# Comparing generalized and customized spread models for nonnative forest pests

Emma J. Hudgins,[1,5] Andrew M. Liebhold,[2,3] and Brian Leung[1,4]

[1]*Biology Department, McGill University, Montreal, Quebec H3A 1B1 Canada*
[2]*Northern Research Station, USDA Forest Service, Morgantown, West Virginia 26505 USA*
[3]*Faculty of Forestry and Wood Sciences, Czech University of Life Sciences Prague, Praha 6 – Suchdol, Czech Republic*
[4]*School of Environment, McGill University, Montreal, Quebec H3A 2A7 Canada*

*Abstract.* While generality is often desirable in ecology, customized models for individual species are thought to be more predictive by accounting for context specificity. However, fully customized models require more information for focal species. We focus on pest spread and ask: How much does predictive power differ between generalized and customized models? Further, we examine whether an intermediate "semi-generalized" model, combining elements of a general model with species-specific modifications, could yield predictive advantages. We compared predictive power of a generalized model applied to all forest pest species (the generalized dispersal kernel or GDK) to customized spread models for three invasive forest pests (beech bark disease [*Cryptococcus fagisuga*], gypsy moth [*Lymantria dispar*], and hemlock woolly adelgid [*Adelges tsugae*]), for which time-series data exist. We generated semi-generalized dispersal kernel models (SDK) through GDK correction factors based on additional species-specific information. We found that customized models were more predictive than the GDK by an average of 17% for the three species examined, although the GDK still had strong predictive ability (57% spatial variation explained). However, by combining the GDK with simple corrections into the SDK model, we attained a mean of 91% of the spatial variation explained, compared to 74% for the customized models. This is, to our knowledge, the first comparison of general and species-specific ecological spread models' predictive abilities. Our strong predictive results suggest that general models can be effectively synthesized with context-specific information for single species to respond quickly to invasions. We provided SDK forecasts to 2030 for all 63 United States pests in our data set.

*Key words: exotic; macroecology; multispecies; nonindigenous; risk assessment; simulation; spatially explicit.*

## Introduction

Identification of generalized models that explain and predict species distributions is of fundamental importance to ecologists. However, while general models "potentially inform about phenomena that exist in many systems...", they "...may not necessarily make good predictions about any individual system" (Evans et al. 2013). This tension between generality and context specificity underlies much of ecology.

The trade-off between generality and ecological prediction also exists within invasion biology, where the focus has been on species-specific models using context-specific information (Liebhold et al. 1992, Gilbert et al. 2004). For instance, the spread phase, a fundamental part of the invasion process, has typically relied on customized models, accounting for life history, ecological complexity and spatial factors such as dominant wind direction and habitat suitability (Liebhold et al. 1992, Koch and Smith 2008, Kovacs et al. 2011). Intuitively, models that are based on a particular invasive species' local context should provide better predictions than general models and should facilitate management. For instance, the gypsy moth Slow-the-Spread (STS) project in the United States has reduced spread rates by >70%, since its inception in 2000, (Sharov et al. 2002, Grayson and Johnson 2018; see also Slow Ash Mortality [SLAM] program; McCullough and Mercader 2012).

While customized models have undeniably been useful, there have been calls for pathway-level analyses, which account for multiple invasive species simultaneously (Lodge et al. 2006, Bradie and Leung 2015). For invasive species, one phenomenon that supports such cross-species generality is the dominant role of humans in transporting species via mechanisms that are analogous across entire suites of species invading different spatial locations (e.g.,

through ballast water transport [Seebens et al. 2015], firewood transport [Haack 2006]). We hypothesized that, across invasions, unique natural dispersal processes are commonly overridden by anthropogenic ones, and that predictable generalities that operate across entire suites of species arise as a consequence of these processes' broad effects (Hudgins et al. 2017). In the case of post-establishment spread, anthropogenic mechanisms such as transport through road networks may mean that conventional ecological processes governing dispersal, which are more idiosyncratic across species (wind direction, flight ability, etc.; Aylor 1990, Taylor et al. 2010), are less important for forecasting spread at large scales. Further, species traits relating to association with anthropogenic dispersal vectors may be most important in determining their spread rates. This anthropogenic dominance can thus provide us with general macroecological predictions for the spread of groups of invasive species within a particular transport network.

