# Predicting distributions of Mexican birds using ecological niche modelling methods

A. TOWNSEND PETERSON*, LISA G. BALL & KEVIN P. COHOON
*Natural History Museum, University of Kansas, Lawrence, Kansas 66045, USA*

We tested the utility of the modelling program *Genetic Algorithm for Rule-set Prediction* (GARP) for modelling ecological niches to make accurate predictions of geographical distributions for 25 bird species across Mexico. Specimen-based point-occurrence data were entered into the algorithm in the form of geographical coordinates, and related to digitized maps of environmental variables, including mean annual precipitation, elevation, mean annual temperature, and potential vegetation. Two Mexican states were used as test areas by withholding their points from model construction; these points were later overlaid on predictions to measure model performance. Statistically, most models (78–90%) were significantly more powerful than random models in predicting occurrences in test states; model failures were most often due to low sample size for testing, rather than an inability to model distributions of particular species. The success of this test indicates that ecological niche modelling approaches such as GARP provide a promising tool for exploring a broad range of questions in ecology, biogeography and conservation.

The distribution, richness and uniqueness of taxa must be assessed if priorities for conservation action are to be established. Only then can areas representing concentrations of rare, endangered or endemic species be accorded high priority in conservation decisions (Wilson 1988). Obtaining detailed and accurate distributional information for such species therefore becomes a critical step in the biodiversity conservation challenge.

Existing biological inventory data, however, are almost universally patchy and incomplete (Peterson *et al.* 1998). Because access to remote areas is costly, sample localities cluster along access routes (Hall & Moreau 1970, Binford 1989). More thorough sampling, as might conform to a uniform grid or random points, can be prohibitively expensive and time-consuming, if not simply impossible. An approach to maximizing the utility of available distributional data for sparsely sampled regions is to employ models that predict distributions based on ecological requirements of species, extrapolating from known points into unknown areas. A variety of such algorithms have been developed (i.e. Nix 1986, Miller *et al.* 1989, Walker & Cocks 1991; Carpenter *et al.* 1993).

An innovative modelling algorithm, the *Genetic Algorithm for Rule-set Production* (GARP; http://biodi.sdsc.edu/), uses some of the same theoretical approaches and algorithms upon which earlier approaches were based (e.g. logistic regression, BIOCLIM) in a dynamic machine-learning environment (Stockwell & Noble 1992). GARP is an artificial-intelligence-based super-algorithm that works by combining sets of rules to build the most accurate prediction possible for the region being considered (Stockwell & Noble 1992). GARP uses known occurrence points and electronic maps of relevant ecological dimensions to produce a model of a species' ecological niche, which when projected onto a landscape provides a prediction of the species' geographical distribution.

The present analysis centred around GARP's ability to predict species' ecological requirements (and hence geographical distributions) in areas for which no input occurrence points are available. We sampled available points a priori to provide independence of model building and testing. This constituted a particularly difficult challenge for the algorithm. We used Mexican birds as a test group, because ample data were available for the region (Peterson *et al.* 1998). Although large portions of the region are sparsely sampled, the states of Jalisco and Oaxaca are well sampled (Peterson *et al.* 1998), and are well separated

*Corresponding author.
Email: town@ukans.edu

geographically; hence, we used them as two independent test areas.

## METHODS

Species' occurrence data were drawn from 32 scientific collections (see Appendix). Specimen records were compiled as part of the *Atlas of Mexican Bird Distributions* project (Peterson *et al.* 1998).

We arbitrarily chose a minimal sample size of 50 unique latitude–longitude records and at least five unique occurrence points in each of the test states. Of the terrestrial, non-migratory species that met these criteria, 25 were selected at random (Table 1). For each species, tests of model predictivity were built for the two states separately. Localities in the test state (Oaxaca or Jalisco) were excluded from predictive analyses; these subsets ('test data') were withheld for measuring accuracy of the species predictions. Data from remaining states ('training data') were used to develop predictive models.

