

Approaches to advance scientific understanding of macrosystems ecology

Ofir Levy^{1*}, Becky A Ball^{2†}, Ben Bond-Lamberty³, Kendra S Cheruvilil⁴, Andrew O Finley⁴, Noah R Lottig⁵, Surangi W Punyasena⁶, Jingfeng Xiao⁷, Jizhong Zhou⁸, Lauren B Buckley⁹, Christopher T Filstrup¹⁰, Tim H Keitt¹¹, James R Kellner¹², Alan K Knapp¹³, Andrew D Richardson¹⁴, David Tcheng⁶, Michael Toomey¹⁴, Rodrigo Vargas¹⁵, James W Voordeckers⁸, Tyler Wagner¹⁶, and John W Williams¹⁷

The emergence of macrosystems ecology (MSE), which focuses on regional- to continental-scale ecological patterns and processes, builds upon a history of long-term and broad-scale studies in ecology. Scientists face the difficulty of integrating the many elements that make up macrosystems, which consist of hierarchical processes at interacting spatial and temporal scales. Researchers must also identify the most relevant scales and variables to be considered, the required data resources, and the appropriate study design to provide the proper inferences. The large volumes of multi-thematic data often associated with macrosystem studies typically require validation, standardization, and assimilation. Finally, analytical approaches need to describe how cross-scale and hierarchical dynamics and interactions relate to macroscale phenomena. Here, we elaborate on some key methodological challenges of MSE research and discuss existing and novel approaches to meet them.

Front Ecol Environ 2014; 12(1): 15–23, doi:10.1890/130019

Many ecological studies are conducted by measuring responses to stressors within populations, communities, or ecosystems. The interactions of basic building blocks of ecological systems – from atoms to organisms, with each other and with the environment – aggregate to form broad-scale ecological patterns. Local-scale research on these interactions is very important for scientists to understand the impacts of environmental change on ecological systems and the processes that shape these phenomena. However, environmental change operates across a range of local to broad scales, forcing ecologists to expand, adapt, and integrate approaches (Heffernan *et al.* 2014).

Improving approaches for prediction is one of the goals

In a nutshell:

- Macrosystems ecology uses new approaches and applies existing methods in novel ways to study ecological processes interacting within and across scales
- These approaches often include multiple scales, diverse data objects, data-intensive methods, cross-scale interactions, and hierarchical relationships
- These studies require large volumes and diverse types of data from many sources, encouraging ecologists to build field and laboratory methods, database objects, and the data infrastructure capable of the joint analysis of multiple large data streams
- Scientists use powerful statistical methods, such as Bayesian hierarchical models, machine learning, and simulations, to find and explain important patterns in complex, multi-scale datasets

of the emerging field of macrosystems ecology (MSE; Heffernan *et al.* 2014). MSE researchers study ecological systems as a whole and ask how processes and patterns at regional to continental scales interact, respond, and emerge from (and with) finer (eg individual) and broader (eg continental) system levels (Peters *et al.* 2007; Evans *et al.* 2012; Heffernan *et al.* 2014). Heffernan *et al.* (2014) describe the important conceptual underpinnings of a macrosystems perspective and its disciplinary foundations. Here, we illustrate the suite of data, approaches, and tools that can be used to address such research questions. Many of these approaches were unavailable 10–20 years ago, and are not commonly used by ecologists today.

The methods we describe here differ from simple upscaling procedures used in early research efforts that laid the foundation for MSE, such as the 1970s International Biological Program, which funded large-scale ecosystem research projects studying the structure and function of key biomes (Hagen 1992; Golley 1993). These studies improved the understanding and development of ecology by refining methods; collecting large amounts of data on ecosystem components, processes, and interactions; and creating many successful, smaller-scale systems models (Golley 1993). However, their upscaling approaches were limited by data resources, analytical tools, and computer capabilities. Ecologists are now able to develop and use technologies to incorporate the complex organization and interactions across scales necessary for interpreting macroscale phenomena (Hagen 1992).

MSE studies explore how broad-scale variation in fine-scale characteristics – such as organismal behavior and fitness, nutrient transformations, and water-use efficiency – relate to broad-scale spatial and temporal processes and patterns such as climate change, landscape alteration, and

¹Arizona State University, Tempe, AZ *(levyofi@gmail.com);

²Arizona State University at the West Campus, Glendale, AZ;

³JGCRI, Pacific Northwest National Lab, College Park, MD;

