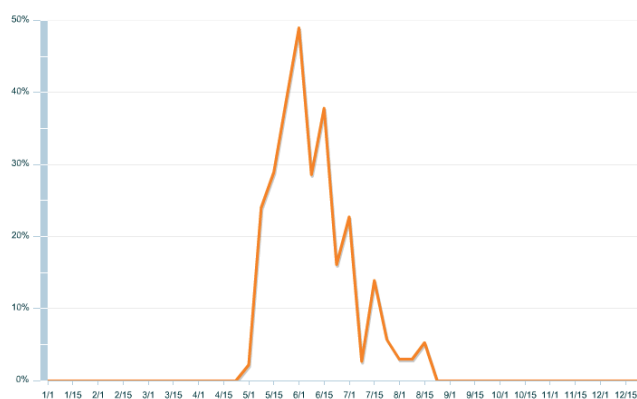


Figure 3. Frequency of occurrence results for Black-billed Cuckoos (*Coccyzus erythrophthalmus*) in upstate New York. The Y-axis is the frequency of eBird checklists that reported this species, and the X-axis is the date. Cuckoos arrive in early May and are detected at high frequencies because they are conspicuous and vocal during their mating season. But after they lay eggs, their detection probability drops dramatically. Most birds leave by mid-August.

In order to decrease the volume of data that needed to be reviewed by the experts, we have implemented a new data quality filter and screening process that automates much of the review process. One of the most powerful calculations performed on eBird data is the frequency in which a particular bird species was reported during a particular period of time (Figure 3). Since each observation contains details of where and when a bird was detected, we can estimate the “likelihood” of observing a specific species at any spatial level (e.g., country, state, county, backyard, or any spatial extent of interest) and for any date. Frequency filters delineate when a species can be reported in a region and determines the validity of an observation.

The eBird database currently holds more than 100 million bird observations. These historical records can be used to filter unusual observations that require review, but allow entry of expected species within the expected times when species should occur. These filters automatically *emerge* from historic eBird data. We have set the emergent filter at 10% of maximum annual frequency of occurrence for every species across the United States. This provides a consistent limit that allows expected observations through the filter but flags for review unusual records. For example, if a common species reaches a maximum frequency of 68% then the filter would identify the day when the filter first crosses the 6.8% threshold. Any record submitted on a date either prior or after the threshold limit, is flagged for review. Similarly, if a rare species reaches an annual peak of 6.5% frequency, the threshold limit would

be .65%. Table 1 shows the number of flagged records the emergent filter identifies for 2 counties in New York State, Jefferson Co. and Tompkins Co (Table 1). These 2 counties were selected because Jefferson Co. has relatively sparse year-round data coverage, while Tompkins Co. is one of the most active regions in eBird.

When the emergent filter is triggered the submitter gets immediate feedback indicating that this was an unusual observation (Figure 1). If they confirm they made the observation, their record is flagged for review, and one of the volunteer experts will review the observation. All records, their flags and their review history are retained in the eBird database.

What is most significant about the emergent filter process is that it identifies key periods during a bird’s life history when their patterns of occurrence change (e.g. during periods of migration when the bird either arrives or departs a specific region). Figure 4 shows those records that are flagged for review by the emergent filter for the 2 New York Counties. The Chipping Sparrow (*Spizella passerina*) is a common breeding bird in upstate New York, but departs the region in the fall and rarely occurs in winter. The emergent filter for each county is different, due to the variation in each county’s respective historic data. The triangles and circles are all records that are flagged for review by the emergent filter. Without the emergent filter it would be difficult to accurately identify arrival and departure dates of when a bird appears in a county. The threshold of occurrence established by the emergent filter allows the determination of arrival and departure and then accurately flags outlier observation for further processing and review.

	Tompkins Co.	Jefferson Co.
Total Observations	704,053	78,745
Total Flagged	50,743	6,082
Percent Flagged	7	8
Total Flagged Expert	38,574	3,787
Total Flagged Novice	12,170	2,295
Percent Expert	5	5
Percent Novice	2	3

Table 1. Results of the Emergent Filter process applied to 2 counties in Upstate New York (upper), and the proportion of flagged records submitted by experts and novices (lower).