Ecological niches (the conjunction of ecological conditions within which a species is able to maintain populations) were modelled using GARP, which includes several distinct algorithms for niche modelling in an artificial-intelligence-based approach (Stockwell & Noble 1992). Input points are divided evenly into datasets for rule generation and model testing. GARP works in an iterative process of rule selection, evaluation, testing and incorporation or rejection: first, a method is chosen from a set of possibilities (e.g. logistic regression, bioclimatic rules), applied to the training data, and a rule developed; rules may be spliced, truncated or otherwise adjusted. Predictive accuracy is evaluated on the basis of 1250 points resampled from the test data and 1250 points sampled randomly from the study region as a whole. Change in predictive accuracy from one iteration to the next is used to evaluate whether a particular rule should be incorporated into the model, and the algorithm runs either 1000 iterations or until convergence. It should be noted that our use of 'test' and 'training' data differs from that of Stockwell and Noble (1992) in that our test dataset represents a further division of available information into subsets.

**Table 1.** Test data sample sizes, chi-square values, and significance levels for distributional models for 25 Mexican bird species.

| Species | Jalisco | | | Oaxaca | | |
|---|---|---|---|---|---|---|
| | *n* | $\chi^2$ | *P* | *n* | $\chi^2$ | *P* |
| *Aratinga canicularis* | 8 | 14.13 | *** | 32 | 12.63 | *** |
| *Atlapetes pileatus* | 14 | 9.27 | *** | 13 | 7.83 | *** |
| *Camptostoma imberbe* | 11 | 4.93 | * | 24 | 16.47 | *** |
| *Campylorhynchus megalopterus* | 5 | 3.94 | * | 8 | 6.34 | * |
| *Certhia americana* | 13 | 6.45 | * | 10 | 5.36 | * |
| *Ciccaba virgata* | 5 | 1.78 | * | 21 | 8.74 | ** |
| *Columba fasciata* | 5 | 5.17 | n.s. | 6 | 8.13 | ** |
| *Contopus pertinax* | 24 | 4.90 | * | 16 | 8.65 | ** |
| *Lampornis amethystinus* | 13 | 7.80 | * | 16 | 13.23 | *** |
| *Lepidocolaptes leucogaster* | 14 | 4.80 | ** | 10 | 3.31 | + |
| *Melanerpes formicivorus* | 18 | 11.62 | * | 22 | 2.86 | + |
| *Mitrephanes phaeocercus* | 8 | 6.71 | *** | 16 | 4.78 | * |
| *Momotus mexicanus* | 24 | 15.97 | ** | 56 | 2.97 | + |
| *Myadestes occidentalis* | 17 | 5.34 | *** | 14 | 4.68 | * |
| *Otus trichopsis* | 7 | 6.64 | * | 9 | 12.85 | *** |
| *Parula pitiayumi* | 6 | 7.47 | ** | 7 | 2.86 | + |
| *Parus wollweberi* | 18 | 4.00 | * | 12 | 10.02 | ** |
| *Picoides villosus* | 8 | 4.09 | * | 21 | 5.65 | * |
| *Piculus auricularis* | 8 | 3.36 | + | 16 | 6.77 | ** |
| *Piranga bidentata* | 16 | 0.94 | n.s. | 11 | 0.20 | n.s. |
| *Piranga erythrocephala* | 12 | 4.61 | * | 6 | 1.20 | n.s. |
| *Tityra semifasciata* | 11 | 4.70 | * | 39 | 18.62 | *** |
| *Trogon mexicanus* | 16 | 10.82 | ** | 15 | 13.38 | *** |
| *Vireo hypochryseus* | 12 | 3.19 | + | 13 | 0.70 | n.s. |
| *Vireolanius melitophrys* | 10 | 7.27 | ** | 16 | 9.79 | ** |

*$P < 0.05$, **$P < 0.01$, ***$P < 0.001$, +$0.05 < P < 0.10$, n.s. not significant.

Geographical data for Mexico were entered into the GARP algorithm in the form of digitized raster data coverages. Geographical data layers used included annual mean temperature, annual mean precipitation, elevation and potential vegetation, at a spatial resolution of 7 × 7-km pixels. These data, based on maps produced by Instituto Nacional de Estadística, Geografía e Informática (INEGI), were kindly provided in digitized form by the Comisión Nacional para el Uso y Conocimiento de la Biodiversidad (CONABIO).

