Testing methods for predicting mammalian species' responses to 20th century climate change in California

Prepared For:

California Energy Commission

Public Interest Energy Research Program

Prepared By:

Adam B. Smith

Michelle Koo

Karen C. Rowe

Jim Patton

Steven Beissinger

Craig Moritz

Museum of Vertebrate Zoology 3160 Valley Life Sciences Building #3160 University of California, Berkeley Berkeley, CA 94720-3160

March 31, 2011

Summary	3
Introduction	5
Methods	6
Species Distribution Models	6
Species Data	9
Environmental Data	10
Grinnell Surveys and Resurveys1	11
Occupancy Modeling1	11
SDM Validation1	12
Results	14
Western US (PO Assessment)1	14
Grinnell Sites (PA Assessment)1	15
Future Ranges1	16
Discussion1	17
Previous Comparable Studies1	17
Which Species Model Well and Why?1	17
Presence-only vs. Presence-absence Assessments	19
Projecting across Time	20
Where Do SDMs Fail?	21
Conclusions	22
Literature Cited	23
Figures	28
Tables	34
Supplemental Tables	38
Supplemental Projection Maps for Each Species	40

Summary

Governmental and conservation organizations commonly use species distribution models (SDMs) to project forward in time and thus anticipate the effects of climate change on species' ranges. However, without data from the future with which to assess model performance, the reliability of models projected across eras can only be inferred using the performance of models projected and tested within the same era. Models could perform differently when projected across time if species' ranges are influenced differently by the same predictors used to train the model, or if factors not included in the model affect the expansion or contraction of a species' range. Even when SDMs are used to predict ranges within the same era as the training data, their accuracy is often assessed using presence-only (PO) data, but this does not indicate how well they predict absences. And even if presence/absence (PA) data is available, the problem of non-detections can bias assessments of model performance.

Here, we validate five SDMs (BIOCLIM, boosted regression trees, generalized additive models, generalized linear models, and MAXENT) for 26 mammalian species in the western United States with data collected in two time periods, from 1900 to 1939 and from 1970 to 2009. Model performance is assessed across the study area using commonly-used metrics of reliability when only PO is available. We also assess model performance for selected sites within the Sierra Nevada first surveyed by Joseph Grinnell and colleagues, and resurveys of the same or similar sites we have conducted in the past decade. For these sites we use occupancy modeling to infer the probability of true absence in order to obtain high-quality PA data.

We find that the SDMs tested here are adequate for differentiating true presences from randomly chosen sites across the study region (mean \pm standard error of AUC across all SDMs and projections = 0.84 \pm 0.02), and that this performance is a function of the type of projection (historic-to-historic, historic-to-modern, modern-to-modern, and modern-to-historic), range size of the species, and species identity. Likewise, model performance within an era predicts performance between eras (mean *r* across both cross-era projections = 0.96, P<0.001). Other factors related to species' experience of the

environment, reproduction, or other life history attributes (hibernation, daily activity cycle, young per year, and body mass) were not important in predicting model performance, even though some species modeled poorly or well, regardless of the SDM or projection type.

In contrast, when tested with high-quality PA data, SDMs performance was mediocre (mean AUC = 0.75 ± 0.02) and depended on the projection and species, but not the SDM. Moreover, within-era performance correlated less strongly with cross-era performance (mean r = 0.56, P<0.01). Projection maps thresholded at the value which maximized the true positive and true negative rate for PA data tended to overpredict presences more than they overpredicted absences. Across all species, SDMs, and projections, the mean overprediction (false positive) rate was 0.35 ± 0.01 , while the mean underprediction (false negative) rate was 0.12 ± 0.01 . Thus, projections within the same time period and projections forward in time tend to be overly optimistic for the viability of the species, even when the SDMs predict range contractions.

Finally, we projected species' ranges to future eras centered on 2050 and 2080 under the IPCC's Special Report on Emissions Scenarios' A2 scenario for three general circulation climate models (CCCMA, CSIRO and Hadley). Maps of each species' historic, modern, and future ranges are included as a supplement. Given the observed performance of SDMs in our study, we suggest that the veracity of these projections be interpreted with caution, especially since they tend to overpredict presences into the future.

Managers should consider the type of error (undue optimism or undue pessimism) that would be more serious in their situation and interpret SDM results in that light. If procedures similar to those here are used, predictions can be expected to reflect species' actual presences and absences with mediocre confidence, and they will tend to do so in a manner that overstates the true range size of a species.

Introduction

Anthropogenic global change promises to rewrite the biogeography of Earth's species, with some expected to gain, some lose, and some shift their current distributions (Parmesan and Yohe 2003). As a result, conservation planners require a reliable set of methods to predict future distributions to identify at-risk species and prioritize efforts (e.g., Carroll et al. 2010). One increasingly popular method uses species distribution models (SDMs), which relate limited observation data of a species with climate variables and other data indicative of habitat quality to produce maps of environmental suitability for the species (Guisan and Zimmermann 2000).

Despite their apparent utility, the reliability of SDMs to predict ranges across time periods relevant to conservation remains relatively unknown (Kharouba et al. 2009). Scores of studies have assessed performance of different algorithms using within-era validation, which tests models against data withheld from the data in the same era that is used to train the models (e.g., Elith et al. 2006). However, data from the same time and region is not independent of training data, meaning that within-era assessments of SDM performance may give overly optimistic estimates of cross-era performance (Araújo et al. 2005). Alternatively, cross-era validation ensures independence between training and test data, but it requires data to be available on both ends of the time period of interest and so is rarely available. Projections across time would be expected to differ from within-era projections if predictors change in their relative influence on species' ranges (Davis and Shaw 2001), the correlational structure changes between predictors (Jackson and Overpeck 2000), or factors not considered in the model are uncorrelated but influence species' ranges (Guisan and Zimmermann 2000). Hence, it would be useful to know if within-era validation indicates reliable performance when models are projected across eras.

Fortunately, museum records provide a unique opportunity to assess SDM performance across time. Here, we profit from the foresight of Joseph Grinnell and colleagues who conducted systematic surveys of vertebrates in California and nearby States between 1900 and 1939 (Grinnell and Storer 1924). Grinnell and his colleagues' meticulous field notes (~50,000 pages) and specimens (>80,000), preserved at the University of California, Berkeley's Museum of Vertebrate Zoology, have allowed us to resurvey the same or similar sites since 2003 (Moritz et al. 2008).

A common challenge faced by those assessing SDM performance is that most available data indicates where species are present, but not where they are absent (Elith et al. 2006). Even when presence/absence data is available, "absence" records can be confounded by the problem of detection: a species may be present at a site even if it is not detected (MacKenzie et al. 2006). One alternative is to use the occupancy modeling framework, which uses the pattern of detection across multiple surveys at the same sites to infer the probability of true absence at sites where the species was not detected (Tingley and Beissinger 2009, Kéry 2011). Grinnell's detailed field notes and our modern resurvey design have allowed us to apply the occupancy framework to historic and modern records, enabling us to have confidence in assigning true absences.

Here, we use within- and cross-era validation to assess the reliability of five SDMs trained on records of 26 mammalian species from 1900 to 1939 (the "historical" era) and 1970 to 2009 (the "modern" era) in the western conterminous United States. We assess model performance using presence/absence metrics from sites located in the Sierra Nevada range, where absences are inferred with high confidence from occupancy modeling. We also assess models using presence-only metrics commonly used in other studies. Withinera performance is compared to cross-era performance to determine if the former serves as a reliable guide to the latter. Finally, in the context of these results we project ranges to a future era centered on 2050 and 2080 under high greenhouse gas emissions scenarios modeled with three general circulation models (GCMs).

Methods

Species Distribution Models

We compared performance of five SDMs: BIOCLIM, boosted regression trees (BRTs), generalized additive models (GAMs), generalized linear models (GLMs), and maximum

entropy (MAXENT). A multi-model, within-era assessment found GAMs, BRTs, and MAXENT to be among the top performers, GLM to be of intermediate performance, and BIOCLIM among the poorest (Elith et al. 2006). All models were implemented using custom code and the dismo (Hijmans et al. 2011), raster (Hijmans and van Etten 2011), MASS (Ripley 2011), mgcv (Wood 2006), gbm (Ridgeway 2007), e1071 (Dimitriadou et al. 2011), and PBSmapping (Fisheries and Oceans Canada 2010) packages in R Ver. 2.12.2 (R Core Development Team 2011).

BIOCLIM assumes that predictors are equally important and do not interact (Busby 1991). In the single-variable case, the environmental suitability for a species at a site is estimated as a linearly decreasing function of the distance of the site in environmental space to the median of the distribution of all sites (Hijmans et al. 2011). When more than one variable is used, each is assumed equally important and the probability of presence is the minimum suitability score across all variables.