⁴Michigan State University, East Lansing, MI; continued on p 23

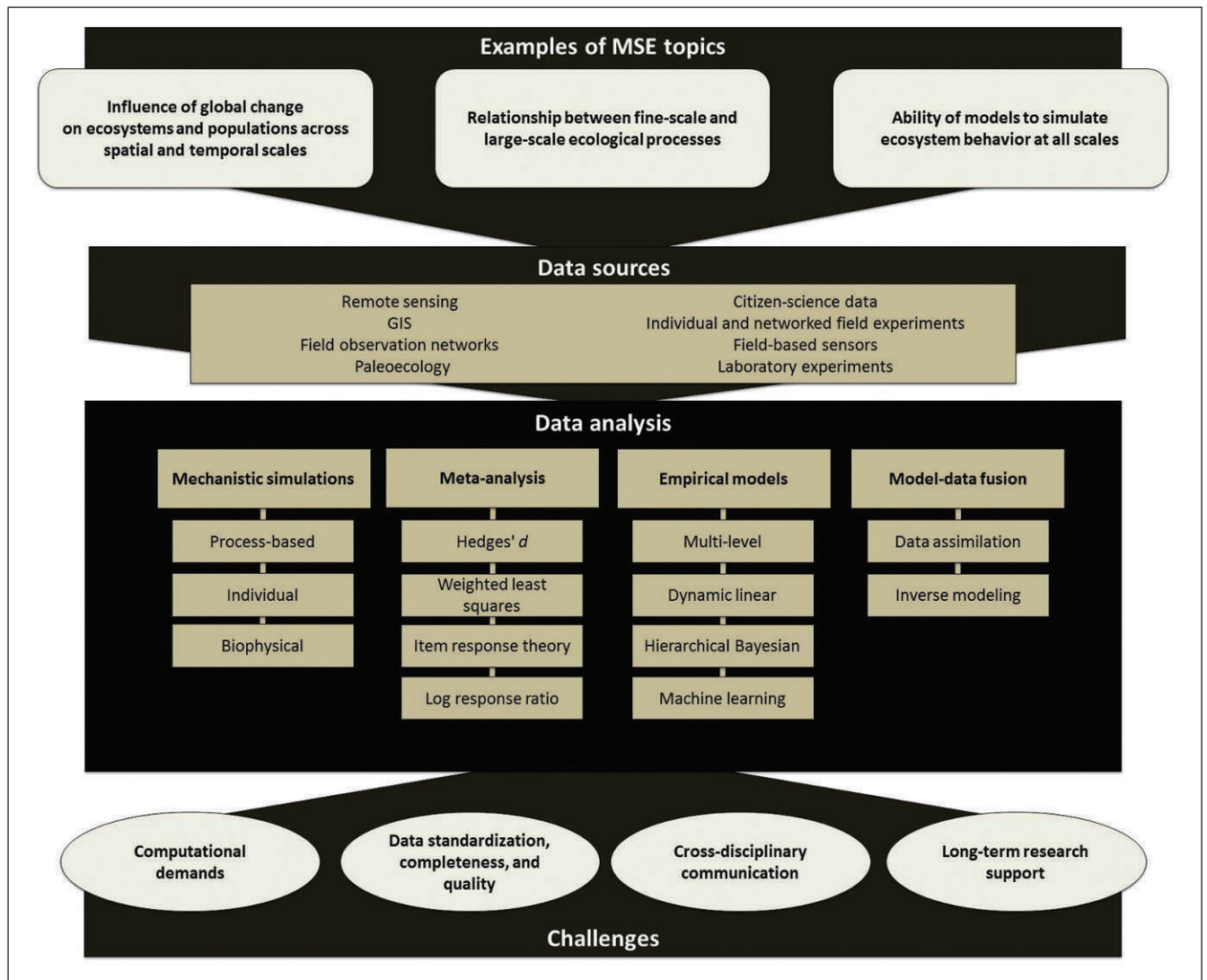


Figure 1. Relationship between the major research questions, data sources and analyses, and existing challenges in MSE research. (Top panel) Three broad areas of research are the focus of MSE. (Middle panel) These questions can be addressed using many of the existing approaches, or through a combination thereof. Although these are the dominant data sources and analyses, the list within each category is not exhaustive. (Bottom panel) All analyses are associated with problems that will need to be resolved as the discipline develops.

topography. Because MSE research questions are defined at fine-to-broad spatial and temporal scales, the data used to examine such questions must also be at such scales (Figure 1). Fortunately, recent technological and methodological advances are making it easier to obtain and distribute data measured across a range of scales, such as remote sensors on satellites or aircraft, compilations from many individual studies and citizen-science programs (Figure 2), or other studies using labor-intensive or long-term traditional methods. Notably, such heterogeneous data streams often require sophisticated, computationally demanding standardization techniques before analysis (Michener and Jones 2012; Rüegg *et al.* 2014). New approaches are emerging that can handle large volumes and diverse types of data, including mechanistic simulations, meta-analyses, empirical models, and model–data fusion (Figure 1). For example, the Paleo-Ecological Observatory Network (PaleON) is using the fusion of

model and data to integrate long-term data with terrestrial ecosystem models to better understand and model forest dynamics (Panel 1 Example 3). Using the approaches described below, MSE practitioners have the potential to make novel contributions to the understanding of broad-scale phenomena, how broad- and local-scale phenomena interact, and how such patterns and processes are likely to respond to environmental changes at multiple scales.

■ Common methodological characteristics and challenges of MSE studies

Here we highlight some important characteristics of, and strategies for meeting the challenges inherent to MSE. These elements are commonly, but not exclusively, a part of MSE studies, and the research question being asked will determine the appropriate methodology to be used and the associated difficulties (Figure 1).