Separate GARP models were produced for the 25 species for each test state, so 50 models were built in all. Models were saved in ASCII raster grid format, and imported into a GIS: ArcView (version 3.1). The proportional area of the test state (Jalisco or Oaxaca) predicted present was calculated for each of the 50 models. Test points were then overlaid on each to determine numbers falling within areas of predicted presence. A goodness-of-fit test was used to calculate how well the test points fitted the expected model, comparing success in predicting test points with that expected from the proportion of the test state predicted to hold the species. It is important to emphasize the difficult conditions under which GARP was tested. Presences and absences of species were predicted in areas *outside* the geographical extent of the training points, and tested with independent data.

## RESULTS

Figure 1 presents an example (Rufous-capped Brush-Finch *Atlapetes pileatus*, test in Jalisco) of testing
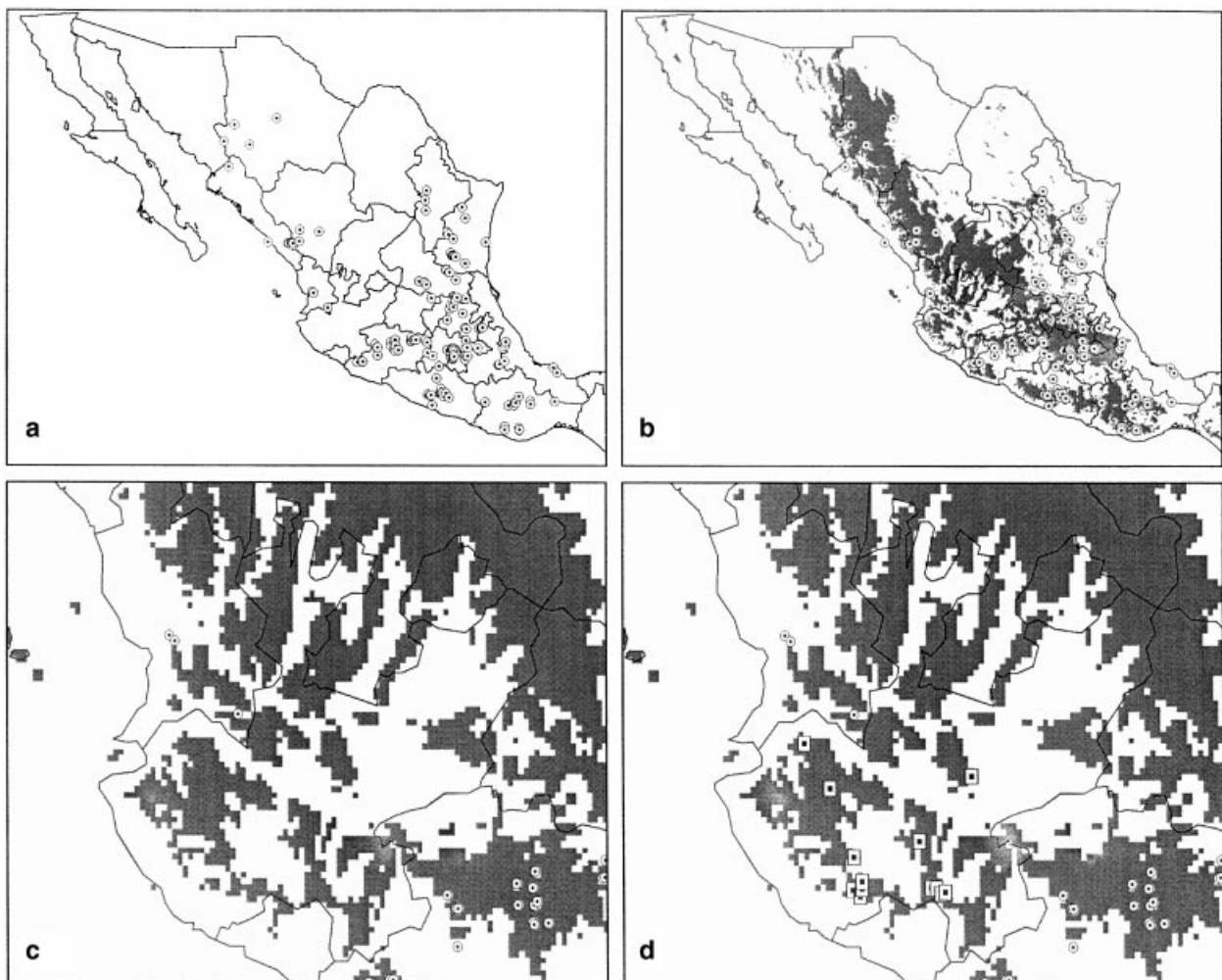