BRTs use a series of linked classification and regression trees, with each tree trained on the residuals from the previous tree (Hastie et al 2001). We implemented BRTs with a modified version of the gbm.step function from Elith et al. (2008), which calculates the optimal number of trees for a model. This function requires users to specify a learning rate, tree complexity, and bag fraction (proportion of data withheld for testing each tree iteration). We searched the model performance space using the range of learning rates (0.0001 to 0.1), tree complexities (2 to 15), and bag fractions (0.5 to 0.7) suggested by Elith et al. (2008) and chose the model with the parameters combination that gave the highest AUC using 5-fold cross-validation tested against random background sites once models with <1000 trees had been removed.

GAMs apply non-linear smoothers on predictor variables before relating them to the dependent variable (Wood 2006). Here, we used cubic splines as smoothers. Models were applied with shrinkage, which allows a smoother term's influence to go to zero if it is unimportant relative to the other terms in the model. Variables were first sorted on AIC when used alone in a GAM, then added to the model if there were \geq 30 presences in the

training set for each variable (preliminary work indicted models with fewer than 30 sites per initial term performed poorly or were unstable).

GLMs are a regression technique that assumes a linear relation between the predictors and species' occurrence. We used linear, quadratic, and two-way interaction terms selected on the basis of forward and backward stepwise model selection using AIC (Anderson et al. 2000). Initial model form was determined by entering each term along in a GLM (forcing inclusion of linear terms when testing quadratic and interaction terms), then building the initial model from lowest to highest AIC, provided there were \geq 20 presences sites per term in the training set. Terms in the initial model were then dropped using automatic forward and backward AIC-based model selection until the optimal model was found.

MAXENT is a machine-learning algorithm that first estimates the probability of each predictor for each presence site given the distribution of the predictor across the study area, and then uses Bayes' Theorem to invert the relation, yielding the probability of presence given the environment (Phillips et al. 2006, Phillips and Dudík 2008). The distribution is derived using maximization of information entropy, which produces the mathematically smoothest distribution possible given constraints (e.g., predictor mean, variance, covariance with other predictors). We implemented MAXENT using Ver. 3.3.3e (Phillips et al. 2009) called from R using the dismo package (Hijmans et al. 2011).

To account for sampling bias in geographical and hence environmental space (Phillips et al. 2009), we used records from all mammals in the study region as target background sites for all SDMs except for BIOCLIM, which does not require background data. For GLMs we used the same number of target background sites as we had records for the species in question to ensure that prevalence (proportion of presences to presences plus background sites) was 50% since imbalanced prevalence can bias model performance (McPherson et al. 2004). We also used an equal number of presence and random background sites for training BRTs in order to reduce computing time.

Species Data

To train the SDMs we used presence-only records from the entire conterminous US west of the eastern border of the Rocky Mountains (103.77 W; Figure 1a and b). We downloaded all mammalian records from Arctos (http://arctos.database.museum/) and MaNIS (www.manisnet.org) in July, 2010. We kept specimens collected from 1900 to 1939 and from 1970 to 2009, inclusive, to comprise our historic and modern eras, respectively. We also included records from 2010 that were part of our resurveys.

We initially retained all records with coordinate uncertainty \leq 5000 m. Graham et al. (2008) found that similar levels of coordinate uncertainties did not appreciably degrade model performance. Point maps for each species were checked by Jim Patton, and records outside the known range of each species were double-checked with the specimens in the Museum of Vertebrate Zoology when possible or discarded. We also checked for outlying records in environmental space. Models are trained in environmental, not geographic space, so spatial error may not affect model performance even if incorrect georeferencing misplaces a record, so long as the site is similar climatically to the one in which the animal was actually captured. Visual inspection suggested that outlying records were generally >2 SD from the mean, so for each species we removed records that were below the 0.25th and above the 99.75th percentiles in either mean annual temperature or precipitation relative to the other sites where the species was located. We also removed records collected before year 2000 with coordinate uncertainties <3 m since these had unlikely accuracies, and we removed records from 2000 onward with coordinate uncertainties >200 m since the widespread use of GPS improved spatial accuracy. For each species, records were further thinned so that no presence points were within 1 km of one another. To ensure a fair comparison between models built for each era, for each species we subsampled the era with the greater number of sites so each had an equal number of presences.

The final data set had 26 mammalian species with \geq 37 records in each era (median = 116 sites per species in an era, minimum = 40, maximum = 1059), which allowed \geq 30 sites for model training with 5-fold cross-validation data, a number of sites found to be adequate in

other assessments (Wisz et al. 2008). Each SDM was trained and tested using the same training and test data folds for each species.

Environmental Data

We used 800-m resolution PRISM climate layers, averaged across the two eras in our study, as predictors (Daly et al. 2000). PRISM is an expert-tuned meteorological interpolation system with predictions based on observed weather measurements. From these layers we derived 19 BIOCLIM variables (Nix 1986), and kept those we expected to be biologically meaningful to the species and that had cross-correlations >-0.7 and <0.7 (Tingley and Herman 2009). In the end, we were left with 9 predictors averaged across the years in each era: mean diurnal temperature range (BIO02), minimum and maximum temperature of the warmest/coldest month (BIO05, BIO06), temperature annual range (BIO07), isothermality (BIO2/BIO07), precipitation of the wettest/driest month and warmest quarter (BIO13, BIO14, BIO18), and the coefficient of monthly precipitation (BIO15).

We also projected the ranges to two 30-yr future eras, centered on the years 2050 and 2080. We examined the A2 scenario from the IPCC's 4th Assessment Report's Special Report on Emissions Scenario in which global regions develop at different paces and greenhouse gas emissions remain relatively unabated (Nakicenovic et al. 2000). This scenario is a "high-emissions" future, but current emissions are higher than originally projected in the A2 scenario (Raupach et al. 2007). We examined sets of climate surfaces for each era produced by three GCMs: the CGCM3 model from the Canadian Centre for Climate Modeling and Analysis (CCCMA); the Mark 3.0 model from the Commonwealth Scientific and Industrial Research Organization of Australia (CSIRO); and one from the Hadley Centre for Climate Prediction and Research (Hadley). We performed these projections using MAXENT since it is currently the most commonly used SDM. Table 1 gives mean values for each of our predictor variables under the historic, modern, and two future eras across our study region for each source. Between the historic and modern era, mean annual temperature and precipitation increased 0.44°C and 28 mm, respectively. Amongst the three models, CCCMA predicts the wettest and coolest future for California (though mean temperature still rises), the Hadley GCM predicts the hottest and driest, and the CSIRO GCM predicts precipitation will remain roughly constant but temperature will rise at first slowly so it's comparable to the CCCMA output but then rise quickly to make it comparable to the Hadley model.

Grinnell Surveys and Resurveys

Between 1900 and 1939 Grinnell and colleagues conducted an extensive census of vertebrate wildlife in California. Here we focus on independent transects that straddle the elevational gradient along Sierra Nevada range in each region and have experienced relatively little development: the Lassen region (surveyed at elevations spanning 80 to 2510 m and centered on what is now Lassen National Park and Lassen National Forest), Yosemite (from 50 to 3280 m; focused on Yosemite National Park), and the Southern Sierras (from 120 to 3640 m; including Sequoia and Kings Canyon National Parks and Sequoia, Sierra, and Inyo National Forests). We mined Grinnell's historical field notes and specimen records to ascertain species caught, trapping effort, and the pattern of captures across nights at each site to use for occupancy modeling (see below). Since 2003 we have been resurveying these and similar sites across the same transects, yielding 89 historical sites and 136 modern sites for occupancy modeling. Grinnell's records and the resurvey data were also incorporated into the larger dataset used to train the SDMs but comprise a small portion of the total records used to train and test them.

Occupancy Modeling

Occupancy modeling uses the pattern of detections and non-detections at a site and across multiple sites to infer the probability that a species is truly absent when it is not detected (MacKenzie et al. 2006). As a simple example, if a species has a "01011" capture history across five nights at a site, where 0 indicates non-detection and 1 indicates detection, we can infer that the probability of detecting the species in any single day is 3/5 = 0.60. Hence, if we were to census an environmentally similar site for five nights within the species' range and obtain all zeros and thus assume the species was absent from that site, we would be correct with a probability of $1 - (1 - 0.60)^5 = 0.99$. Here we calculated the probability of true absence for each site as

Eq. 1
$$p^* = 1 - \prod_{n=1}^{N} (1 - p_n)$$

where *N* is the total number of nights a sites was censused and *p* the probability of detecting a species at that site in night *n*, which was modeled with covariates that accounted for the era in which the trapping was performed and for trapping effort (number of traps and log(number of traps); Moritz et al. 2008).

We implemented the occupancy framework using MARK (White and Brunham 1999), the RMark package, and R code modified from Royale and Dorazio (2008). We first constructed dectability models using all possible combinations of linear and quadratic transformations of the covariates and then used AIC-weighted model selection to obtain values for *p* (Moritz et al. 2008).