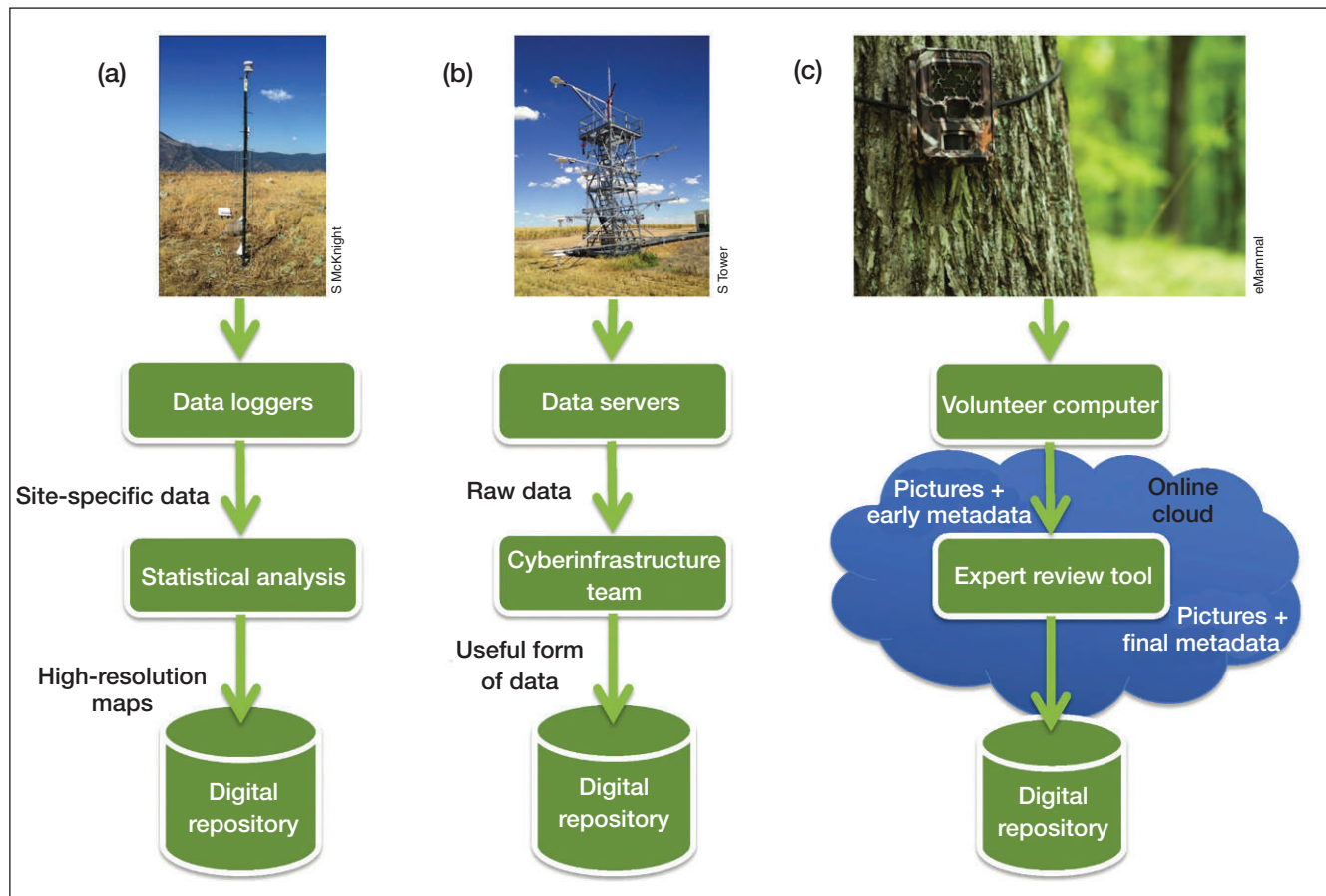


Figure 2. Networked sensors and technologies enable researchers to collect real-time, standardized, local- to continental-scale data. (a) A microclimate sensors network is automatically collecting local climate data at 10-min resolution to estimate ground-surface temperatures across entire landscapes and to test the efficacy of different methods for climate interpolation and downscaling. (b) Sensor towers of the US National Ecological Observatory Network (NEON) are collecting ecologically relevant data at the continental scale. These data are transferred automatically to NEON servers where they are standardized and transformed (including calibration, quality assurance/quality control, and data flagging) into a usable format before being shared through an online repository. (c) The eMammal network of camera traps helps researchers monitor animal populations and document the effect of recreational use on conservation areas. The motion-sensitive cameras are operated by citizen scientists and record voucher photographs of all animals as they pass by. All photos are uploaded to an online cloud and must pass expert review for data quality assurance before being archived and used in scientific research.

First, we describe the characteristics of data collection in MSE, which should allow analysis across scales and could include (1) multi-scaled and (2) diverse data objects. The large and diverse amounts of information requires researchers to adapt (3) data-intensive approaches. Finally, macrosystems analysis must incorporate the complex organization and interactions (4) across scales and the (5) hierarchy among scales. We distinguish the challenges posed by each of these characteristics and the available approaches to address them, while acknowledging that novel techniques will emerge as MSE continues to develop.

Multiple scales

Macrosystems research can include ecological processes that occur not only at local and short-term scales, but also at the spatial and temporal macroscale (hundreds to a few thousand square kilometers and temporally from days to decades and beyond; Heffernan *et al.* 2014). For instance,

studies are being conducted to identify which macroscale characteristics (eg temperature or rainfall) influence local site responses to global change (Panel 1 Example 1) and how the impact of local disturbances, such as fire, extend beyond the immediate vicinity in time and space (Goulden *et al.* 2006; Miao *et al.* 2009). The multi-scale nature of studies often requires high-resolution data: regional data are needed to examine broad-scale phenomena, but ecological processes often occur at fine spatial and temporal scales. For example, hourly environmental data may be necessary to test how environmental extremes affect regional patterns (Kearney *et al.* 2012). Moreover, when ecological information, such as land cover, varies at a spatial frequency that is finer than the data grain, aggregation effects may lead to analytic biases based on the more common finer-scaled landscape features (eg Nol *et al.* 2008; reviewed in Verberg *et al.* [2011]). On the basis of available data, researchers may choose upscaling and/or downscaling approaches to transfer data

Panel 1. Introduction of five case studies that demonstrate novel approaches to MSE

Refer to WebPanels 1–5 for the full description of each example.

Example 1

Most studies investigating ecosystem responses to climate change are conducted in a single ecosystem type; consequently, scientists lack knowledge of how (or if) site-level mechanisms – ones that explain ecological responses to climate change – may scale regionally where environmental context also varies. A geographically distributed drought experiment is being conducted in grasslands in New Mexico, Colorado, Wyoming, and Kansas that differ strongly in their ecological attributes, to test predictions of how environmental context and site-level mechanisms interact to determine regional responses.

Example 2

To provide a regional climate forecast, Salazar *et al.* (2011) proposed a hierarchical Bayesian model that assimilates different climate model simulations while accounting for discrepancies between the simulations and historical weather data. Their model acknowledges multiple sources of data and uncertainty, captures complex space–time dependence structures to improve prediction, reduces dimensionality and computational burden, and delivers full uncertainty assessment at all space and time coordinates. The results can be used to explore hypotheses related to climate change.

Example 3

Many ecological processes operate at spatiotemporal scales not amenable to direct observation and experimentation (eg the effect of decadal- to centennial-scale climate variability on tree population dynamics, legacies of historical land use, cultural eutrophication of lakes, lake acidification). Broad-scale macrosystems research thus requires the tight integration of contemporary ecological observations with geohistorical data streams and close collaborations among paleoecologists, modelers, and statisticians. The PaEON team is integrating long-term data with terrestrial ecosystem models to better understand and model forest dynamics at annual to millennial timescales.

Example 4

A data-driven approach has been used to upscale carbon (C) fluxes from the AmeriFlux network to the continental scale and to produce gridded C fluxes with high spatial (1-km) and temporal (8-day) resolutions for temperate North America. The resulting continuous gridded flux dataset – EC-MOD – has been used to assess the magnitude, distribution, and interannual variability of ecosystem C fluxes at regional and continental scales (Xiao *et al.* 2008, 2012).