**Figure 1.** Summary of model building and testing for the Rufous-capped Brush-Finch *Atlapetes pileatus*: (a) training data (excluding points in Jalisco); (b) predictive model for the species across Mexico, with areas predicted present in black; (c) close-up of predictive model in test state; and (d) overlay of test data on prediction.
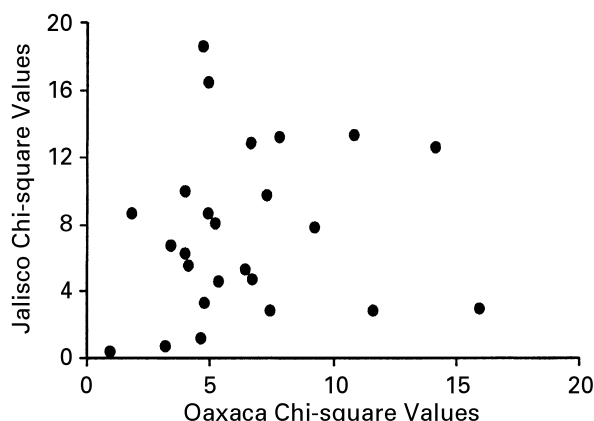
**Figure 2.** Comparison of chi-square statistics for 25 bird species between Jalisco and Oaxaca.

procedures employed in this study, showing training data (Fig. 1a), the model developed (Fig. 1b,c), and application of the test data to the prediction (Fig. 1d). The Jalisco test for Rufous-capped Brush-finch correctly predicted 14 of 14 test points available, whereas only 60.2% of the state was predicted present, yielding a chi-square statistic of 9.27 ($df$ = 1, $P < 0.001$).

For Jalisco, 21 of 25 models were significant ($P < 0.05$), and for Oaxaca, 18 of 25 models were significant ($P < 0.05$) (Table 1). Two further species in

Jalisco and four in Oaxaca approached statistical significance ($0.1 > P > 0.05$). Hence, 78–90% of distributional models were statistically significantly better than random at predicting distributions of test points.

We inspected aspects of species and sampling that might be related to predictive ability. A first suspicion was that particular species might have ecological characteristics that would be difficult to model, in which case a positive association of chi-square values for species between the two states would be expected. However, bivariate plots revealed no significant correlation ($P > 0.05$, Fig. 2). Hence, the data did not support the idea that particular species were more or less able to be modelled accurately.

Regarding effects of sample size, we tested for influences of two aspects of sample size on predictive efficiency. Training sample sizes and associated chi-square statistics were not significantly associated in either Jalisco or Oaxaca ($P > 0.05$, in both states; Fig. 3). Test sample size, however, showed a positive relationship that approached statistical significance for Jalisco ($r^2 = 0.325$, $P = 0.07$; Fig. 3d). In Oaxaca, although the overall association was not significant ($P > 0.05$), removal of one striking outlier (Russet-crowned Motmot *Momotus mexicanus*, large test dataset, low chi-square value; Fig. 3c) left a significant association ($r^2 = 0.589$, $P = 0.002$). Hence,
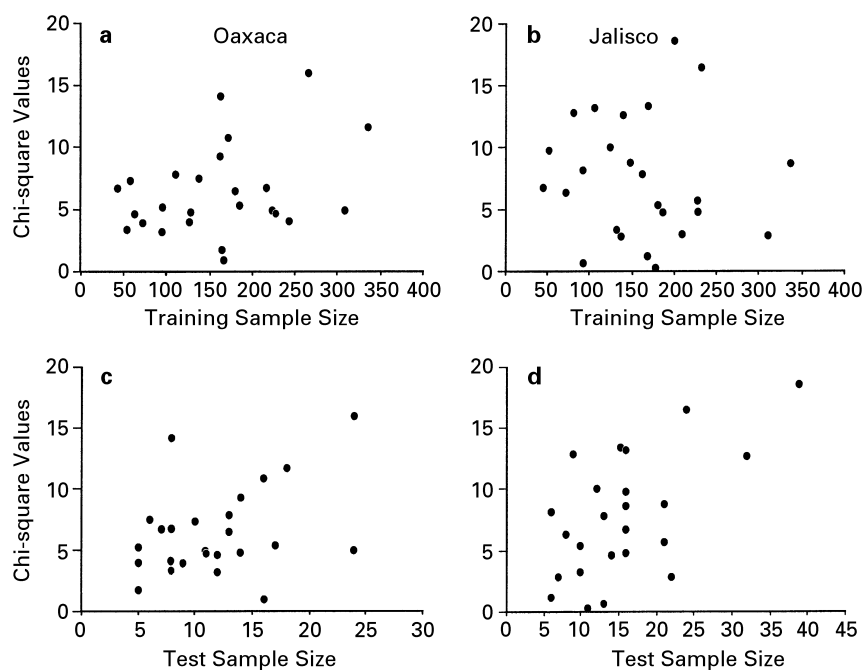