SDM Validation

For each species we evaluated SDM performance for two within-era projections (historicto-historic, abbreviated "HH"; modern-to-modern, "MM") and two cross-era projections (historic-to-modern, "HM"; and modern-to-historic, "MH").

SDMs were evaluated in two ways, one using presence-only data (PO) and the other presence-absence (PA) data. The first assess each SDM's ability to discriminate between sites randomly drawn from the entire study area and sites where the species were known to be present (PO assessments). We measured model performance using the area under the receiver-operator curve (AUC) and the True Skill Statistic (TSS). AUC ranges between 0 and 1, with 0.5 indicating the model performs no better than random and values close to 1 indicating reliable fit (Fielding and Bell 1997). While not as informative a measure of model performance as a metric using true absences, AUC measured against random background points is used in many studies for lack of validated absence data (Phillips et al. 2006). In this context AUC yields the probability that a model will give a higher score to an average presence site than an average background site (Elith et al. 2006). We calculated AUC using a number of randomly selected sites equal to the number of test presence sites to obviate AUC's sensitivity to prevalence (McPherson et al. 2004). We also calculated the true skill statistic (TSS), a measure of performance independent of prevalence and equal to the proportion of correctly classified presences plus the proportion of correctly classified pseudoabsences minus 1 (Allouche et al. 2007). For each species-SDM-projection

combination we calculated mean AUC or TSS across the five-fold test sets. The PO assessment was conducted for all points across the Western US.

To ascertain whether projection type (MM, MH, etc.) or SDM affected accuracy, we used a linear mixed model with species as a random intercept. We also included as covariates traits which we expected to affect model performance because they relate to species' experience of environmental extremes, have been shown in other studies to affect model performance, or relate to other aspects of species' autecology: adult body mass, number of young per breeding female per year, hibernation (hibernator/non-hibernator), activity cycle (nocturnal/diurnal/both), and range area (approximated by the area of the convex hull of training sites, averaged across all k-folds). These data were collated from Moritz et al. (2008), the PanTHERIA database (Jones et al. 2009), the American Society of Mammalogists' Mammalian Species series, and Jim Patton. The full model was tested against a model with the same fixed structure but without species as a random effect; AIC of the former was substantially lower (Δ AIC=-303.3), so species was retained as a random effect in all subsequent models (Zuur et al. 2009). The best model was then selected with manual, AIC-based, backwards stepwise evaluation (Zuur et al. 2009). AUC was transformed with $\sin^{-1}(\sqrt{x})$ and TSS with (x/2 + 0.5) before analysis to reduce heteroscedascity; visual inspection of residuals indicated the best-fitting models behaved well.

Second, we evaluated the ability of models to differentiate between known presences and absences in the Sierra Nevada (the PA assessment). Here, we inferred absence of a species at each of the Grinnell survey or resurvey sites using Eq. 1, and assumed true absence if p^* was ≤ 0.10 (Rubidge et al. 2010). For each combination of species, projection, and SDM, we calculated AUC and TSS as above but using absences instead of random background sites. We analyzed results with a linear mixed model using the same procedures described in the paragraph above.

To see if within-era performance predicts cross-era performance, for each SDM we calculated the Pearson correlation coefficient for each within- vs. cross-era AUC.

Finally, we projected the ranges of each species using MAXENT from the modern era to the 2050 A2 climate scenario for each of the three GCMs. We thresholded the maps to convert them into presence/absence maps based on the threshold that maximized the sum of sensitivity and specificity from the PO assessment for the within-era projection from the modern era. For these projections we used MAXENT models trained on all possible points to calculate range size (the modern era was not subsampled if it had more records than the historic era).

Results

AUC values for both PO and PA assessments for all species, all SDMs, and all projections are given in Supplemental Tables 1 and 2 but are discussed separately in the following two sections.

Western US (PO Assessment)

Averaged across species and SDMs, mean AUC of within-era projections was 0.85 ± 0.01 (± standard error) and cross-era projections was 0.83 + 0.01, but varied widely between SDMs (Figure 2). AUC values between 0.9 and 1 are indicative of excellent discriminative ability, between 0.9 and 0.8 of good ability, between 0.7 and 0.8 of mediocre ability, and below 0.7 of poor ability (Swets 1988). The highest ranked linear mixed model included SDM, projection, and range size as fixed effects and species as a random effect (Table 2). None of the life history traits or range size occurred in the best model. TSS did not qualitatively differ in any of the analyses, so it will not be discussed further.

In an analysis of variance on AUC that used only SDM and projection but did not control for species, SDM performance and projection were significant, but the interaction was not (Table 3). In order from highest to lowest performance, Tukey's HSD test ranked BRTs, GAMs, and MAXENT as statistically equivalent in performance and BIOCLIM and GLMs as equivalent, with the latter set two performing less well than the former. Though projection was significant in the model, only MM vs. MH was marginally significant (0.05<P<0.06) in a Tukey's HSD test.

In Table 4 we report the Pearson correlation coefficients between within- vs. cross-era AUC and TSS, which indicate whether within-era performance is a good predictor of cross-era performance. In every case within-era performance significantly and positively correlated with cross-era performance. Within- vs. cross-era correlations were highest for GLM (0.97 for HH vs. MH and 0.96 for MM vs. HM).

Grinnell Sites (PA Assessment)

In contrast to the PO data, we did not equilibrate test presence/absence points for the PA assessment because some species had extremely low absence and/or presence rates in each era (e.g., in the historic era *Otospermophilus beecheyi* was present at 8 sites and absent at 1). Prevalence had a negative effect on AUC scores for tests on the historic data, but not the modern data (Pearson correlation coefficients between AUC and for each projection prevalence were: HH: -0.42*, MM r=-0.13, HM: 0.01, MH: r=-0.44*, where * indicates P<0.05). Hence, for species with high prevalence in the historic era our results may be downwardly biased, and for species in the historic era with low prevalence, our results may be upwardly biased (species with historic prevalence ≤ 0.2 and ≥ 0.8 are marked in Table 6 where PA AUC is reported for each projection-species-SDM combination). Under our criteria for confidence in assigning an absence to sites where the species was not detected, *Peromyscus californicus* had no absences in the historic era, and *Sorex trowbridgii* had no absences in the modern era, disallowing calculation of AUC in the respective era for these species.

Average AUC for the Grinnell sites was substantially lower than for the PO assessment, with a mean of 0.76 ± 0.01 for within-era projections and 0.73 ± 0.01 for cross-era projections (Figure 3). Several species had AUC values ≤ 0.50 for particular projection-SDM combinations, meaning that they were worse than random (Figure 3, Supplemental Table 2). Poor performance was especially common for *Peromyscus maniculatus* (eight projection-SDM combinations with AUC ≤ 0.50), *Neotoma fuscipes* (5 combinations), and *O. beecheyi* (4 combinations), and *Microtus californicus* (1 combination). Of these species, only *P. maniculatus* had AUC < 0.60 in the PO assessment (Supplemental Table 1); the remainder had scores ≥ 0.90 . In the PA assessment each SDM-species combination had at least one projection with AUC ≤ 0.50 (Figure 3, Supplemental Table 2). GAMs were notable in having the fewest poor performers (*P. maniculatus* for the HH projection), and GLMs the worst (one to three species with AUC ≤ 0.50 for each projection).

In contrast to the PO assessment, the top-ranked models did not include SDM as a factor (Table 2). Like the PO assessment, they did include projection as a fixed effect and species as a random effect. There was less separation by AIC between the top PA models compared to the top PO models. The other top models included range area, hibernation, and activity cycle as important covariates.

An analysis of variance on AUC using SDM and projection by themselves and interacting, while ignoring species, found no significant differences between models and projections (Table 3), in contrast to the same analysis for the PO assessment where SDM and projection were significant.

Correlations between within- and cross-era performance for each model and for all combinations of species and SDM were lower than for the PO assessment (Table 4). BIOCLIM had the lowest correlations, whereas GLMs had the highest (≥0.71 in both cases).

Future Ranges

Attached to this report are MAXENT projections for each species for HH, MH, MM, HM, and the six modern-to-2050 and -2080 A2 scenarios for the CCCMA, CSIRO, and Hadley GCMs. These maps are thresholded at the value that most maximizes sensitivity plus specificity using the PO assessment in the era in which the model is trained. Though not necessarily as accurate as the PA assessment for some species, we used the PO assessment for thresholding because the low number of presences and absences for some species may have made the threshold unstable. We used MAXENT because it is currently the most commonly used SDM.