Example 5

Integrating spatial and temporal data to quantify drivers of temporal patterns is a key issue for some MSE research. Dynamic linear models (DLMs; Pole *et al.* 1994) provide a framework for understanding how ecological patterns and relationships change over time (Hampton 2005) and are often more representative of the underlying data structure than traditional approaches (Lamon *et al.* 1998). DLMs have the potential to be particularly effective in MSE because they incorporate uncertainty estimates and are sensitive to changes in relationships through time.

optimally between scales, such as upscaling locally measured carbon dioxide (CO₂) fluxes from the AmeriFlux network to the continental scale (Panel 1 Example 4) and downscaling species distribution data to the grain of biological processes, to work at a scale at which management decisions can be made (Keil *et al.* 2013).

Careful selection of the most appropriate scale(s) to study and the measurements to be made at each scale is challenging because of incomplete previous knowledge, complexities that scientists are not yet able to predict (eg treatment effects that may extend beyond the treatment site), and logistical and financial constraints. Moreover, statistical detection of processes may become difficult as more locations in space and time are studied and the natural variation encountered by the study is increased. To address these issues, we argue that identifying the scales of processes by which organisms interact with the environment, resources, and other organisms is necessary. For example, when studying migratory birds, telemetry data may give insight on extent of the ecological system (eg migration limits) while also pointing to processes and patterns at local scales (eg stopover locations; Taylor *et al.* 2011). Moreover, the use of previous knowledge, such as

information available from historical records and national resource inventories, may help guide the selection of sites and time ranges for study (eg Goulden *et al.* 2006; reviewed by Hewitt *et al.* [2007]). Additionally, variables measured should represent the most likely explanatory factors needed to test the study's chosen hypothesis, which may be difficult to identify prior to the study being conducted (Hewitt *et al.* 2007). To allow statistical inference for explanatory variables, scientists must carefully select study site locations along gradients and scales (ie extent, lag, grain, and resolution; Figure 3).

Increased availability and use of automated sensors, instruments, and remote-sensing platforms enable researchers to gather data on multiple scales; these data, together with novel data-analysis approaches, can help identify underlying ecological patterns. For example, high temporal, moderate spatial resolution measurements from the Moderate Resolution Imaging Spectroradiometer (MODIS), a satellite-based instrument, can reveal regional temporal patterns (eg forest loss) that can be further investigated spatially using the high spatial resolution, low temporal resolution Landsat measurements (Potapov *et al.* 2008). An additional challenge of working

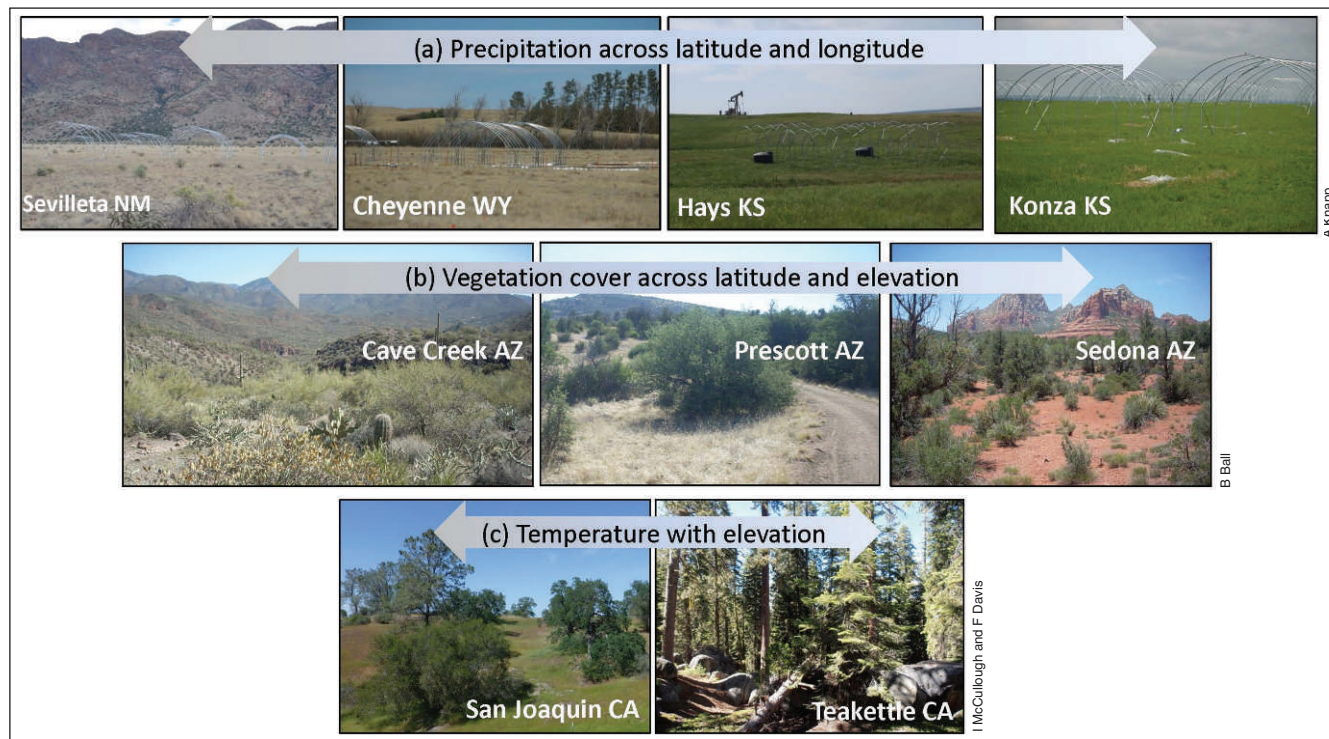


Figure 3. Examples of gradients across (a) broad, (b) intermediate, and (c) fine scales that can be used in macrosystems studies. (a) Grasslands spanning a gradient of temperature and precipitation give insight into how environmental context and ecosystem attributes may interact to determine regional patterns of response to climate change (Panel 1 Example 1). (b) The change in vegetation that occurs from low to high desert in Arizona reveals how regional environmental variability influences biological communities. (c) Elevation- and slope-related environmental gradients in mountainous areas may indicate how local climate affects the vulnerability of tree species – ones that currently dominate warm, dry foothill woodlands versus those in cool, moist montane forests – to regional climate change.

across multiple scales is that instruments, remote-sensing platforms, and climate datasets differ in resolution. Current approaches to combine spatial data with different resolutions include georectification, resampling, data fusion, and Bayesian models (Panel 1 Example 2).