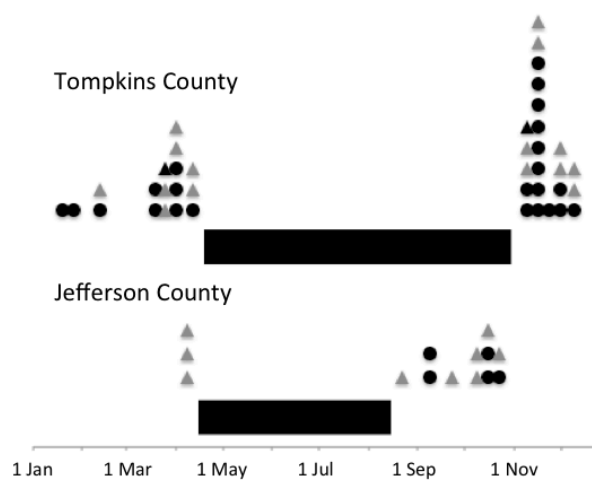


Figure 4. The acceptable date range (dark bars) for the occurrence of Chipping Sparrow in 2 counties in New York. All records that fall outside of the acceptable date range are plotted either as circles (novices) or triangles (experts).

### Can we identify observer variability in their ability to detect objects?

eBird data are contributed by observers with a wide range of expertise in identifying birds. At one extreme observers with high identification skill levels contribute “professional grade” observations to eBird, whereas at the other extreme less-skilled participants contribute data of more variable quality. This inter-observer variation must be taken into account during analysis to determine if outlier observations (i.e., those observations that are unusual) are true occurrences of a rare species, or the misidentification of a common species. Since eBird engages a significant number of skilled observers who are motivated to detect rare species or are skilled in detecting elusive and cryptic species, being able to automatically and accurately distinguish their observations from those of less-skilled observers is crucial. This is because skilled observers are more likely to submit observations of unusual species that get flagged by the regional emergent filters (i.e., skilled birders like to find rare birds). What is required is an objective measure of observer expertise that would automatically classify unusual observations.

To better understand observer variability in eBird we have applied a probabilistic machine learning approach called the Occupancy-Detection-Experience (ODE) model to provide an objective measure of expertise for all eBird observers [17]. The ODE model extends existing ecological models that measure the viability of a site as

suitable habitat for a species, by predicting site occupancy by a particular species.

We can use the ODE model to distinguish the difference between expert observers, who will find more birds and are more likely to find them outside of the emergent filter limits, and novice birders, who are more likely to misidentify common birds. Table 1 (bottom) shows the total number of observations by experts and novices that are flagged. As expected, expert observers had a greater number of flagged records, because of their enhanced bird identification skills, and their desire to find unusual birds. We can use the ODE model results for experts in the data filtering process by automatically accepting their expert observations, which dramatically reduces the total number of flagged records that need to be reviewed (Table 1 bottom). Finally, to test the accuracy of the ODE model we analyzed all observations that fell outside of the emergent filter for more than a dozen species that easily confuse novices, and show results for Chipping Sparrow (Figure 4). We did this by engaging the current reviewers for the 2 counties in New York, who confirmed that the ODE model properly categorized the observer as either an expert or novice and validated more than 95% of the expert observations that fell outside of the emergent filters.

We have found that the combination of the emergent checklist filters with the ODE model provides the best strategy for both improving data quality and streamlining the review process in eBird. This two-step approach, where the emergent data filters are used to identify outliers, and the ODE model allowed us to identify valid outliers, identifies unusual records more accurately than previous methods. The result is that we can now provide accurate occurrence probabilities, which are based on existing eBird data to allow the quick identification and classification of outliers.

### How can we address the spatial bias in citizen-science projects?

An inherent liability with many citizen-science projects is that observation locations are highly biased towards regions with high human populations. If this inequity is ignored, the spatial bias will produce results in which regions with the most data have excessive influence on the overall results accuracy and regions with the least data are under represented [8]. We address this issue using a mediated optimization strategy to identify areas that if



sampled would most improve eBird spatial coverage and improve analysis results.