Discussion

Previous Comparable Studies

Despite the prominent use of SDMs in conservation planning, to date only three other published studies have validated SDM performance across time scales relevant to conservation. Araújo et al. (2005) used GAMs, GLMs, artificial neural networks, and classification and regression trees to model the range dynamics of over 100 British birds across 20 yr. They found that models that performed well within-era also performed well across eras, but performance was always lower when projecting across time. In their study neural networks and GAMS performed best and second-best, respectively, when projecting across eras. Kharouba et al. (2009) used MAXENT to model ranges of almost 140 butterflies across nearly the same time period as our study. They found that within-era accuracy predicted cross-era accuracy, but less so for widely-ranging species or species with low sampling density. Finally, Dobrowski et al. (*in press*) modeled the past and current distributions of 133 plant species in the Sierra Nevada over a 75-yr time span. They found that species' traits (dispersal, endemicity, and fire-adaptation) mattered more to model performance than type of SDM (they examined BRTs, GAMs, GLMs, and random forests).

Which Species Model Well and Why?

In our study the performance of SDMs depended on the species being modeled, as well as the type of assessment (PO vs. PA). The PA assessment should be more indicative of real model performance because it allow direct estimation of false presence and false absence rates, but we conducted the PO assessment in addition because most users of SDMs do not have high-quality presence/absence data.

Both assessments suggest that species identity matters more than SDM to accuracy of a projection (Table 2). This can be seen in Figure 6, where the mean AUC for each model is plotted for each species. In general, there is more variation between species than between models for the same species (GLMs tend to be the exception, having erratic performance for some species). For example, it matters more that *Microtus longicaudus* is being

modeled and less that BRTs, GAMs, or another model is being used. Thus, different species model poorly or well, regardless of the model used. The top PO and PA models included species as a random factor, which is a statistical means to account for differences between species even if the relevant differences are unknown (Table 2). The importance of species as a random factor in these models remained despite including traits which a priori would seem to affect species' exposure to environmental extremes (hibernation and daily activity cycle), relate to reproductive activity (young per year), or are broadly indicative of the range of environmental conditions the species may experience (range size). Body mass also did not affect projection accuracy, despite its correlation with a large suite of life history traits (Brown and West 2000). The top five models for PA AUC (bottom half of Table 2) did include some of these traits (range area, hibernation status, and daily cycle), but the top model included only projection and a random term for species. We interpret these results as weakly indicative of a relationship between SDM performance and the traits examined here.

These findings compare and contrast with those of others who have examined relationships between SDM performance and variables related to species' traits. McPherson and Jetz (2007), examining over 1000 avian species, found range size, some types of habitat preference, and endemism to be related to performance, but other traits like mass and trophic rank were not. Similarly, we found a significant association between range size and PO model performance, but range size did not influence PA performance (Table 2). Dobrowski et al. (in press), projecting distributions of Californian plants across a time span similar to the one in this study, found endemicity, dispersal capacity, and fireadaptation more important to projections than SDM. Without detailed study it is not possible for us to ascertain dispersal capacity of the mammals studied here, but it should generally scale with body size (Brown and West 2000), a variable found not to be important in our study. Likewise, endemicity should be correlated with range size, which was in the top model in the PO assessment and the second-top model in the PA assessment. In these models large range size negatively impacted PO and PA performance, perhaps because species with larger ranges are less restrained by climate. (Fire adaptation is uncommon amongst mammals and so does not apply to them as it does plants.) Tellingly,

Moritz et al. (2008), examined elevational range shifts of 21 of the 26 mammals examined here and found weak associations between the same life history traits we included in our models and whether or not the species shifted elevational range in and around Yosemite National Park.

Hence, our failure to find strong correlations between species' traits and model performance, despite the fact that some species model well and others model poorly independently of SDM, suggests that there is something as of yet unquantified about each species that is makes it good or poor at being modeled. Also notable is the inconsistent relationship between species in the same genus. For example, the three species of woodrats, *Neotoma cinerea, N. fuscipes,* and *N. macrotis,* respectively have good, poor, and average performance, averaged across all SDMs and projections (Table 6). Hence, whatever it is that makes some species' models perform well or poorly does not seem to be phylogenetically conserved.

Within-era performance did not strongly correlate with cross-era PA performance (Table 4), meaning that the accuracy of projections to a new time period cannot be reliably ascertained without knowledge of species' distributions in the target time period. Of the SDMs examined here, GLMs had the highest within- vs. cross-era performance correlation (Table 4).

Presence-only vs. Presence-absence Assessments

If we did not possess high-quality absence data, we would conclude that BRTs, GAMs, and MAXENT are the most reliable models (Figure 2) and that within-era performance predicts cross-era performance (Table 4). Many analyses using PO data or PA data with no confidence assigned to absences have come to similar conclusions. For example, Elith et al. (2006) conducted a within-era performance assessment of 16 SDMs trained on PO data and tested on PA data. They found BRTs, GAMs, and MAXENT worked best, but this result is predicated on the quality of their absence data. Since they did not use occupancy modeling to assign confidence to assignations of absence, it is possible that some of their "absence" points were actually occupied by the species in question. Further work testing SDMs against high-quality PA data is required (Kéry 2011).

Projecting across Time

Projecting across time is fundamentally different from projecting within the same time period for several reasons. First, species can adapt to changing environmental conditions, meaning that the functional relationship estimated in one time period may not apply to another (Davis and Shaw 2001). Second, species can disperse into previously unoccupied regions or become extirpated from areas where they once were. If such range shifts occur independently of predictors used in a SDM, then the SDM will either over- or underpredict occurrence.

We did not find that within-era performance accurately predicted cross-era performance (bottom half of Table 4). Although some of the correlations between within- and cross-era performance for the PA assessment are statistically significant, they are also low (r<0.5), meaning that despite a positive association between within- and cross-era reliability, one cannot be used to predict the other with much confidence. In contrast, within- and cross-era PO performance are highly correlated (top half of Table 4), but this result is less relevant to actual model performance since the PO assessment cannot differentiate between presences and absences.

Visually most of the MM and HM projections tend to match one another except in fine detail (see supplemental maps). Similarly, all three modern-to-2050 projections tended to agree visually, though when they were different the Hadley projection was most dissimilar to the other projections. For example, *Callospermophilus lateralis* experienced dramatic contraction in the Sierras and from northwestern California under Hadley but not so much under CCCMA and CSIRO. In contrast, the modern-to-2080 projections suggest stark changes from modern distributions for some species, with the differences between the three GCMs much more accentuated. For example, *Chaetodipus californicus* is expected to "retreat" to the southern part of the Great Valley and the Transverse ranges under the CCCMA scenario, whereas it retreats from the southern part of the Great Valley under the CSIRO and Hadley scenarios.

These results are predicated on our choice of threshold to designate presence/absence (Thuiller 2003). We chose to use the threshold that maximized the sum of the proportion

of true positive and true negative (randomly-located sites) predictions, a method that a priori does not weigh errors toward over- or underprediction but may nevertheless produce under- or overpredictions. We could have used a threshold that equalized the rates, which would have generally increased the threshold value and thus predicted smaller area of occupancy in historic, modern, and future eras. Likewise, we could have used the lowest value of a presence site as a threshold (e.g., Algar et al. 2009, Kharouba et al. 2009, Rubidge et al. 2011), but owing to the possibility that some sites may have been mis-georeferenced (see the "Species data" section in "Methods"), we decided to use a threshold that put less weight on a single presence point.

Our projection maps are also dependent on the SDMs' ability to extrapolate into climate regimes under which they were not trained. Future climates do indeed differ from the current and historic climates (Table 1). When training the models we used the "clamping" option, which keeps the species' response to a variable constant at the last estimated value when the variable passes outside the range on which it was trained. In reality we would expect some species' responses to increase or decrease as temperature and precipitation went above or below values observed in the historic and modern eras.

Our results are also somewhat affected by the negative correlation between prevalence and AUC in the historic period. However, we believe this to be of less concern where the same true for the modern period (which it is not). Since managers are most commonly interested in projecting to the future and not the past, our projections in this direction are less likely to be biased relative to prevalence. Although SDMs do not "care" from a statistical point of view whether they are projecting forward or backward in time, differences between the processes that cause range contractions and expansions make time irreversible in this context, and so we should expect differences between predictions and retrodictions. Hence, our results should be robust for projections forward in time.

Where Do SDMs Fail?

SDMs can fail by either predicting presences where the species is absent (false positives or overprediction of range) or predicting absences where the species is present (false negatives or underprediction of range). Figures 4 and 5 respectively show the false

positive and false negative rates (FPR and FNR) for all SDMs against the PA data (the true negative and true positive rates are 1 minus these values, respectively). The mean false positive rate across all species, SDMs, and projections was 0.35 ± 0.02 , while the false negative rate was substantially lower, 0.14 ± 0.01 . Hence, models tend to be more generous in allocating presences than absences. This result could be in part due to our choice of threshold, which we set at the value which maximized the FPR plus FNR for all presences in the western US versus random background points for "absences." Since the FPR is higher than the FNR, it seems that this threshold is too high. Nonetheless, were we not in the position to have high-quality PA data, the threshold we chose would be a defensible choice.

BIOCLIM consistently overpredicts areas of presence (Figure 4), whereas the other models perform roughly equally. GLMs tend to overpredict absences (Figure 3), though the differences between models is not a striking.