Diverse data objects

Understanding complex macroscale phenomena from the systems perspective requires ecologists to look for ways to expand their data resources at both local and regional scales. With increased availability of data, researchers are able to use existing data resources. However, these datasets are usually from different thematic areas, such as population studies, geology, meteorology, and hydrology. Moreover, relevant macrosystems data may include (but are not limited to) remotely sensed imagery, citizen-science data, on-the-ground sampling data, laboratory-derived data, and reconstructed historical records, all differing in collection protocols, temporal and spatial resolution, format, quantity, quality, and costs of capture, curation, and analysis. For instance, to study how CO₂ exchange and evapotranspiration change during secondary succession, Goulden *et al.* (2006) compared high-frequency eddy covariance measurements, low-frequency tree inventories, and tree-ring analyses extending over decades.

Integrating such data objects into one unified dataset is a frequent challenge in MSE, given that traditional ecological datasets are characterized by single-investigator studies in which future applications were not considered during data collection. Recently, scientists, professional societies, and research sponsors are recognizing the value of data as a product of the scientific enterprise and placing increased emphasis on data stewardship, data sharing, openness, and supporting study repeatability (reviewed by Michener and Jones [2012]). When sharing data, ecologists need to provide complete metadata that includes such information as a full description of the methods (reviewed by Rüeegg *et al.* [2014]). In addition, if data collectors use standardization protocols as well as quality assurance (QA) and quality control (QC) procedures, their data variables can be easily converted to common variables in another database (Figure 1; Panel 1 Examples 3 and 4; reviewed by Rüeegg *et al.* [2014]). Networked sensors and technologies, such as the tower sensors of the US National Ecological Observatory Network (NEON, www.neoninc.org), the PhenoCam network (<http://phenocam.sr.unh.edu>), and the camera traps of the citizen-science eBird (www.ebird.org) and eMammal (www.facebook.com/eMammal) projects deliver regional- to continental-scale arrays of real-time data (Figure 2). Because data segments arrive from the same equipment, standardization, QA, and database compilation are relatively straightforward.

High-dimensional data (ie data containing a large number of variables relative to the number of observations) are also common in MSE, requiring specific approaches for data exploration (eg visualization), statistical inference (eg model selection and parameter estimation; Johnstone and Titterton 2009), and intensive computational power. The use of dynamic graphics enables the display of many two-dimensional projections of data (Johnstone and Titterton 2009). Promising approaches for statistical inference include machine-learning algorithms (eg Random Forests; Breiman 2001) and parameter-space sampling optimizations (eg Markov chain Monte Carlo and piece-wise Laplacian representations) that can use Distributed Computing frameworks to process large datasets (eg Hadoop; Shvachko *et al.* 2010).

Data-intensive approaches

Ecological studies usually gather discrete pieces of information over only a few years. Conversely, today's technologies are producing exponentially increasing volumes of broad-scale scientific data with networks that allow fast sharing, accessing, and collecting. In MSE, such data objects are often used as input to statistical and simulation models that themselves generate large amounts of data (ie model output). Such a volume of data poses new challenges for ecologists at various stages of study, from collecting to validating the data, to building statistical and simulation models (Kelling *et al.* 2009; reviewed by Michener and Jones [2012]), and finally to documenting and sharing data. Many of these datasets contain corrupted, missing, or meaningless sections, making it hard to obtain the relevant information about the measured variables. Efficient information management practices are therefore required to facilitate data consistency and completeness. Ecoinformatics and information management practices (and programs such as DataONE, www.dataone.org) continue to be developed to help ecologists efficiently process, store, share, integrate, and synthesize their data, while reducing data gaps and noise (Rüegg *et al.* 2014). For example, the launch of the MODIS sensor (Justice *et al.* 1998), with its near-daily global coverage and wide spectral range, catapulted the use of large remote-sensing datasets into ecosystem process models and in upscaling approaches (eg Xiao *et al.* 2012). The archiving of datasets with temporal and spatial consistency has enabled ecologists to take advantage of MODIS (Justice *et al.* 1998), without having to deal with the burdens of volume, noise, and gaps in the raw data.

In most ecological studies, data are collected and tested against specific hypotheses. However, broad-scale, multi-dimensional datasets may contain unknown (sometimes even unexpected) complexities and relationships. When seeking to understand whole-system processes, the challenges in analyzing large datasets are pushing ecology (as well as other scientific fields), toward “data-driven” approaches (Xiao *et al.* 2008;

Kelling *et al.* 2009) as opposed to the more traditional, hypothesis-testing techniques. In data-driven models, most knowledge is extracted from the data while minimizing the cost and time of model formation as well as maximizing the accuracy, speed, reliability, and comprehensibility of the models produced (Vargas *et al.* 2011). Machine-learning algorithms are able to manage multi-dimensional data with missing observations and to identify complex interactions among variables. The machine-learning approach has shown great promise in species distribution modeling. However, when data are imbalanced, these models are often biased toward selections of variables with more observations, and it is necessary to use methods – such as Cost Sensitive Learning (Zhou and Liu 2010) and Active Machine Learning (Settles 2012) – to artificially balance the data. These algorithms can be highly computer-intensive when dealing with an extensive amount of input data and may require parallel-processing to decrease execution time (Xiao *et al.* 2008). Importantly, once ecological knowledge is found, new hypotheses can be generated and tested using hypothesis-driven data collection and confirmatory analysis (Kell and Oliver 2004; Kelling *et al.* 2009).

Cross-scale interactions

In ecological systems, processes that occur at one scale may affect processes at others. For example, broad-scale precipitation regime and fine-scale soil properties jointly determine plant-available water both spatially and temporally (Browning *et al.* 2012). Similarly, warm weather may be the proximate cause of a wildfire event but factors such as tree properties and the composition and spacing within the forest determine longer-term fire dynamics (Peters *et al.* 2007). By studying multiple scales, MSE research helps reveal which interactions among scales are important features of ecological systems (Peters *et al.* 2007; Soranno *et al.* 2014). These “cross-scale interactions” can result in nonlinear dynamics and produce thresholds with pronounced implications for macrosystems behavior (Peters *et al.* 2007). However, to date, only a few examples of these interactions have been quantified. To explore these interactions, ecologists are carefully planning field studies (see study design scheme in Peters *et al.* [2008]) and developing and exploiting both statistical and process-based models.