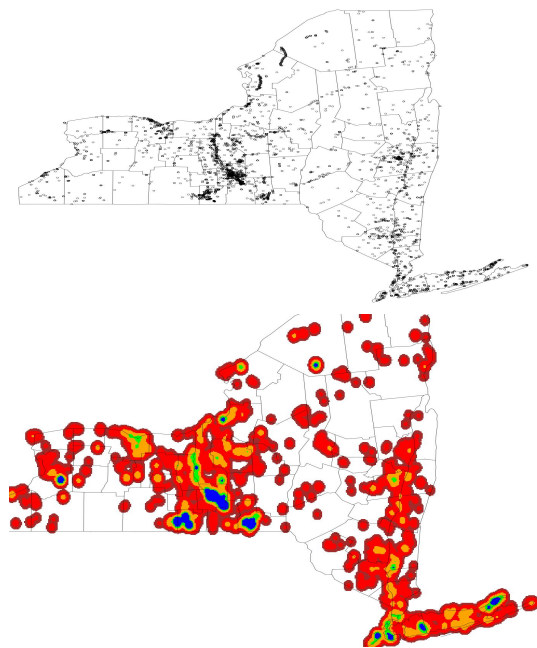


Figure 5. Top: locations in New York where submissions were made to eBird in 2009. Bottom: Results showing areas with sufficient data density (colored regions) and those requiring more data (white regions).

Machine learning algorithms can improve the predictive performance of eBird by guiding the sampling process. Consider the locations where eBird observations were made in New York (Figure 5 top). It is clear that spatial sampling biases are present as the majority of the observations come from a small subset of geographical locations. Active learning applied to eBird improve the resulting predictive models by providing a context to advise participants *where to sample next*. A first strategy, as displayed in Figure 5 (bottom), has been to aim for a uniform sampling coverage in geographical space, by concentrating data collection efforts to the areas of highest model uncertainty and low density. This is accomplished through a novel active learning approach that combines density information and information-theoretic measures [19].

Already, our research in offering optimal sampling strategies is paying off. We display maps similar to Figure 5 (bottom) on the eBird website, and provide rewards for individuals who report checklists from under sampled regions. Eventually, such sampling trajectories will be employed throughout eBird, to enhance the overall birding experience. For example, it is straightforward to propose

paths that have the highest probability of detecting birds. Hence one can envision educating observers by proposing appropriate paths that trains their detection capabilities on specific species or increases the probability of them recording a species they have never observed before.

#### Challenge 4: Species Distribution Models

The effective management and conservation of biodiversity requires knowledge of a species' geographic distribution throughout the year. Until the inception of eBird, detailed data documenting a species' distribution were often not available for the entire range of a species, particularly for widely distributed species or species not closely studied. eBird provides broad-scale survey data that allows researchers to analyze and interpret a specie's distribution across broad spatial extents and for any time of year.

One major area of analysis of eBird data is to explore the continent-wide inter-annual patterns of occurrence of North American birds. To do this we have developed new Spatial-temporal Exploratory Models (STEM) of species distributions, that allow us to automatically discover patterns in spatiotemporal data [8].

We designed our statistical models specifically to discover seasonally- and regionally-varying patterns in eBird data. Spatiotemporal variation in habitat associations are captured by combining a series of separate submodels, each describing the distribution within a relatively small area and time window. The approach is semiparametric; yielding a highly automated predictive methodology that allows an analyst to produce accurate predictions without requiring a detailed understanding of the underlying dynamic processes. This makes STEMs especially well suited for exploring distributional dynamics arising from a variety of complex dynamic ecological and anthropogenic processes. STEMs can be used to study how spatial distributions of populations respond over time, both seasonally (Figure 6) as well as to broad-scale changes in their environments (i.e., changes in land-use patterns, pollution patterns, or climate change).