Thus, in some cases correlative SDMs like those used here may overpredict species' ranges, even when they are experiencing contractions. As a result, prediction of range retraction may actually understate the degree to which the range is collapsing, and predictions of expansion should be viewed with skepticism. Hence, the projection maps we present in the supplemental information may overstate the ranges of species.

Conclusions

Our results suggest that using SDMs to project across time fares worse than projecting within the same time period, and that within-era performance is not a good predictor of cross-era performance. We also show that presence-only data gives a more optimistic picture of model performance than presence/absence data, where absences are assigned with statistical confidence and not just assumed from lack of observations of the species.

These results are disappointing in light of the use of SDMs by governmental and conservation organizations to predict species' future responses to climate change. Our exercise suggest that SDMs are more reliable than a random assignation of presences and

absences (for most species), but different species perform well or poorly independently of the SDM used, probably because species differ in how much climate determines their ranges. We also note that if avoiding poor performance is more important than achieving high performance, GAMs may be the best modeling option because they produced fewer poor models than other SDMs.

In conclusion, we recommend care be taken when projecting SDMs across time and that managers consider whether over- or underprediction is a more egregious mistake would occur given situational goals. Using procedures similar to those here, models will tend to overpredict ranges than underpredict them, meaning that predictions of range collapse actually understate the degree to which the species are reduced. Regardless, model performance across time is not strongly related to performance within the same time period.

Literature Cited

- Algar, A.C., Kharouba, H.M., Young, E.R., and Kerr, J.T. 2009. Predicting the future of species diversity: Macroecological theory, climate change, and direct tests of alternative forecasting methods. Ecography 32:22-33.
- Allouche, O., Tsoar, A., and Kadmon, R. 2007. Assessing accuracy of species distribution models: prevalence, kappa and the true skill statistic (TSS). Journal of Applied Ecology 43:1223-1232.
- Anderson, D.R., K.P. Burnham, and W.L. Thompson. 2000. Null hypothesis testing: Problems, prevalence, and an alternative. Journal of Wildlife Management 64:912-923.
- Araujo [Araújo], M.B., R.G. Pearson, W. Thuiller, and M. Erhard. 2005. Validation of speciesclimate impact models under climate change. Global Change Biology 11:1504-1513.
- Brown, J.H. and West, G.B. 2000. Scaling in Ecology. Oxford University Press, Oxford.
- Busby, J. R. 1991. BIOCLIM a bioclimate analysis and prediction system. In: Margules, C. R. and Austin, M. P. (eds.), Nature conservation: Cost effective biological surveys and data analysis. CSIRO, pp. 64–68.

- Carroll, C., Dunk, J.R., and Moilanen, A. 2010. Optimizing resiliency of networks to climate change: Multispecies conservation planning in the Pacific Northwest, USA. Global Change Biology 16:891-904.
- Daly, C., G.H. Taylor, W. P. Gibson, T.W. Parzybok, G. L. Johnson, P. Pasteris. 2001. Highquality spatial climate data sets for the United States and beyond. Transactions of the American Society of Agricultural Engineers, 43:1957-1962.
- Davis, M.B and R.G. Shaw. 2001. Range shifts and adaptive responses to Quaternary climate change. Science 292:673-679.
- Dimitriadou, E., Hornik, K., Leisch, F., Meyer, D., and Weingessel, A. 2011. "e1071", ver. 1.5-25, package for R. http://cran.r-project.org/.
- Dobrowski, S.Z., Thorne, J.H., Greenburg, J.A., Safford, H.D., Mynsberge, A.R., Crimmins, S.M., and Swanson, A.K. *In press*. Modeling plant ranges over 75 years of climate change in California, USA: Relating transferability to species traits. Ecological Monographs.
- Elith, J., C.H. Graham, R.P. Anderson, M. Dudi'k [Dudík], S. Ferrier, A. Guisan, R.J. Hijmans, F. Huettmann, J.R. Leathwick, A. Lehmann, J. Li, L.G. Lohmann, B.A. Loiselle, G. Manion, C. Moritz, M. Nakamura, Y. Nakazawa, J.McC. Overton, A.T. Peterson, S.J. Phillips, K. Richardson, R. Scachetti-Pereira, R.E. Schapire, J. Sobero'n [Soberón], S. Williams, M.S. Wisz, and N.E. Zimmermann. 2006. Novel methods improve prediction of species' distributions from occurrence data. Ecography 29:129-151.
- Elith, J., J.R. Leathwick, and T. Hastie. 2008. A working guide to boosted regression trees. Journal of Animal Ecology 77:802-813.
- Fisheries and Oceans Canada. 2010. "PBSmapping", ver. 2.61.9, package for R. http://cran.r-project.org/.
- Graham, C.H., J. Elith, R.J. Hijmans, A. Guisan, A.T. Peterson, B.A. Loiselle, and the NCEAS Predicting Species Distributions Working Group. 2008. The influence of spatial errors in species occurrence data used in distribution models. Journal of Applied Ecology 45:239-247.

- Grinnell, J. and Storer, T. 1924. Animal Life in the Yosemite. University of California Press, Berkeley.
- Guisan, A. and N.E. Zimmermann. 2000. Predictive habitat distribution models in ecology. Ecological Modelling 135:147-186.
- Hastie, T., Tibshirani, R. and Friedman, J.H. 2001. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer-Verlag, New York.
- Hijmans, R.J. and van Etten, J. 2011. "raster", ver. 1.8-9, package for R. http://cran.rproject.org/.
- Hijmans, R.J., Phillips, S., Leathwick, J., and Elith, J. 2011. "dismo", ver. 0.6-3, package for R. http://cran.r-project.org/.
- Jackson, S.T. and Overpeck, J.T. 2000. Responses of plant populations and communities to environmental changes of the late Quaternary. Paleobiology 26 (Suppl. 4):194-200.
- Jones, K.E., et al. 2009. PanTHERIA: a species-level database of life history, ecology, and geography of extant and recently extinct mammals. Ecology 90:2648.
- Kéry, M. 2011. Towards the modelling of true species distributions. Journal of Biogeography 38:617-618.
- Kharouba, H.M., Algar, A.C., and Kerr, J.T. 2009. Historically calibrated predictions of butterfly species' range shift using global change as a pseudo-experiment. Ecology 90:2213-2222.
- MacKenzie, D.I., Nichols, J.D., Royle, J.A., Pollock, K.H., Bailey, L.L., and Hines, J.E. 2006. Occupancy Estimation and Modeling: Inferring Patterns and Dynamics of Species Occurrence. Elsevier, Amsterdam. 324 pp.
- McPherson, J.M. and W. Jetz. 2007. Effects of species' ecology on the accuracy of distribution models. Ecography 30:135-151.
- McPherson, J.M., W. Jetz, and D.J. Rogers. 2004. The effects of species' range sizes on the accuracy of distribution models: Ecological phenomenon or statistical artifact? Journal of Applied Ecology 41:811-823.

- Moritz, C., Patton, J.L., Conroy, C.J., Parra, J.L., White, G.C., and Beissinger, S.R. 2008. Impact of a century of climate change on small-mammal communities in Yosemite National Park, USA. Science 322:261-264.
- Nakicenovic, N., Alcamo, J., Davis, G., de Vries, B., Fenhann, J., Gaffin, S., Gregory, K., Grübler,
 A. et al. 2000. Special Report on Emissions Scenarios: A Special Report of Working
 Group III of the Intergovernmental Panel on Climate Change, Cambridge University
 Press, Cambridge, U.K., 599 pp
- Nix, H.A. 1986. A biogeographic analysis of Australian Elapid Snakes. In. Atlas of Elapid Snakes of Australia. (ed.) R. Longmore pp. 4-15. Australian Flora and Fauna Series
 Number 7. Australian Government Publishing Service: Canberra.
- Parmesan, C. and G. Yohe. 2003. A globally coherent fingerprint of climate change impacts across natural systems. Nature 421:37-42.
- Phillips, S.J., Anderson, R.P., and Schapire, R.E. 2006. Maximum entropy modeling of species geographic distributions. Ecological Modelling 190:231-259.
- Phillips, S.J. and Dudík, M. 2008. Modeling species distributions with Maxent: New extensions and a comprehensive evaluation. Ecography 31:161-175.
- Phillips, S.J., Dudík, M., Elith, J., Graham, C.H., Lehmann, A., Leathwick, J., and Ferrier, S.
 2009. Sample selection bias and presence-only distribution models: Implications for background and pseudo-absence data. Ecological Applications 19:181-197.
- R Development Core Team. 2011. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna.
- Raupach, M.R., Marland, G., Ciais, P., Le Quere, C., Canadell, J.G., Klepper, G. & Field, C.B.
 2007. Global and regional drivers of accelerating CO₂ emissions. Proceedings of the National Academy of Sciences USA 104:10288-10293.
- Ridgeway, G. 2007. "gbm," ver. 1.6-3.1, package for R. http://cran.r-project.org/.
- Ripley, B. 2011. "MASS", ver. 7.3-11, package for R. http://cran.r-project.org/.