Although general invasion models are in the minority, the interest in multispecies predictions for the spread of invasive species arises because of their potential advantages. In order to prepare for and limit invasive species impacts across space, which accrue immense costs (Vilà and Hulme 2017), managers need to know where these species will invade next. Further, the sooner they can take action, the more effective their control measures will be (Lovett et al. 2016). The lower the data requirements of a given model, the sooner it can be implemented to inform management. As such, a highly general model could be rapidly applied to many species, potentially including species that have not yet established. Thus, in summary, there are potential benefits from using a general model and logical reasons to expect generality in the spread of invasive species.

Applied ecological models can be viewed along a continuum from specific to general. At the specific end of the spectrum, structure, predictors and parameters may all be fit to each separate species (i.e., customized models). At the most general end, a model may be applied to many species, using the same model structure, predictive factors, and parameters. In the middle of the spectrum, parameters can be added or rescaled to different values within a generalized model "backbone" in order to incorporate additional layers of customization (we term these "semi-generalized models"), without the collection of as much species-specific data (e.g., time series for each species). These intermediate models can be worthwhile to consider, if the reduction in generality is offset by a large gain in predictive ability. Additionally, semi-generalized models that do not rely only on single-species data could conceivably make better predictions relative to customized models, if there are strong commonalities in the spread process across species (e.g., human-mediated vectors), since they are able to "borrow" information from a broad pool of species.

In this paper, we compare a suite of models with varying levels of generality in terms of their ability to predict the spread of invasive forest pests. For context specificity, we designed customized single-species models for three pest species for which time-series data exist (beech bark disease (*Cryptococcus fagisuga*), gypsy moth (*Lymantria dispar*), and hemlock woolly adelgid (*Adelges tsugae*), using species-specific predictors and functional forms (Liebhold et al. 1992, Morin et al. 2007, 2009). These were compared against a general model fit across all forest pest species currently known in the United States, using a previously published generalized dispersal kernel (GDK) (Hudgins et al. 2017). At the intermediate level, we examined whether we could use GDK as a starting point, and incorporate species-specific knowledge (semi-generalized models, SDK), and whether doing so improved predictions compared to GDK and customized models.

## MATERIALS AND METHODS

### Dispersal kernel formulation

Dispersal kernels estimate the probability of pest dispersal across space based on the distance, $d$, between source and destination locations (Kot et al. 1996). In the GDK, we moderated dispersal though spatial predictors affecting the dispersal kernel. We fit our model using discrete time simulations, where at each time step, pests dispersed to surrounding patches according to

$$T_{i,j} = \frac{e^{-d_{i,j}, f(Z)}}{\sum_j e^{-d_{i,j}, f(Z)}} \tag{1}$$

$$f(Z) = 2\alpha \frac{e^{(Z_S+Z_I+Z_O)}}{(1 + e^{(Z_S+Z_I+Z_O)})} \tag{2}$$

where $T_{i,j}$ is the proportion of pests dispersing from cell $i$ to cell $j$, normalized such that the value of the dispersal kernel across all locations $j$ sums to 1 (denominator of Eq. 1), $d_{i,j}$ is the distance between sites $i$ and $j$, and $f(Z)$ is a combination of all fitted species ($Z_S$) and cell (dispersal into a cell = $Z_I$, dispersal out of a cell = $Z_O$) specific predictors influencing the dispersal probabilities, scaled to have a mean value of $\alpha$ (i.e., dispersal occurs at rate $\alpha$ for sites with predictor variables at their mean levels). For the GDK, our distributional data were limited to each species' final distribution at the end of the fitting period, plus data on their reported first year of establishment in the United States.

The GDK is made up of both a dispersal and a growth component, where the relative propagule pressure in cell $i$ at time step $t + 1$ is equal to the relative propagule pressure at time $t$, minus emigration to all cells $j$, plus immigration from all cells $k$, multiplied by the growth rate $\delta$:

$$X_{i,t+1} = (X_{i,t} - \sum_j T_{i,j} X_{i,t} + \sum_k T_{k,i} X_{k,t}) \delta. \qquad (3)$$

Cells are considered "presences" capable of being a source of propagules when they are above a threshold population size $\phi$ with a maximum relative propagule pressure in a cell of 1. We assumed that cells that were invaded remained invaded.