The STEM visualizations are now being employed in a number of research and conservation initiatives. For example, bird distribution information used in the *2011 State of the Birds Report* prepared for the U. S. Department of Interior by the North American Bird Conservation Initiative (NABCI), was based on STEM model results. Additionally, other federal (i.e., Bureau of Land Management and U.S. Forest Service) and non-governmental agencies (i.e., The Nature Conservancy) are using STEM distribution estimates to study placement of

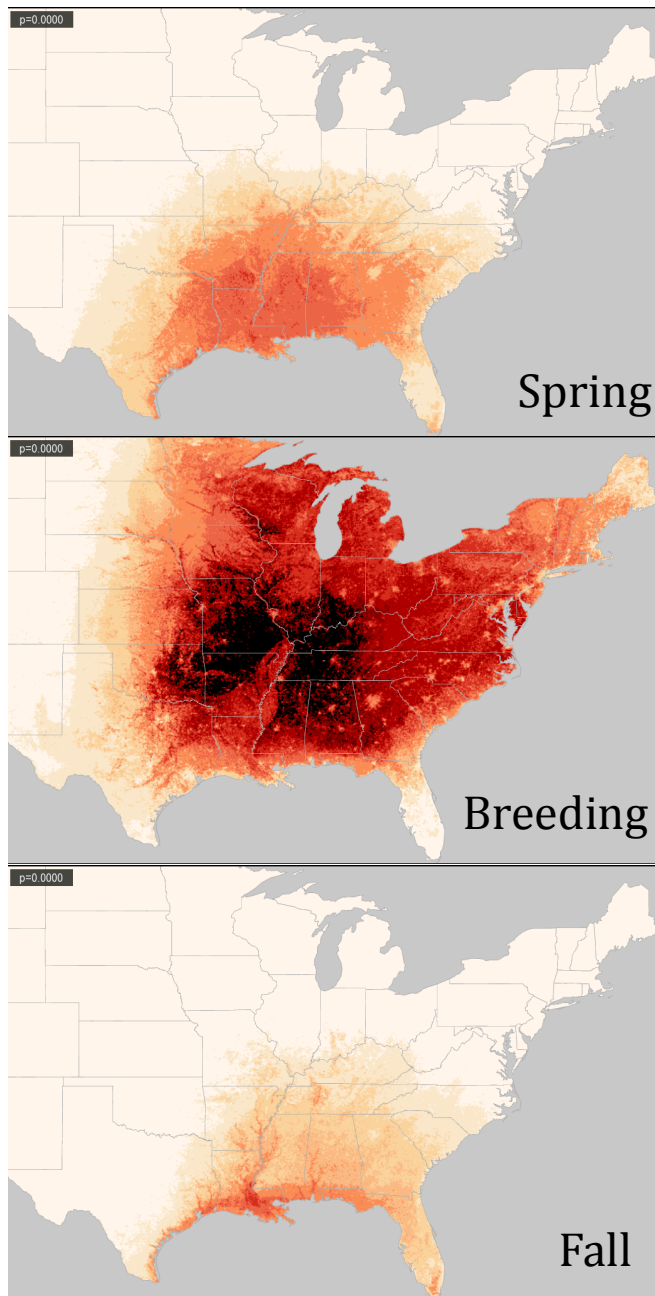


Figure 6. This series of maps illustrate the seasonal patterns of occurrence of the Indigo Bunting (*Passerina cyanea*) throughout the United States. The maps illustrate the seasonal distribution estimates from a spatiotemporal exploratory model (STEM) during spring migration (top), the breeding season (middle) and during fall migration (bottom). Indigo Buntings are Neotropical migrants, wintering in Central America and returning to the United States annually to breed. The occurrence maps show the probability of encountering the species on an early morning 1-hour birding walk, with darker colors indicating higher probabilities. These maps provide continental-scale distribution estimates that allow the quick assessment of the rate of arrival and departure from breeding grounds, and over time will allow researchers to identify and quantify changes in bird populations. More STEM maps can be viewed on the eBird website (<http://www.ebird.org>).

wind farms for sustainable energy production, identifying and prioritizing areas for avian conservations and the Pacific Northwest.