- Royle, J.A. and R.M. Dorazio. 2008. Hierarchical Modeling and Inference in Ecology: The Analysis of Data from Populations, Metapopulations, and Communities. Academic Press, San Diego, CA. 444 pp.
- Rubidge, E., Monahan, W., Parra, J.L,. Cameron, S.E., and Brashares, J.S. 2011. The role of climate, habitat, and species co-occurrence as drivers of change in small mammal distributions over the past century. Global Change Biology 17:696-708.
- Swets, J.A. 1988. Measuring the accuracy of diagnostic systems. Science 240:1285-1293.
- Thuiller, W. 2003. BIOMOD optimizing predictions of species distributions and projecting potential future shifts under global change. Global Change Biology 9:1353-1362.
- Tingley, M.W. and S.R. Beissinger. 2009. Detecting range shifts from historical species occurrences: New perspectives on old data. Trends in Ecology and Evolution 24:625-633.
- Tingley, R. and T.B. Herman. 2009. Land-cover data improve bioclimatic models for anurans and turtles at a regional scale. Journal of Biogeography 36:1656-1672.
- White, G. C., and K. P. Burnham. 1999. Program MARK: survival estimation from populations of marked animals. Bird Study 46 Supplement:120-138.
- Wisz, M.S., Hijmans, R.J., Peterson, A.T., Graham, C.H., Guisan, A., and NCEAS Predicting Species Distributions Working Group. Effects of sample size on the performance of species distribution models. Diversity and Distributions 14:763-773.
- Wood, S.N. 2006. Generalized Additive Models: An Introduction with R. Chapman and Hall, Boca Raton.
- Zuur, A.F., N.N. Ieno, N.J. Walker, A.A. Saveliev, and G.M. Smith. 2009. Mixed Effects Model and Extensions in Ecology with R. Springer, New York.

Figures



Figure 1 Point maps of all species used in the analysis. a and b) Maps for historic and modern species' sites, respectively. For each species all points in this region were used for training the SDMs. b and c) Historic Grinnell survey and modern resurvey sites, respectively.



Figure 2 Presence-only AUC for all sites averaged across each k-fold of each species for SDMs trained on: (a) historic data and tested with historic data, (b) modern data and tested with historic data, (c) modern data and tested with modern data, and (d) historic data and tested with modern data. Here AUC is the probability that the SDM can differentiate between a randomly selected presence site and a randomly selected site on the landscape. Performance varies by SDM, projection, range area, and species (Table 2).



Figure 3 AUC for the Grinnell sites where presences are recorded detections and absences are inferred from occupancy modeling. SDMs were trained on a) historic data and tested against historic data; b) modern data against historic data; c) modern data against modern data; and d) historic data against modern data. If the performance of different species is taken account of, none of the SDMs differ in performance. Some of the species-SDM combinations perform worse than random (AUC ≤ 0.50).



Figure 4 The false positive rate for each SDM and each projection using the threshold that maximized the sum of the proportion of sites correctly classified as present and absent in the PA assessment. No model overwhelmingly overpredicts presences, save GLMs for MM (c) and HM (d) projections.



Figure 5 The false negative rate for each SDM and each projection using the threshold that maximized the sum of the proportion of sites correctly classified as present and absent in the PA assessment. BIOCLIM consistently predicts more false absences than the other SDMs



Figure 6 AUC across all species for the PA assessment. AUC is averaged across all projections for each species-SDM combination. Variation between species is generally more pronounced than variation between models within a species (discounting the occasionally erratic performance of GLMs). Thus, some species model well and some poorly, regardless of the SDM employed. Error bars have been removed to aid visual clarity.

Tables

Table 1 Summary of mean annual temperature and precipitation for California acrosshistorical, current, and the 12 future scenarios investigated here.Values are means with rangesin parentheses.

Era/Scenario (GCM)	Mean annual temperature (°C)	Mean annual precipitation (mm)
Historic (1900-1939)	13.91 (-7.31 to 23.97)	592 (0 to 4377)
Modern (1970-2009)	14.35 (-7.14 to 24.94)	620 (0 to 4408)
2040-2069 A2 (CCCMA)	15.94 (-3.50 to 26.40)	708 (43 to 2639)
2040-2069 A2 (CSIRO)	15.90 (-3.30 to 26.50)	580 (30 to 2413)
2040-2069 A2 (Hadley)	16.26 (-3.10 to 27.40)	496 (28 to 1980)
2070-2099 A2 (CCCMA)	17.39 (-2.20 to 28.00)	974 (63 to 3446)
2070-2099 A2 (CSIRO)	17.61 (-1.70 to 28.30)	575 (32 to 2396)
2070-2099 A2 (Hadley)	17.89 (-1.50 to 29.30)	483 (25 to 1976)

Table 2 The best linear models for AUC for the PO and PA assessments. The top five models foreach assessment are shown. Variables in italics were included as random effects. Lower AIC valuesindicate a more parsimonious model, and values different by more than 10 units indicateimportantly different parsimonies. "Projection" indicates the era of the training and test data sets(historic-to-historic, historic-to-modern, etc.). "Activity Cycle" indicates diurnal/nocturnal/both,and "Hibernation" is hibernator/non-hibernator.

Rank	Model	AIC
All Sites	s (PO Assessment)	
1	AUC ~ SDM + Projection + log_{10} (Range Area) + <i>Species</i>	-1208.3
2	AUC ~ SDM + Projection + log_{10} (Range Area) + Hibernation + <i>Species</i>	-1201.4
2	AUC ~ SDM + Projection + log_{10} (Range Area) + Young per Year + <i>Species</i>	-1197.6
4	AUC ~ SDM + Projection + <i>Species</i>	-1196.8
5	AUC ~ SDM + Projection + log_{10} (Range Area) + <i>Species</i>	-1196.2
Grinnel	l Sites (PA Assessment)	
1	AUC ~ Projection + Species	-519.64
2	AUC ~ Projection + log_{10} (Range Area) + <i>Species</i>	-517.80
3	AUC ~ Projection + Hibernation + <i>Species</i>	-514.44
4	AUC ~ Projection + Hibernation + log_{10} (Range Area) + <i>Species</i>	-513.55
5	AUC ~ Projection + Daily Cycle + log_{10} (Range Area) + <i>Species</i>	-507.49

Table 3 ANOVA on AUC for the PO and PA assessments. SDM and projection (e.g., historiuc-tohistoric, historic-to-modern, etc.) were the only factors in the model. Sums of Squares are Type III. Unlike Table 2, species is not controlled for in these analyses. * $P \le 0.05$, ** $P \le 0.01$.

All Sites (PO Assess	ment)			
Source	SS	df	MS	F
SDM	0.38	4	3.35	3.64**
Projection	0.22	3	4.41	2.76*
SDM × Projection	0.14	12	1.17	0.45
Total	13.00	500		
Grinnell Sites (PA A	ssessment)			
Source	SS	df	MS	F
SDM	0.22	4	0.05	1.91
Projection	0.56	3	0.19	1.80
SDM × Projection	0.04	12	0.00	0.14
Residual	13.81	480		

Table 4 Pearson correlation coefficients between withinand cross-era AUC and TSS. For the PO assessment, n=26 except for the analysis using all species-SDM combinations, for which n=130. For the PA assessment, n=24, except for the analysis using all species-SDM combinations, for which n=120. Two species were excluded from the PA assessment because they had no absences in at least one period. *** P ≤ 0.001 , ** P ≤ 0.01 , * P ≤ 0.05

	A	UC	Т	SS
SDM	HH vs. HM	MM vs. MH	HH vs. HM	MM vs. MH
	All Site	s (PO Assess	ment)	
BIOCLIM	0.89***	0.80***	0.88***	0.79***
BRT	0.95***	0.92***	0.96***	0.94***
GAM	0.96***	0.92***	0.97***	0.93***
GLM	0.97***	0.96***	0.96***	0.97***
MAXENT	0.97***	0.94***	0.97***	0.96***
All	0.95***	0.91***	0.95***	0.91***
	Grinnell S	ites (PA Asse	essment)	
BIOCLIM	0.35	0.47**	0.23	0.41*
BRT	0.55**	0.48**	0.57**	0.61***
GAM	0.44*	0.56**	0.57**	0.56**
GLM	0.75***	0.71***	0.57**	0.62***
MAXENT	0.58**	0.48*	0.54**	0.62***
All	0.58**	0.55**	0.53**	0.57**