For the GDK, we considered predictors, including propagule pressure proxies, habitat invasibility proxies and pest life history traits (sources fully described in Hudgins et al. [2017]). The best-fitting model retained four predictor variables, wherein sites with greater forested land area and human population density are attractors to invasive pests, and sites with greater tree density and human population density are sinks from which pests do not disperse as much, relative to sites with lower values of these predictors. We modeled 5-yr time steps, to achieve finer-scale forecasting (the original model used 10-yr steps, but was shown to be robust; Hudgins et al. 2017).

### Allowing context specificity

We designed customized dispersal models for each of three highly damaging invasive forest pests: beech bark disease (BBD), gypsy moth (GM), and hemlock woolly adelgid (HWA; Fig. 1). BBD is a disease complex made up of the introduced beech scale insect *C. fagisuga* and (most likely native) fungi (one of two species of *Neonectria*) first detected in Halifax, Nova Scotia in 1890, with potential additional introductions around Boston and New York City (Houston 1994). GM is a highly poly-phagous (i.e., having many host tree species) defoliator introduced from France to Medford, Massachusetts around 1869 (Liebhold et al. 1989). HWA is a sap-feeding insect that was first detected in 1951 in Richmond, Virginia (Ward et al. 2004). These species span three of the four feeding guilds of the broader set of 63 species used to fit the GDK (Appendix S1; included pathogens, foliage feeders, sap feeders; missing borers; Hudgins et al. 2017). While time-series data exist for emerald ash borer *Agrilus planipennis*, its detection records begin in 2002, which was after our fitting year (2000).

Across the customized models, we tested the inclusion of four additional levels of complexity compared to the GDK: testing additional dispersal kernel shapes, pest entry points, additional species-specific predictor variables, and time series of pest spread.

First, in addition to the negative exponential dispersal kernel employed in the GDK, a leptokurtic kernel was explored (sensu Kot et al. 1996):

$$T_{i,j} = \frac{e^{-\sqrt{d_{i,j}} f(Z)}}{\sum_j e^{-\sqrt{d_{i,j}} f(Z)}}. \qquad (4)$$

Leptokurtic dispersal kernels allow for nonlinear spread rates and increased dispersal to distant locations

(Shigesada et al. 1995, Kot et al. 1996). Spatial predictor variables were analogously incorporated via Eq. 2, but the leptokurtic kernel has more density in its tails and therefore leads to a higher chance of long-distance dispersal. The dispersal kernel that resulted in the best model fit was selected for each species separately.

Second, we simulated the best-known starting location of each pest species (Ward et al. 2019) and the host centroid as a starting point for each pest's dispersal. While the use of best-known starting points did not improve the overall fit of the GDK (Hudgins et al. 2017), given that these three species are some of the most well studied, these starting points are likely more reliable than for other pests. If a starting location was not within our known host range for a given pest (e.g., first detection in an urban area), we chose the closest grid cell in the host range. As with the dispersal kernel, the starting point that resulted in the best model fit was chosen for each species.

Third, we tested additional predictors mined from the literature in our forward selection models. We tested firewood and campground-related variables, which were frequently included in spread models of gypsy moth (Bigsby et al. 2011). We sourced these predictors from the U.S. Census' American Housing Survey (homes fueled by wood, campground density, seasonal homes), and tested for all three pests. Additionally, HWA is known to be highly climatically limited, with high mortality when exposed to low winter temperatures (Paradis et al. 2008, Morin et al. 2009). We modeled climatic limitation for HWA using minimum temperature of the coldest month (bio6) from WorldClim (Fick and Hijmans 2017 ) (Appendix S2), and setting the pest density to zero for any patch below a fitted threshold (climate data *available online*).[6] Any predictor that substantially improved fit was included in the final customized model for a given species.

Finally, the customized models were each fit to time-series of species dispersal patterns, using historical discovery records by county available for the above three species, while the GDK was constructed using only the final distributions (but many more species).