## Conclusion

In this paper, we have demonstrated the implementation of a novel network that links machine learning methods and human observational capacity to address several unique challenges inherent in a broad-scale citizen-science project. By exploring the synergies between mechanical computation and human computation, which we call a Human/Computer Learning Network we can leverage emerging technologies that integrate the speed and scalability of AI, with human computation to solve computational problems that are currently beyond the scope of existing AI algorithms.

eBird uses a broad-scale survey design to maximally engage volunteers to gather bird observations following a basic protocol for data collection. Designing such broad-scale surveys to maximize the information obtained for use in analysis depends on finding the proper balance between data quantity and data quality. If we can engage a large number of participants to collect data through eBird's very basic protocols a sufficiently large volume of data can be gathered and effectively analyzed. While eBird data has relatively lower per-datum information content, we have found that eBird data can contain more information for broad-scale distribution estimates than a smaller amount of data with higher per-datum quality [20].

The appropriate design of data input and management procedures is critical to maintain the balance between data quantity and data quality in broad-scale citizen science projects. The additional implementation of novel AI functionality provides incentives for encouraging surveyors to contribute even more data while simultaneously limiting errors and providing opportunities for dramatically improved data review and validation procedures.

Although our discussion has focused on one citizen-science project, eBird, the general HCLN approach are more widely applicable. Specifically, by implementing an uncomplicated protocol via web-based and wireless handheld devices and providing appropriate rewards for participation, citizen-science projects can recruit large numbers of participants to submit massive quantities of meaningful data. By taking an adaptive learning approach

for both humans and computers we can improve the quality and scope of the data that the volunteers provide. Finally, new analysis techniques that bridge the gap between parametric and non-parametric processes provide extremely accurate estimates of species occurrence at continental levels.

In conclusion, broad-scale citizen-science projects can recruit extensive networks of volunteers, who act as intelligent and trainable sensors in the environment that gather observations across broad spatial (e.g., globally) and temporal (e.g., any time) extents. However, there is much variability in the observations volunteers make. Artificial Intelligence processes can dramatically improve the quality of the observational data by filtering inputs using emergent filters based on aggregated historical data, and on the observers' expertise. By guiding the observers with immediate feedback on observation accuracy, the Artificial Intelligence processes contribute to advancing expertise of the observers, while simultaneously improving the quality of the training data on which the Artificial Intelligence processes make their decisions. The outcome is improved data quality that can be used for research and analysis.

### Acknowledgments

This work was funded by the Leon Levy Foundation, Wolf Creek Foundation and the National Science Foundation (Grant Numbers OCI-0830944, CCF-0832782, ITR-0427914, DBI-1049363, DBI-0542868, DUE- 0734857, IIS-0748626, IIS-0844546, IIS-0612031, IIS-1050422, IIS-0905385, IIS-0746500, AGS-0835821, CNS-0751152, CNS-0855167).