	BIOCLIM	BRT	GAM	GLM	MAXENT	Grand
Species	HH MM MH HM Mean	Mean				
Callospermophilus lateralis	0.74 0.75 0.72 0.74 0.74	0.82 0.89 0.83 0.83 0.84	0.85 0.88 0.86 0.87 0.87	0.81 0.76 0.72 0.79 0.77	0.86 0.89 0.86 0.87 0.87	0.82
Chaetodipus californicus	0.90 0.90 0.81 0.90 0.88	0.97 0.92 0.96 0.95 0.95	0.97 0.98 0.99 0.99 0.98	0.96 0.99 0.95 0.98 0.97	0.97 0.97 0.96 0.97 0.97	0.95
Dipodomys agilis	0.92 0.88 0.86 0.80 0.86	0.99 0.97 0.98 0.94 0.97	0.97 0.99 0.98 0.97 0.98	0.96 0.98 0.95 0.94 0.96	0.99 0.98 0.98 0.98 0.98 0.98	0.95
Dipodomys heermanni	0.80 0.87 0.64 0.88 0.80	0.98 0.99 0.94 0.95 0.97	0.99 0.99 0.94 0.96 0.97	0.97 0.97 0.98 0.95 0.97	0.99 0.97 0.98 0.98 0.98	0.94
Microtus californicus	0.93 0.90 0.80 0.91 0.89	0.96 0.96 0.95 0.94 0.95	0.96 0.98 0.95 0.94 0.96	0.93 0.83 0.88 0.95 0.90	0.96 0.96 0.96 0.94 0.96	0.93
Microtus longicaudus	0.71 0.73 0.71 0.74 0.72	0.81 0.86 0.81 0.87 0.84	0.85 0.86 0.85 0.88 0.86	0.69 0.77 0.73 0.69 0.72	0.85 0.88 0.84 0.86 0.86	0.80
Microtus montanus	0.59 0.68 0.54 0.61 0.61	0.70 0.78 0.75 0.61 0.71	0.74 0.80 0.69 0.65 0.72	0.54 0.59 0.62 0.51 0.57	0.74 0.84 0.70 0.70 0.75	0.67
Neotoma cinerea	0.61 0.65 0.58 0.64 0.62	0.61 0.81 0.64 0.58 0.66	0.64 0.82 0.66 0.67 0.70	0.65 0.53 0.47 0.57 0.55	0.65 0.80 0.63 0.71 0.70	0.65
Neotoma fuscipes	0.87 0.84 0.64 0.81 0.79	0.98 0.95 0.96 0.93 0.95	0.98 0.93 0.87 0.92 0.92	0.88 0.86 0.92 0.83 0.87	0.95 0.96 0.98 0.94 0.96	0.90
Neotoma macrotis	0.90 0.93 0.83 0.88 0.88	0.98 0.95 0.96 0.97 0.96	0.96 0.98 0.98 0.97 0.97	0.95 0.92 0.93 0.94 0.94	0.95 0.99 0.96 0.97 0.97	0.94
Otospermophilus beecheyi	0.87 0.87 0.80 0.90 0.86	0.89 0.91 0.93 0.91 0.91	0.94 0.95 0.96 0.96 0.95	0.86 0.88 0.90 0.82 0.87	0.93 0.94 0.91 0.90 0.92	0.90
Peromyscus boylii	0.79 0.78 0.80 0.76 0.78	0.87 0.90 0.89 0.84 0.87	0.88 0.89 0.88 0.85 0.87	0.75 0.84 0.85 0.79 0.81	0.85 0.92 0.90 0.84 0.88	0.84
Peromyscus californicus	0.91 0.92 0.77 0.81 0.85	0.98 0.95 0.96 0.99 0.97	0.96 0.99 0.99 0.97 0.98	0.97 0.97 0.96 0.95 0.96	0.97 0.99 0.99 0.98 0.98	0.95
Peromyscus maniculatus	0.57 0.59 0.54 0.60 0.57	0.56 0.67 0.61 0.58 0.60	0.59 0.66 0.62 0.58 0.62	0.40 0.58 0.59 0.38 0.49	0.59 0.64 0.59 0.59 0.61	0.58
Peromyscus truei	0.79 0.78 0.65 0.77 0.75	0.84 0.82 0.78 0.80 0.81	0.84 0.85 0.81 0.79 0.82	0.81 0.79 0.82 0.75 0.79	0.84 0.84 0.79 0.83 0.83	0.80
Reithrodontomys megalotis	0.71 0.71 0.61 0.67 0.67	0.82 0.85 0.84 0.80 0.83	0.79 0.82 0.83 0.82 0.81	0.77 0.81 0.75 0.74 0.77	0.83 0.84 0.83 0.80 0.82	0.78
Sorex monticolus	0.76 0.74 0.62 0.78 0.72	0.84 0.87 0.80 0.84 0.84	0.86 0.88 0.85 0.84 0.85	0.83 0.76 0.73 0.83 0.79	0.86 0.89 0.86 0.88 0.87	0.82
Sorex ornatus	0.88 0.87 0.72 0.80 0.82	0.98 0.95 0.99 0.99 0.98	0.97 0.97 0.99 0.98 0.97	0.97 0.95 0.99 0.97 0.97	0.98 0.96 0.91 0.99 0.96	0.94
Sorex palustris	0.67 0.77 0.63 0.63 0.67	0.71 0.85 0.73 0.77 0.76	0.76 0.88 0.81 0.79 0.81	0.65 0.81 0.74 0.69 0.72	0.79 0.87 0.82 0.86 0.83	0.76
Sorex trowbridgii	0.88 0.86 0.73 0.88 0.84	0.90 0.94 0.91 0.95 0.93	0.94 0.93 0.91 0.94 0.93	0.91 0.94 0.93 0.93 0.93	0.92 0.93 0.94 0.94 0.93	0.91
Sorex vagrans	0.69 0.71 0.68 0.72 0.70	0.76 0.83 0.75 0.81 0.79	0.78 0.86 0.80 0.81 0.81	0.74 0.68 0.73 0.67 0.70	0.83 0.83 0.86 0.86 0.84	0.77
Tamias amoenus	0.79 0.81 0.68 0.78 0.77	0.87 0.85 0.83 0.89 0.86	0.87 0.90 0.80 0.89 0.87	0.65 0.62 0.61 0.72 0.65	0.90 0.89 0.86 0.90 0.89	0.81
Tamias senex	0.85 0.90 0.76 0.88 0.85	0.91 0.99 0.94 0.92 0.94	0.94 0.97 0.90 0.95 0.94	0.92 0.89 0.90 0.85 0.89	0.98 1.00 0.97 0.98 0.98	0.92
Tamias speciosus	0.93 0.94 0.90 0.86 0.90	1.00 0.98 1.00 0.98 0.99	0.97 0.99 0.99 0.98 0.98	0.94 0.86 0.85 0.91 0.89	1.00 0.99 1.00 0.99 1.00	0.95
Urocitellus beldingi	0.82 0.77 0.68 0.71 0.74	0.73 0.79 0.63 0.75 0.72	0.83 0.82 0.81 0.87 0.83	0.64 0.67 0.58 0.59 0.62	0.86 0.87 0.87 0.82 0.85	0.76
Zapus princeps	0.79 0.79 0.76 0.79 0.78	0.81 0.84 0.78 0.86 0.82	0.81 0.86 0.85 0.83 0.84	0.74 0.67 0.59 0.73 0.68	0.83 0.85 0.81 0.87 0.84	0.79
Means	0.79 0.80 0.71 0.78 0.77	0.86 0.89 0.85 0.86 0.86	0.87 0.90 0.87 0.87 0.88	0.80 0.80 0.79 0.80	0.88 0.90 0.88 0.88 0.89	0.84

Supplemental Table 1 AUC measured against known presences and randomly located "pseudoabsences" across the study region.

	Mean	SE
Historic-to-Historic	0.84	0.01
Modern-to-Modern	0.86	0.01
Modern-to-Historic	0.82	0.01
Historic-to-Modern	0.84	0.01

Supplemental Table 2 AUC measured against known presences and true absences at Grinnell re/survey sites. AUC is weakly but negatively related to prevalence in the historic era, meaning that for low-prevalence species it may be upwardly biased, and for high-prevalence species it may be downwardly biased. Species with historic prevalence ≥ 0.80 are denoted with "(\uparrow H)" and those with historic prevalence ≤ 0.20 with "(\downarrow H)." AUC was unrelated to prevalence in the modern era. HH (historic training data/historic test data), MM (modern/modern), MH (modern/historic), HM (historic/modern).