Statistical models that rely on a multi-scaled dataset can be used to determine the operating scales (eg units of time or space) for the macrosystem of interest as well as the interactions that occur across those scales (see Rüegg *et al.* [2014] for an example of database compilation and integration). For example, hierarchical models allow the incorporation of variables at multiple spatial and temporal extents (Qian *et al.* 2010); in particular, Bayesian hierarchical models that use quantitative inference to accommodate unbalanced data across space and/or through time (Cressie *et al.* 2009) have recently been

applied to quantify cross-scale interactions and describe their nonlinear dynamics (Panel 1 Example 5; Soranno *et al.* 2014).

A major and critical challenge in ecology is to understand the processes behind these interactions, especially for forecasting future dynamics. Biophysical niche models – which combine the morphology, physiology, and behavior of an organism – are being used to predict species distributions that are solely based on climate conditions (Figure 4). Such models, for example, have shown how macroscale climate limits species distribution (eg Buckley *et al.* 2010) and activity times (eg Sears *et al.* 2011), or how diel cycles in ambient temperature may have shaped activity patterns in small mammals (Levy *et al.* 2012). These models may allow for the explicit assessment of how plasticity or evolutionary changes at the individual level affect ecological communities at coarser scales and may be a way to determine when and how cross-scale interactions are shaping species ranges and behaviors. Modeling across time, space, and levels of biological organization is an exciting new direction for MSE research, one that is needed in order to meet the pressing needs of global change.

Hierarchy among scales

Macrosystems can be viewed as one “level” in a hierarchical system that includes levels from local to global spatial extents (Heffernan *et al.* 2014). There is widespread consensus that ecological complexity (ie biocomplexity) emerges from the interactions between organisms and their biotic and abiotic environments (Anand *et al.* 2010). In a bottom-up process, for example, spatiotemporal patterns of population and community dynamics are often emergent properties that can only be captured by studying much finer levels of ecological detail. On the other hand, in a top-down process, high fitness costs will be caused by range retractions that decrease the genetic pool and lead to increased inbreeding. Studying these kinds of hierarchical interactions is not straightforward; the multi- and cross-scaled nature of the data is further complicated by the possible interactions among levels of ecological organization, posing serious statistical challenges (eg Finley *et al.* 2009). Moreover, practical constraints of time and space may limit the ability to observe and manipulate interactions and emergent processes that occur between ecological hierarchical levels. Currently, modeling tools, such as hierarchical Bayesian methods for statistical analysis, and individual-based models (IBMs) serving as “virtual laboratories” may help solve these problems.

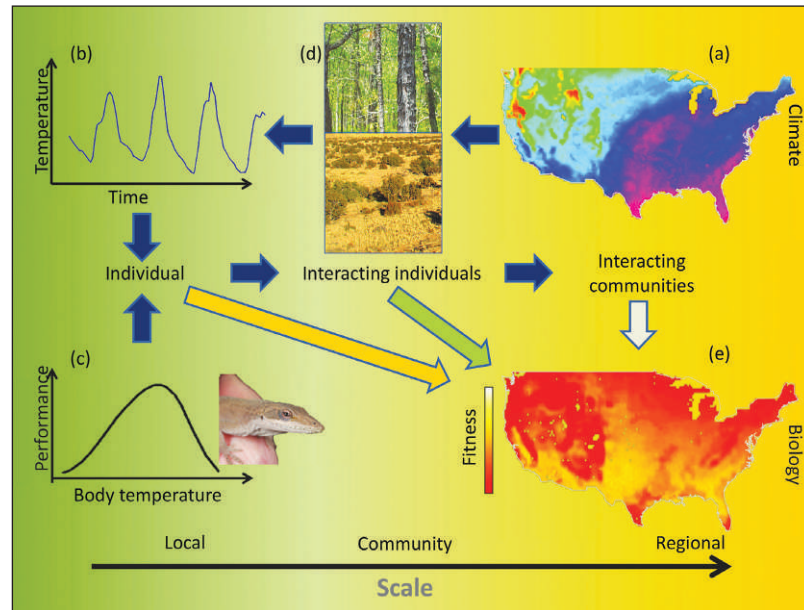


Figure 4. Biophysical niche models are used to study how climate affects animal survival, growth, and reproduction at macroscales. Because climate and individuals operate at different scales, information flow between the scales is necessary. Regional climate datasets (a) are downscaled to local climate datasets (b), which are combined with knowledge of the morphology, physiology, and behavior of an organism to predict organismal fitness in one location (ie parameters such as activity times, growth rate, survival, and reproduction). Through the use of individual-based models, it is possible to study climate effects at different levels of ecological organizations; simulating ecological communities (d) with both climate and interactions among individuals (eg competition and predation) and allowing movements of organisms between adjacent communities can help ecologists study how community-level interactions and broad-scale processes (eg gene flow) may affect individuals’ fitness. At each hierarchical level, fitness maps can be drawn from model results for each location (e). Orange arrow = individual-level model; green arrow = community-level model; white arrow = metacommunity-level model.

The use of hierarchical Bayesian methods is particularly well suited to deal with complex dependence structures in statistical modeling and thus represents a valuable analytical framework for making inferences at macroscales (Panel 1 Example 2). Finley *et al.* (2009) used plot-based estimates of the National Forest Inventories, such as tree species composition in the US, together with environmental predictors such as climate variables, to model regional forest tree species composition and to gain insight into forest ecosystem sustainability, biodiversity, and productivity. Using spatial multinomial hierarchical Bayesian models, the authors were able to improve prediction of species composition by taking into account the spatial proximity between measurements and showed that space-varying relationships exist between species occupancy and environmental predictors. This approach presents many difficulties, including specifying valid probability models, implementation, and high computational demands (eg Banerjee *et al.* [2008] and references therein).