### Semi-generalized models

For the SDK, we tested the inclusion of three additional layers of species-specific information that went beyond the original GDK, but did not use time-series information (in contrast to the customized models), as these data are relatively rare. First, we utilized an "intercept correction" to offset each single-species spread trajectory such that it minimized fitted GDK residuals. Second, we tested whether incorporating information on the best-known initial invasion location improved predictive ability for each pest. Third, we tested whether
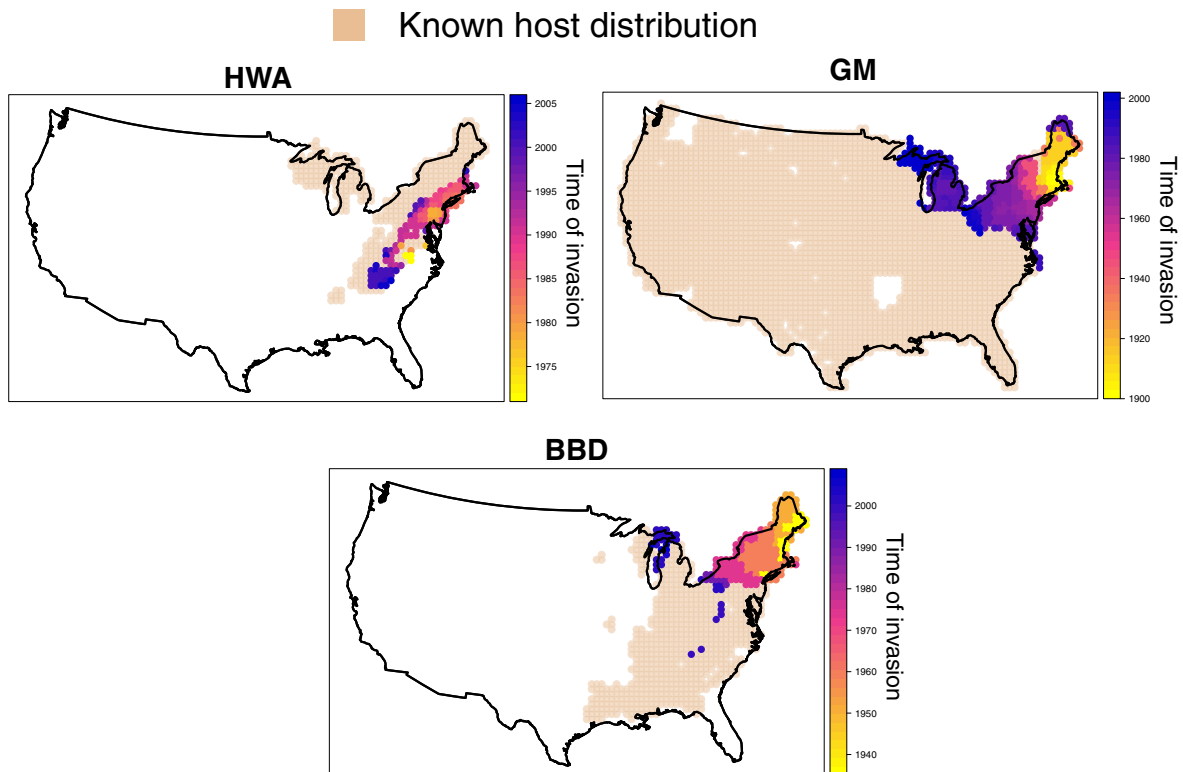
---

[6] www.worldclim.org.

Fig. 1. Historical spread patterns of the three focal species: gypsy moth (GM), hemlock woolly adelgid (HWA), and beech bark disease (BBD). Older invasions are shown in yellow while more contemporary invasions are shown in blue. Known host distribution is shown in beige.

incorporating information on a species' known niche limitations could improve our forecasts.

*GDK intercept correction.*—The earliest invasive spread model is Skellam's seminal work on reaction-diffusion (RD, Skellam 1951). RD uses the physics of diffusion to predict radial spread emanating from a single source, where the size of the invaded range expands uniformly over time (Shigesada et al. 1995), and is a core model in theoretical ecology (Morin et al. 2007, 2009, Skuhravá et al. 2007).