**Figure 3.** Comparison of chi-square statistics for 25 species with sample sizes in training (panels a and b) and test (panels c and d) datasets for Oaxaca and Jalisco, respectively.

evidence suggests that sample size for *testing* model accuracy may influence the predictive accuracy of distributional models used.

## DISCUSSION

Our tests of the accuracy of models and predictions of geographical distributions developed with the GARP algorithm demonstrate its ability to predict species' distributions in real-world situations. We tested the algorithm under difficult conditions, in which predictions were tested in areas outside the spatial extent of the training data. This general approach thus offers a fascinating possibility for understanding species' distributions in poorly known regions.

Although most models were statistically significant, we considered several potential sources of error that may have caused some models to fail. One possibility is that the resolution of the data used in the study may have contributed to poor performance of some models. If key environmental characteristics for a species are not 'visible' at the resolution of the environmental data used, accurate modelling of ecological niches and geographical distributions would not be possible. For example, in Mexico, the Acorn Woodpecker *Melanerpes formicivorus* occurs most abundantly in stands of dead trees within oak and pine–oak woodlands (Short 1982, Koenig *et al.* 1995, Peterson pers. obs.), which would not be detectable at the resolution of the present study; in fact, the model for this species' distribution in Oaxaca failed to achieve statistical significance. Nonetheless, the absence of an association of predictivity of species' models between the two test states suggests that this complication does not constitute an overriding problem for predicting species' geographical distributions.

The ability to predict species' distributions accurately clearly can be only as good as the geographical and occurrence data that enter the algorithm. The electronic maps used in this study are certain to contain errors, as are some of our interpretations of historical museum specimen locality data. Transmission of these error components can cause compounding of errors in synthetic analyses such as those developed here (Veregin 1989). With geographical data at finer resolutions, additional information sources such as satellite imagery, and occurrence data based on coordinates from global positioning systems, we expect the accuracy of this algorithm to improve.

Comparisons of accuracy of models between states suggested that it is not the case that certain species are difficult or impossible to model. Moreover, accuracy of our models was not sensitive to sample size in the training dataset, suggesting that the sample sizes used herein for model development (43–334) did not approach minimum thresholds for model development. Rather, our results suggested that sample sizes for testing models were the primary factors influencing our measures of predictive accuracy. Under this view, the *models* were probably good, but our statistical tests of their accuracy were not sufficiently powerful to detect their significant predictive ability. Hence, our confidence in distributional predictions developed using GARP was bolstered considerably by the results of these analyses.

In conclusion, our results suggest that ecological niche modelling approaches to predicting geographical distributions of species – GARP in particular – constitute promising tools for a wide variety of applications. As demonstrated here, they can be used to model individual species, and these models can be important in studies of autecology, in locating rare and poorly known species, and in focusing conservation efforts for endangered species (e.g. Peterson *et al.* 2000). In biogeographical applications, distributional modelling can provide valuable insights into historical origins, relations with other species and potential colonizing ability (e.g. Peterson *et al.* 1999). Potential applications to conservation biology are numerous: GARP models have already been used to develop optimal reserve designs for rare and endangered species (Egbert *et al.* 1998, Peterson *et al.* 2000) and to predict species invasions (Peterson & Vieglais in press).

## REFERENCES

**Binford, L.C.** 1989. A distributional survey of the birds of the Mexican state of Oaxaca. *Ornith. Monogr.* **43**: 1–418.

**Carpenter, G., Gillison, A.N. & Winter, J.** 1993. DOMAIN: a flexible modeling procedure for mapping potential distributions of plants and animals. *Biodiv. Conserv.* **2**: 667–680.