IBMs are also well suited to study emergent properties

between different organizational levels (Figure 4). These models can be directly and relatively simply parameterized and have the intrinsic ability to include both temporal and spatial scales, allowing researchers to observe the outcome on a population of individuals (Anand *et al.* 2010). Regional IBMs can be used to study how different levels of organization, from genes to individuals to populations, can survive, grow, evolve, and interact to shape species distributions. In such instances, improving landscape realism using geographic information systems and remote-sensing data will enhance our understanding of the processes shaping communities (eg Wallentin *et al.* 2008). Moreover, multiple species simulations can provide insights into the functional roles of organisms in an ecological system and how interspecific interactions that occur locally may have a broader impact on ecological communities. Alternatively, comparisons between complex and simplified models (eg by excluding organizational levels, reducing spatial resolution, relaxing environmental stochasticity) may help identify the most important levels and interactions of an ecological system. Although individuals operate at scales of hours and meters, data at these scales are not yet readily available, making simplifications inevitable in many cases.

■ Conclusions and future directions

To investigate how long-term and broad-scale phenomena influence or interact with ecological patterns and processes at other scales, ecologists need to collect sufficient data and use robust techniques of data standardization and analysis (Figure 1). Many ongoing broad-scale data collection and integration efforts will provide valuable, standardized data to support such studies. However, there are many challenges – from study design, to data collection, to analysis – that need to be considered. During the study design stage, there is often incomplete information regarding which factors operating from global to local scales need to be measured to understand the process of interest. At the data collection stage, it is necessary to discover the most relevant data resources that come with various resolutions and collection techniques; these data must then be combined, validated, and standardized. Innovative statistical and simulation techniques provide flexible approaches for explaining cross-scale and hierarchical dynamics and interactions in the ecological system.

MSE is in an early stage of development. Many innovative techniques, such as Bayesian hierarchical models, machine learning, mechanistic simulations, meta-analysis, and model–data fusion, are currently used. Still, there is much room for development of novel approaches for data collection and analysis. For example, to observe natural multi-scale processes and interactions, ecologists need to evolve field techniques for multi-scale observational and experimental studies (eg automated comprehensive field data collection across networks, experiments like those in

Panel 1 Example 1). In most cases, these approaches require cross-disciplinary communication among scientists from many fields, including statistics, geophysics, climatology, and computer and information science (Goring *et al.* 2014; Rüegg *et al.* 2014).

Macrosystems research is a resource-intensive undertaking that requires sufficient time and funding, typically scaling beyond traditional single-investigator experimental work. These requirements can be a substantial limitation to realizing the potential of larger-scale and more integrative studies. Support from funding agencies and research institutions for data documentation and long-term access will be a key to the success of MSE. Scientists, scientific organizations, and institutions should promote a culture of data sharing – for example, by giving credit for publishing data (and metadata) and contributing to data libraries – and scientists should get into the habit of providing open access to both raw and processed data (Goring *et al.* 2014).

In summary, practitioners of MSE studies must use a suite of approaches and methods to answer questions across increased scales and levels of complexity, while dealing with the difficulties inherent in MSE. Continued innovation in methodologies will allow for the development and testing of exciting new hypotheses and theories across broad spatial and temporal extents.

■ Acknowledgements

This paper is the result of the efforts of Working Group 2 at the MacroSystems Biology PI meeting (March 2012) in Boulder, CO. We thank the MacroSystems Biology program in the Emerging Frontiers Division of the Biological Sciences Directorate at NSF for support, as well as all of the meeting and working group participants for their valuable input. Use of trade names does not imply endorsement by the US Government. For author contributions, see WebPanel 6.

■ References

- Anand M, Gonzalez A, Guichard F, *et al.* 2010. Ecological systems as complex systems: challenges for an emerging science. *Diversity* 2: 395–410.
- Banerjee S, Gelfand AE, Finley AO, and Sang H. 2008. Gaussian predictive process models for large spatial data sets. *J Roy Stat Soc B* 70: 825–48.
- Breiman L. 2001. Random forests. *Mach Learn* 45: 5–32.
- Browning DM, Duniway MC, Laliberte AS, and Rango A. 2012. Hierarchical analysis of vegetation dynamics over 71 years: soil-rainfall interactions in a Chihuahuan Desert ecosystem. *Ecol Appl* 22: 909–26.
- Buckley LB, Urban MC, Angilletta MJ, *et al.* 2010. Can mechanism inform species' distribution models? *Ecol Lett* 13: 1041–54.
- Cressie N, Calder CA, Clark JS, *et al.* 2009. Accounting for uncertainty in ecological analysis: the strengths and limitations of hierarchical statistical modeling. *Ecol Appl* 19: 553–70.
- Evans MR, Norris KJ, and Benton TG. 2012. Predictive ecology: systems approaches introduction. *Philos T Roy Soc B* 367: 163–69.
- Finley AO, Banerjee S, and McRoberts RE. 2009. Hierarchical spa-