Theoretically, we recognized that the GDK spread intercept ($\alpha$) is related to the RD diffusion coefficient ($D$), as it similarly acts to scale the relationship between dispersal distance and probability (Appendix S3), and hypothesized that, holding all spatial predictors constant, this intercept could be rescaled to adjust species spread, thereby improving forecasts. Using only the fitted values of the GDK to forecast spread neglects additional information contained in the mismatch between these fitted values and the observed distribution in the fitting year. If GDK residuals reflect constant, unmeasured species-specific factors (such as probability of uptake by humans, fertility rates, etc., that are difficult to obtain for an entire community of species), accounting for these deviations in spread trajectories could

improve predictions. We thus refit the spread intercept for each species, but otherwise maintained the proportional relationship with other spatial predictors, and general structure of the GDK previously fit (i.e., Eqs. 1, 2; with dispersal based on forested land area, human population, and tree density as in Hudgins et al. 2017). The data requirements for this adjustment are simply the locations at one time-point in a species' distribution, something that is presently available for all known damaging invasive pests.

*GDK starting point correction.*—Second, as one of the simplest correction factors, we tested whether adding information on our best estimate of a pest's initial invasion site within the United States improved the predictions of our general model, when there is reasonable confidence in those estimates. We note that such estimates may not improve the predictions of all species if the starting point is not well known, but as we have mentioned, all three of these species are well studied. For this correction, we updated our GDK simulation to begin spread from these sites rather than the host range centroid.

*Niche limitation correction.*—Third, we tested the inclusion of species' niche constraints. While the first two

correction factors require very little context-specific information, which is likely to be known for any future invader, niche limitations are more idiosyncratic species information. For one of our studied species, HWA, it is clear that climatic limitation plays a strong role in limiting its northward dispersal (Morin et al. 2009). Just like in the customized models, we tested the addition of a minimum temperature threshold for HWA persistence (Appendix S2).

We chose the SDK corrections for each species that maximized fit. Importantly, all three customizations were added to the basic structure of the GDK, holding all other fitted parameter values constant for all predictor variables in the published model.

### Data preparation

For the customized models, historical county-level spread records were assembled for GM, HWA, and BBD. Records of historical GM spread were obtained from the United States Code of Regulations (Title 7, Chapter III, Section 301.45) which annually designates quarantined counties that are part of the "generally infested area." Federal quarantines do not exist in the United States for BBD and HWA, however, similar county-level records were obtained from other sources (Morin et al. 2007, 2009). County records were overlaid on a 50 × 50 km grid in order to control for county size, where a detection anywhere in a grid cell was considered a valid presence. Five-year time steps within historical spread data sets with less than two new detections were not considered in our spread models, because there are inherent delays between the detection of pests in surveys and the incorporation of that information into range databases. These years likely correspond to times where monitoring was not adequately performed (resulting in minimal apparent spread even if new invasions were occurring), and are not a good indicator of the spread trajectory. Once these low-detection years were removed, the number of independent fitting years for each species was 15 for GM, 6 for BBD and 5 for HWA. All three historical spread data sets included data beginning only at the first date of multi-county range for each pest, but dates of initial discovery/introduction are known for each pest, so simulated spread was adjusted to include the period between initial discovery/introduction to the first record of multi-county spread. HWA was adjusted from 1971–2005 to 1950–2005, GM was adjusted from 1902–2005 to 1865–2005, and BBD was adjusted from 1935–2005 to 1890–2005.

### Customized model fitting

To maintain consistency with the GDK, a forward selection procedure based on the same metric (the minimum energy test, MET; Aslan and Zech 2005) and using the same threshold for parameter inclusion (5 km) as in Hudgins et al. (2017) was employed to build the customized models. MET accounts not only for exact spatial matches of predicted and observed presences (similar to measures such as accuracy), but also apportions better scores to "close" matches than presences predicted very far away from the observed presences (for a further discussion of MET, see Hudgins et al. [2017]). Rather than taking the average across 63 species, in the customized models, this 5-km MET threshold was applied on average across all fitting years for a single species. We chose the best single-species forward selection model among the two dispersal kernel shapes and two possible starting locations for each pest species. In the case of HWA, the temperature threshold was applied to all four possible customized models, to remain consistent with the literature on niche limitation and to ensure the methodology was comparable to the fitting of the SDK.