		J	BIOCL	IM				BRT	l.					GAM						GLM					MAXE	NT		Grand
Species	HH	MM	MH	HM	Mean	HH	MM	MH	HM	Mean	Η	H I	MM	MH	HM	Mean	H	H N	MM	MH	HM	Mean	HH	MM	MH	HM	Mean	Mean
Callospermophilus lateralis	0.73	0.66	0.72	0.64	0.69	0.72	0.65	0.80	0.65	0.70	0.	78 (0.68	0.77	0.65	0.72	0.7	79 0	.64	0.79	0.63	0.71	0.78	0.69	0.83	0.67	0.74	0.71
Chaetodipus californicus	0.80	0.84	0.82	0.84	0.83	0.88	0.87	0.88	0.88	0.88	0.	31 (0.83	0.82	0.86	0.83	0.7	74 0	.85	0.76	0.76	0.78	0.76	0.85	0.82	0.85	0.82	0.83
Dipodomys agilis (↓H)	0.79	0.75	0.91	0.70	0.79	0.98	0.84	0.73	0.85	0.85	0.	98 (0.93	0.93	0.90	0.93	0.8	35 0	.92	0.75	0.90	0.85	0.98	0.93	0.95	0.85	0.93	0.87
Dipodomys heermanni (↓H)	0.76	0.94	0.64	0.70	0.76	0.94	0.94	0.95	0.93	0.94	0.) 5 (0.92	0.87	0.92	0.91	0.9	95 0	.92	0.88	0.92	0.92	0.93	0.93	0.92	0.92	0.92	0.89
Microtus californicus	0.79	0.85	0.69	0.91	0.81	0.80	0.77	0.50	0.88	0.74	0.	34 (0.86	0.66	0.85	0.80	0.6	58 C	.79	0.69	0.78	0.74	0.85	0.88	0.69	0.88	0.82	0.78
Microtus longicaudus	0.62	0.64	0.61	0.64	0.63	0.66	0.62	0.66	0.60	0.64	0.	55 C	0.61	0.66	0.59	0.63	0.6	65 0	.61	0.65	0.61	0.63	0.65	0.63	0.65	0.57	0.63	0.63
Microtus montanus	0.74	0.68	0.68	0.70	0.70	0.73	0.72	0.67	0.70	0.70	0.	72 (0.70	0.68	0.70	0.70	0.7	74 0	.69	0.70	0.70	0.71	0.76	0.75	0.69	0.72	0.73	0.71
Neotoma cinerea	0.68	0.92	0.64	1.00	0.81	0.77	1.00	0.78	1.00	0.89	0.	70 1	1.00	0.73	1.00	0.86	0.6	59 1	.00	0.75	1.00	0.86	0.76	1.00	0.74	1.00	0.87	0.86
Neotoma fuscipes (↓H)	0.71	0.59	0.50	0.66	0.61	0.83	0.73	0.59	0.72	0.72	0.	75 (0.87	0.73	0.84	0.80	0.2	24 0	.43	0.44	0.16	0.32	0.88	0.91	0.86	0.85	0.87	0.66
Neotoma macrotis	0.75	0.91	0.80	0.91	0.84	0.80	0.86	0.67	0.87	0.80	0.	77 (0.87	0.67	0.81	0.78	0.5	58 0	.76	0.67	0.67	0.67	0.70	0.92	0.67	0.84	0.78	0.78
Otospermophilus beecheyi (↑H)	0.75	0.71	0.69	0.71	0.71	0.75	0.77	0.38	0.66	0.64	0.	53 (0.69	0.75	0.61	0.67	0.6	59 C	.49	0.31	0.59	0.52	0.75	0.72	0.25	0.67	0.60	0.63
Peromyscus boylii	0.67	0.78	0.68	0.78	0.73	0.65	0.79	0.70	0.75	0.72	0.	58 (0.80	0.65	0.75	0.72	0.6	51 0	.71	0.62	0.68	0.66	0.68	0.81	0.72	0.76	0.74	0.71
Peromyscus californicus (↑H)	NA	0.92	NA	0.92	0.92	NA	0.92	NA	0.95	0.93	N	A (0.95	NA	0.93	0.94	N.	A C	.59	NA	0.56	0.58	NA	0.77	NA	0.88	0.83	0.84
Peromyscus maniculatus	0.50	0.47	0.50	0.46	0.49	0.51	0.54	0.51	0.55	0.53	0.4	49 (0.54	0.51	0.54	0.52	0.4	48 0	.54	0.49	0.51	0.51	0.49	0.53	0.49	0.53	0.51	0.51
Peromyscus truei	0.67	0.78	0.69	0.77	0.73	0.70	0.78	0.69	0.64	0.70	0.	79 (0.74	0.77	0.70	0.75	0.7	71 0	.55	0.64	0.59	0.62	0.78	0.81	0.73	0.67	0.75	0.71
Reithrodontomys megalotis	0.70	0.78	0.71	0.75	0.74	0.80	0.82	0.79	0.84	0.81	0.	31 (0.89	0.83	0.83	0.84	0.8	30 0	.77	0.77	0.79	0.78	0.82	0.91	0.81	0.83	0.84	0.80
Sorex monticolus	0.71	0.61	0.61	0.61	0.63	0.79	0.57	0.73	0.57	0.66	0.	77 (0.54	0.74	0.59	0.66	0.7	77 0	.59	0.80	0.55	0.68	0.79	0.58	0.77	0.59	0.68	0.66
Sorex ornatus	0.61	1.00	0.58	0.96	0.79	0.64	0.79	0.61	0.95	0.75	0.	52 (0.90	0.73	1.00	0.79	0.7	77 1	.00	0.77	1.00	0.88	0.73	0.97	0.67	1.00	0.84	0.81
Sorex palustris	0.75	0.61	0.66	0.59	0.65	0.69	0.80	0.75	0.70	0.73	0.	72 (0.55	0.70	0.55	0.63	0.6	62 0	.66	0.62	0.66	0.64	0.72	0.73	0.74	0.70	0.72	0.68
Sorex trowbridgii	0.89	NA	0.81	NA	0.85	0.93	NA	0.83	NA	0.88	0.	33	NA	0.71	NA	0.77	0.7	79	NA	0.82	NA	0.80	0.86	NA	0.83	NA	0.85	0.83
Sorex vagrans (↓H)	0.76	0.85	0.85	0.76	0.81	0.94	0.80	0.78	0.88	0.85	0.	96 (0.72	0.63	0.91	0.81	0.6	59 C	.52	0.59	0.60	0.60	0.96	0.88	0.92	0.89	0.91	0.80
Tamias amoenus (↓H)	0.87	0.87	0.82	0.87	0.86	0.92	0.92	0.87	0.91	0.90	0.	92 (0.93	0.92	0.91	0.92	0.8	30 0	.67	0.66	0.74	0.72	0.92	0.91	0.91	0.91	0.91	0.86
Tamias senex	0.86	0.84	0.74	0.79	0.81	0.75	0.86	0.78	0.85	0.81	0.	38 (0.87	0.77	0.87	0.85	0.7	76 0	.85	0.74	0.84	0.80	0.89	0.88	0.83	0.88	0.87	0.83
Tamias speciosus	0.63	0.70	0.59	0.67	0.65	0.68	0.74	0.62	0.72	0.69	0.	70 (0.77	0.68	0.75	0.72	0.7	70 0	.70	0.68	0.77	0.71	0.69	0.78	0.66	0.73	0.71	0.70
Urocitellus beldingi (↓H)	0.68	0.71	0.55	0.60	0.64	0.84	0.79	0.89	0.71	0.80	0.	36 (0.84	0.89	0.79	0.85	0.8	33 0	.75	0.82	0.75	0.79	0.87	0.85	0.85	0.75	0.83	0.78
Zapus princeps	0.69	0.60	0.61	0.64	0.64	0.69	0.69	0.62	0.68	0.67	0.	70 (0.66	0.66	0.71	0.68	0.6	66 0	.69	0.65	0.68	0.67	0.75	0.68	0.69	0.69	0.70	0.67
Means	0.73	0.76	0.68	0.74	0.73	0.78	0.78	0.71	0.78	0.77	0.	77 0).79	0.74	0.78	0.77	0.7	70 0	.71	0.68	0.70	0.70	0.79	0.81	0.75	0.79	0.79	0.75

	Mean	SE
Historic-to-Historic	0.75	0.01
Modern-to-Modern	0.77	0.01
Modern-to-Historic	0.71	0.01
Historic-to-Modern	0.76	0.01

Models for Callospermophilus lateralis



Models for Chaetodipus californicus



Models for Dipodomys agilis



Models for Dipodomys heermanni



Models for Microtus californicus



Models for Microtus longicaudus



Models for Microtus montanus





Models for Neotoma cinerea





Models for Neotoma fuscipes



Models for Neotoma macrotis

Models for Otospermophilus beecheyi

Models for Peromyscus boylii

Models for Peromyscus californicus

Models for *Peromyscus maniculatus*

Models for Peromyscus truei

Models for Reithrodontomys megalotis

Models for Sorex monticolus

Models for Sorex ornatus

Models for Sorex palustris

Models for Sorex trowbridgii

Models for Sorex vagrans

Models for Tamias amoenus

Models for Tamias senex

Models for Tamias speciosus

Models for Tamiasciurus douglasii

Models for Urocitellus beldingi

Models for Zapus princeps