### References

- [1] E. Law and L. v. Ahn, "Human Computation," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 5, pp. 1-121, 2011/06/30 2011.
- [2] Y. Man-Ching, C. Ling-Jyh, and I. King, "A Survey of Human Computation Systems," in *2009 International Conference on Computational Science and Engineering*, 2009, pp. 723-728.
- [3] S. Cooper, F. Khatib, A. Treuille, J. Barbero, J. Lee, M. Beenen, A. Leaver-Fay, D. Baker, Z. Popovic, and F. players, "Predicting protein structures with a multiplayer online game," *Nature*, vol. 466, pp. 756-760, 2010.
- [4] C. J. Lintott, K. Schawinski, S. Anze, K. Land, S. Bamford, D. Thomas, M. J. Raddick, R. C. Nichol, A. Szalay, D. Andreescu, P. Murray, and J. Vandenberg, "Galaxy Zoo: morphologies derived from visual inspection of galaxies from the Sloan Digital Sky Survey," *Monthly Notices of the Royal Astronomical Society*, vol. 389, pp. 1179-1189, 2008.
- [5] B. L. Sullivan, C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling, "eBird: A citizen-based bird observation network in the biological sciences," *Biological Conservation*, vol. 142, pp. 2282-2292, 2009.
- [6] C. Wood, B. Sullivan, M. Iliff, D. Fink, and S. Kelling, "eBird: Engaging Birders in Science and Conservation," *PLoS Biol*, vol. 9, p. e1001220, 2011.
- [7] S. Kelling, J. Yu, J. Gerbracht, and W. K. Wong, "Emergent Filters: Automated Data Verification in a Large-scale Citizen Science Project," *Proceedings of the IEEE eScience 2011 Computing for Citizen Science Workshop. To Appear.*, 2011.
- [8] D. Fink, W. M. Hochachka, D. Winkler, B. Shaby, G. Hooker, B. Zuckerberg, M. A. Munson, D. Sheldon, M. Riedewald, and S. Kelling, "Spatiotemporal Exploratory models for Large-scale Survey Data," *Ecological Applications*, vol. 20, pp. 2131-2147, 2010.
- [9] W. M. Hochachka, R. Caruana, D. Fink, A. Munson, M. Riedewald, D. Sorokina, and S. Kelling, "Data-mining discovery of pattern and process in ecological systems," *Journal of Wildlife Management*, vol. 71, pp. 2427-2437, 2007.
- [10] M. J. Iliff, L. Salas, E. R. Inzunza, G. Ballard, D. Lepage, and S. Kelling, "The Avian Knowledge Network: A Partnership to Organize, Analyze, and Visualize Bird Observation Data for Education, Conservation, Research, and Land Management," presented at the Proceedings of the Fourth International Partners in Flight Conference: Tundra to Tropics, McAllen Texas, USA, 2009.
- [11] N. A. B. C. I. U.S., "The State of the Birds 2011 Report on Public Lands and Waters," Washington, D.C.2011.
- [12] J. E. McCormack, A.J. Zellmer, and L. L. Knowles, "Does niche divergence accompany allopatric divergence in *Aphelocoma* jays as predicted under ecological speciation?: insights from tests with niche models.," *Evolution*, pp. 1-13, 2009.



- [13] A. H. Hurlbert and Z. Liang, "Spatiotemporal Variation in Avian Migration Phenology: Citizen Science Reveals Effects of Climate Change," *PLoS ONE*, vol. 7, p. e31662, 2012.
- [14] J. Klicka, G. M. Spellman, K. Winker, V. Chua, and B. T. Smith, "A Phylogeographic and Population Genetic Analysis of a Widespread, Sedentary North American Bird: The Hairy Woodpecker (*Picoides villosus*)," *The Auk*, vol. 128, pp. 346-362, 2011/04/01 2011.
- [15] J. Howe, *Crowdsourcing. Why The Power of the Crowd is Driving the Future of Business*. New York: Crown Business, 2008.
- [16] S. Kelling, "The Significance of Observations in Biodiversity Studies," *GBIF Best Practices Series* (<http://www.gbif.org/communications/news-and-events/showsingle/article/now-available-white-paper-on-significance-of-organism-observations/>), 2008.
- [17] J. Yu, W. K. Wong, and R. Hutchinson, "Modeling experts and novices in citizen science data for species distribution modeling," presented at the IEEE International Conference on Data Mining, Sydney, Australia, 2010.
- [18] D. I. MacKenzie, J. D. Nichols, J. A. Royle, K. H. Pollock, L. L. Bailey, and J. E. Hines, *Occupancy Estimation and Modeling*. Amsterdam: Elsevier, 2006.
- [19] B. Dilkina, T. Damoulas, C. Gomes, and Daniel Fink, AL2: Learning for Active Learning. Workshop "Machine Learning for Sustainability" in the 25<sup>th</sup> Conference of Neural Information Processing Systems (NIPS), Granada, Spain 2011.
- [20] M. A. Munson, R. Caruana, D. F. Fink, W. M. Hochachka, M. I. Iliff, K. V. Rosenberg, D. R. Sheldon, B. L. Sullivan, C. L. Wood, and S. Kelling, "A Method For Measuring the Relative Information Content of Data From Different Monitoring Protocols," *Methods In Ecology and Evolution*, vol. 1, pp. 263-273, 2010.