### Predictive validation metric

To explore predictive ability with greater ease of interpretation, we derived a novel, simple pseudo $R^2$ value, based on optmatch, an algorithm originally used to match treatment to control subjects in clinical trials (R package optmatch; Hansen 2007). The optmatch tool uses a global optimization approach to match two sets of points, minimizing the total multivariate distance between the sets. We wished to have a metric that takes its maximum value when two distributions have the same number of points, with the points in the correct spatial locations.

We first used optmatch to perform a one-to-one match between our predicted and observed presence points for a given pest. Next, the leftover points caused by differences in predicted and observed range size were then used to penalize the distance score. To do this, we assigned these leftover points the mean distance between that point and all other points in the opposing distribution. We used the mean of this entire vector of distances (optimal matching mean squared error, omMSE) and converted it to a pseudo $R^2$ ($R_{om}^2$) by comparing the observed omMSE value to a spatial null expectation, using 10,000 random points from the host distribution (Appendix S4).

### Community forecast

We used the best-fitting SDK to forecast the distribution of all 63 pest species from 2005 to 2030. We used the fitted MET score applied to each individual species' snapshot of dispersal in 2005 in order to determine the SDK layers to include for each species (intercept, starting point, and niche limitation corrections where there was evidence from the literature that they were necessary, see Appendix S5). We reset pest distributions to known distributions at 2005 (setting false absences to $\phi$, false presences to 0, and maintaining the simulated propagule pressure of true presence sites) before

simulating spread using each species' SDK parameters to 2030. We included projected human population estimates from ProximityOne as an updated human population predictor in the GDK-based models, which all included this term (population estimates *available online*).[7]

To model uncertainty, we considered future climate and human population size projections, and uncertainty in fitted model parameters (see Appendix S6 for full details). In brief, for climate change, we used rcp2.6 and rcp8.5 climate scenarios from BIOCLIM, and for population size, we used two scenarios based on the Shared Socioeconomic Pathways: SSP3 or "Regional Rivalry," and SSP5 or "Fossil-Fueled Development" (Hauer 2019), representing the extremes for both climate and population size, respectively. We note that there was only evidence of climate limitations for two species, but we nonetheless considered climate scenarios for completeness (Appendix S5). For the parameter uncertainty in the SDK model, we conducted sensitivity analysis, randomly perturbing model parameters, and using the threshold for parameter inclusion in our model fitting process as our criterion to retain parameter sets (i.e., MET within 5 km of the best-fitting model). We examined the combined effect of uncertainty on the range of predicted future pest richness.

## RESULTS

### Customized model selection and predictive validation

Customized models were highly predictive on average ($R^2_{om} = 0.74$), though predictions were weakest for HWA ($R^2_{om} = 0.45$). The best-fitting customized model had a very simple functional form for each of the three species (Table 1), with fewer predictors than the GDK. For GM, only per capita income was important, showing a negative effect on spread into sites. For BBD, the best model included only an intercept term. In the HWA model, which contained the minimum temperature threshold, human population density displayed the same relationship as it did in the GDK, increasing spread into sites. In all cases, using the hypothesized initial introduction location as a starting point led to better fits than using the centroid of the host range. For BBD, the leptokurtic dispersal kernel fit better than the negative

exponential kernel, while the negative exponential outperformed the leptokurtic model for HWA and GM.

### GDK predictive validation

The strength of the uncorrected model's predictions varied across the three species, from being extremely predictive for GM ($R^2_{om} = 0.87$), to highly predictive for BBD ($R^2_{om} = 0.55$), to more moderately predictive for HWA ($R^2_{om} = 0.30$). The uncorrected GDK overestimated spread for these three species, but predictions were still substantially better than random expectations from our null model, and mean spatial variation explained was $R^2_{om} = 0.57$.