- tial models for predicting tree species assemblages across large domains. *Ann Appl Stat* 3: 1052–79.
- Golley FB (Ed). 1993. A history of the ecosystem concept in ecology. New Haven, CT: Yale University Press.
- Goring SJ, Weathers KC, Dodds WK, et al. 2014. Improving the culture of interdisciplinary collaboration in ecology by expanding measures of success. *Front Ecol Environ* 12: 39–47.
- Goulden ML, Winston GC, McMillan AMS, et al. 2006. An eddy covariance mesonet to measure the effect of forest age on land–atmosphere exchange. *Global Change Biol* 12: 2146–62.
- Hagen JB (Ed). 1992. An entangled bank: the origins of ecosystem ecology. New Brunswick, NJ: Rutgers University Press.
- Hampton SE. 2005. Increased niche differentiation between two *Conochilus* species over 33 years of climate change and food web alteration. *Limnol Oceanogr* 50: 421–26.
- Heffernan JB, Soranno PA, Angilletta MJ, et al. 2014. Macrosystems ecology: understanding ecological patterns and processes at continental scales. *Front Ecol Environ* 12: 5–14.
- Hewitt JE, Thrush SF, Dayton PK, and Bonsdorff E. 2007. The effect of spatial and temporal heterogeneity on the design and analysis of empirical studies of scale-dependent systems. *Am Nat* 169: 398–408.
- Johnstone IM and Titterton DM. 2009. Statistical challenges of high-dimensional data. *Philos T Roy Soc A* 367: 4237–53.
- Justice CO, Vermote E, Townshend JRG, et al. 1998. The Moderate Resolution Imaging Spectroradiometer (MODIS): land remote sensing for global change research. *IEEE T Geosci Remote* 36: 1228–49.
- Kearney MR, Matzelle A, and Helmuth B. 2012. Biomechanics meets the ecological niche: the importance of temporal data resolution. *J Exp Biol* 215: 922–33.
- Keil P, Belmaker J, Wilson AM, et al. 2013. Downscaling of species distribution models: a hierarchical approach. *Method Ecol Evol* 4: 82–94.
- Kell DB and Oliver SG. 2004. Here is the evidence, now what is the hypothesis? The complementary roles of inductive and hypothesis-driven science in the post-genomic era. *Bioessays* 26: 99–105.
- Kelling S, Hochachka WM, Fink D, et al. 2009. Data-intensive science: a new paradigm for biodiversity studies. *BioScience* 59: 613–20.
- Lamon EC, Carpenter SR, and Stow CA. 1998. Forecasting PCB concentrations in Lake Michigan salmonids: a dynamic linear model approach. *Ecol Appl* 8: 659–68.
- Levy O, Dayan T, Kronfeld-Schor N, and Porter WP. 2012. Biophysical modeling of the temporal niche: from first principles to the evolution of activity patterns. *Am Nat* 179: 794–804.
- Miao S, Carstenn S, Thomas C, et al. 2009. Integrating multiple spatial controls and temporal sampling schemes to explore short- and long-term ecosystem response to fire in an Everglades wetland. In: Miao S, Carstenn S, and Nungesser M (Eds). *Real world ecology*. New York, NY: Springer.
- Michener WK and Jones MB. 2012. Ecoinformatics: supporting ecology as a data-intensive science. *Trends Ecol Evol* 27: 85–93.
- Nol L, Verburg PH, Heuvelink GBM, and Molenaar K. 2008. Effect of land cover data on nitrous oxide inventory in fen meadows. *J Environ Qual* 37: 1209–19.
- Peters DPC, Bestelmeyer BT, and Turner MG. 2007. Cross-scale interactions and changing pattern–process relationships: consequences for system dynamics. *Ecosystems* 10: 790–96.
- Peters DPC, Groffman PM, Nadelhoffer KJ, et al. 2008. Living in an increasingly connected world: a framework for continental-scale environmental science. *Front Ecol Environ* 6: 229–37.
- Pole A, West M, and Harrison J (Eds). 1994. *Applied Bayesian forecasting and time series analysis*. New York, NY: Chapman & Hall.
- Potapov P, Hansen MC, Stehman SV, et al. 2008. Combining MODIS and Landsat imagery to estimate and map boreal forest cover loss. *Remote Sens Environ* 112: 3708–19.
- Qian SS, Cuffney TF, Alameddine I, et al. 2010. On the application of multilevel modeling in environmental and ecological studies. *Ecology* 91: 355–61.
- Rüegg J, Gries C, Bond-Lamberty B, et al. 2014. Completing the data life cycle: using information management in macrosystems ecology research. *Front Ecol Environ* 12: 24–30.
- Salazar E, Sanso B, Finley AO, et al. 2011. Comparing and blending regional climate model predictions for the American Southwest. *J Agr Biol Envir S* 16: 586–605.
- Sears MW, Raskin E, and Angilletta MJ. 2011. The world is not flat: defining relevant thermal landscapes in the context of climate change. *Integr Comp Biol* 51: 666–75.
- Settles B (Ed). 2012. *Active learning*. San Rafael, CA: Morgan & Claypool.
- Shvachko K, Hairong K, Radia S, and Chansler R. 2010. The Hadoop distributed file system. In: *IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST)*; 3–7 May 2010; Incline Village, NV. doi:10.1109/MSST.2010.5496972.
- Soranno PA, Cheruvilil KS, Bissell EG, et al. 2014. Cross-scale interactions: quantifying multi-scaled cause–effect relationships in macrosystems. *Front Ecol Environ* 12: 65–73.
- Taylor PD, Mackenzie SA, Thurber BG, et al. 2011. Landscape movements of migratory birds and bats reveal an expanded scale of stopover. *PLoS ONE* 6: e27054.
- Vargas R, Carbone M, Reichstein M, and Baldocchi D. 2011. Frontiers and challenges in soil respiration research: from measurements to model–data integration. *Biogeochemistry* 102: 1–13.
- Verburg PH, Neumann K, and Nol L. 2011. Challenges in using land use and land cover data for global change studies. *Glob Change Biol* 17: 974–89.
- Wallentin G, Tappeiner U, Strobl J, and Tasser E. 2008. Understanding alpine tree line dynamics: an individual-based model. *Ecol Model* 218: 235–46.
- Xiao JF, Zhuang QL, Baldocchi DD, et al. 2008. Estimation of net ecosystem carbon exchange for the conterminous United States by combining MODIS and AmeriFlux data. *Agr Forest Meteorol* 148: 1827–47.
- Xiao JF, Chen JQ, Davis KJ, and Reichstein M. 2012. Advances in upscaling of eddy covariance measurements of carbon and water fluxes. *J Geophys Res-Biogeophys* 117: G00J01.
- Zhou ZH and Liu XY. 2010. On multi-class cost-sensitive learning. *Comput Intell* 26: 232–57.

⁵University of Wisconsin Center for Limnology, Boulder Junction, WI; ⁶University of Illinois, Urbana, IL; ⁷University of New Hampshire, Durham, NH; ⁸University of Oklahoma, Norman, OK; ⁹University of North Carolina, Chapel Hill, NC; ¹⁰Iowa State University, Ames, IA; ¹¹University of Texas at Austin, Austin, TX; ¹²Brown University, Providence, RI; ¹³Colorado State University, Fort Collins, CO; ¹⁴Harvard University, Cambridge, MA; ¹⁵University of Delaware, Newark, DE; ¹⁶US Geological Survey, Pennsylvania Cooperative Fish & Wildlife Research Unit, Pennsylvania State University, University Park, PA; ¹⁷University of Wisconsin-Madison, Madison, WI; †these authors contributed equally to this work