**Egbert, S.L., Peterson, A.T., Sánchez-Cordero, V. & Price, K.P.** 1998. Modeling conservation priorities in Veracruz, Mexico. In: Morain, S. (ed.) *GIS in Natural Resource Management: Balancing the Technical-Political Equation*: 55–63. Santa Fe, New Mexico: High Mountain press.

**Hall, B.P. & Moreau, R.E.** 1970. *An Atlas of Speciation in African Passerine Birds*. London: British Museum (Natural History).

**Koenig, W.P., Stacey, P.B., Stanback, M.T. & Mumme, R.L.** 1995. Acorn Woodpecker. *Birds North Am.* **194**: 1–24.

**Miller, R.I., Stuart, S.N. & Howell, K.N.** 1989. A methodology for analyzing rare species distribution patterns utilizing GIS technology: The rare birds of Tanzania. *J. Landscape Ecol.* **2**: 173–189.

**Nix, H.A.** 1986. A biogeographic analysis of Australian elapid snakes. *Aust. Flora Fauna Series* **8**: 4–15.

**Peterson, A.T., Egbert, S.L., Sánchez-Cordero, V. & Price, K.P.** 2000. Geographic analysis of conservation priorities using distributional modelling and complementarity: Endemic birds and mammals in Veracruz, Mexico. *Biol. Conserv.* **93**: 85–94.

**Peterson, A.T., Navarro-Sigüenza, A.G. & Benítez-Díaz, H.** 1998. The need for continued scientific collecting: a geographic analysis of Mexican bird specimens. *Ibis* **140**: 288–294.

**Peterson, A.T., Soberón, J. & Sánchez-Cordero, V.** 1999. Conservatism of ecological niches in evolutionary time. *Science* **285**: 1265–1267.

**Peterson, A.T. & Vieglais, D.A.** in press. Predicting species invasions using ecological niche modeling. *Bioscience* in press:.

**Short, L.L.** 1982. *Woodpeckers of the World*. Greenville, Delaware: Delaware Mus. Nat. Hist.

**Stockwell, D.R.B. & Noble, I.R.** 1992. Induction of sets of rules from animal distribution data: a robust and informative method of data analysis. *Math. Comp. Simul.* **32**: 249–254.

**Veregin, H.** 1989. Error modeling for the map overlay operation. In: Goodchild, M. & Gopal, S. (eds) *Accuracy of Spatial Databases*: 3–18. New York: Taylor & Francis.

**Walker, P.A. & Cocks, K.D.** 1991. HABITAT: a procedure for modeling a disjoint environmental envelope for a plant or animal species. *Global Ecol. Biogeogr. Lett.* **1**: 108–118.

**Wilson, E.O.** 1988. *Biodiversity*. Washington, DC: National Academy Press.

## APPENDIX

Institutions supplying collection locality data for the present project: American Museum of Natural History, New York; Academy of Natural Sciences, Philadelphia; Bell Museum of Natural History, Minneapolis; British Museum (Natural History), Tring; Carnegie Museum of Natural History, Pittsburg; Canadian Museums of Nature, Ottawa; Denver Museum of Natural History; Delaware Museum of Natural History; Fort Hays State College; Field Museum of Natural History, Chicago; Iowa State University, Ames; University of Kansas Natural History Museum, Lawrence; Los Angeles County Museum of Natural History; Natuurhistorische Museum, Amsterdam; Louisiana State University Museum of Zoology, Baton Rouge; Museum of Comparative Zoology, Harvard University, Cambridge; Moore Laboratory of Zoology, Occidental College, Los Angeles; Museum Nationale d'Histoire Naturelle, Paris; Museum of Vertebrate Zoology, Berkeley; Museo de Zoología, Facultad de Ciencias, Universidad Nacional Autónoma de México; University of Nebraska, Lincoln; Royal Ontario Museum, Toronto; San Diego Natural History Museum; Southwestern College, Kansas; Texas Cooperative Wildlife Collections; University of Arizona; University of British Columbia Museum of Zoology; University of California at Los Angeles; Universidad Michoacana de San Nicolás de Hidalgo; U.S. National Museum of Natural History; Western Foundation of Vertebrate Zoology; and Peabody Museum, Yale University.