### SDK model selection and predictive validation

The $R^2_{om}$ improvement ranged from 0.11 to 0.55 between the uncorrected GDK and the best SDK, and from 0.03 to 0.40 between the customized model and the best SDK for the validation year (mean SDK $R^2_{om} = 0.91$). The best SDK for BBD and GM included the intercept and starting point corrections, and had $R^2_{om} = 0.89$ and $R^2_{om} = 0.98$, respectively (Table 2). HWA required a third level of complexity, where the model with the starting point, intercept and niche limitation corrections resulted in the best fit and had $R^2_{om} = 0.85$. For GM and BBD, corrected intercepts were larger in magnitude than the uncorrected GDK intercept, consistent with a reduction in spread extent. Conversely, for HWA, SDK had a smaller intercept, indicating a higher spread rate. However, this spread rate was offset by pest mortality upon dispersal into the northernmost parts of its range, thereby leading to a lower extent of spread overall.

### Model comparison: spatial details

Although the GDK retained moderate to high predictive power, and performed similarly to the customized model for GM, it was weaker than the customized models for BBD and HWA. For GM, both the customized model and the SDK explained over 92% of spatial variation in pest distributions. However, while the customized model performed well in terms of the $R^2_{om}$, based on visual

TABLE 1. The best-fitting single-species models for hemlock woolly adelgid (HWA), gypsy moth (GM), and beech bark disease (BBD).

| Species | Kernel | ϕ | δ | α | Predictor(s) | MET (km) | $R^2_{om}$ |
|---|---|---|---|---|---|---|---|
| HWA | negative exponential | 0.0020 | 2.48 | 2.93 | human population (+ in) = 0.19; bio6 = −9.33°C | 11.10 | 0.45 |
| GM | leptokurtic | 0.0046 | 4.60 | 5.14 | income (− in) = 1.19 | 9.82 | 0.92 |
| BBD | negative exponential | 0.0008 | 1.27 | 1.37 | NA | 11.37 | 0.86 |

*Notes:* Predictor variables labelled "in" represent predictors of dispersal into sites. In all cases, the best model simulated spread initiating at the most likely initial invasion of the pest rather than the centroid of the host range. Bio6 represents a fitted minimum temperature threshold for HWA mortality, and NA indicates no additional predictors were included.

[7] www.proximityone.com

TABLE 2. The results of the GDK validation for both uncorrected and intercept-corrected models.

| Species | α | MET (km) | $R^2_{om}$ |
|---|---|---|---|
| Uncorrected | | | |
| HWA | 1.74 | 33.22 | 0.30 |
| GM | 1.74 | 41.26 | 0.87 |
| BBD | 1.74 | 110.01 | 0.55 |
| Intercept-corrected | | | |
| HWA | 1.65 | 31.17 | 0.20 |
| GM | 2.25 | 8.94 | 0.98 |
| BBD | 3.95 | 10.13 | 0.87 |
| Starting-point-corrected | | | |
| HWA | 1.74 | 32.14 | 0.73 |
| GM | 1.74 | 18.16 | 0.92 |
| BBD | 1.74 | 1.70 | 0.89 |
| Starting-point, intercept-corrected | | | |
| HWA | 1.48 | 24.84 | 0.35 |
| GM | 1.61 | 2.55 | 0.98 |
| BBD | 1.78 | 1.35 | 0.89 |
| Starting-point, intercept, climate-corrected | | | |
| HWA | 1.47; bio6 = −7.64°C | 1.75 | 0.85 |

*Notes:* All GDK models have parameter values of forested land area (+ in) = 0.53, tree density (− out) = 15.61, human population density (+ in) = 0.16, human population density (− out) = 0.32, threshold population size φ = 0.00054, local population growth rate δ = 1.30. HWA, hemlock woolly adelgid; GM, gypsy moth; BBD, beech bark disease; GDK, generalized dispersal kernel; MET, minimum energy test. Bio6 represents a fitted minimum temperature threshold for HWA mortality.

inspection, it produced spatial patterning incongruent with the true pest distribution, likely due to its fitted relationship with income (Fig. 2, top left panel). The model's leptokurtic kernel and negative relationship with income produce a discrete patch of invaded sites around South Dakota and Nebraska. GM's extensive host range could allow incongruent distributions to have high $R^2_{om}$ if predicted distributions are of approximately the correct range size and close geographically to observed distributions. In contrast, SDK predicted a distribution that overlapped with the observed distribution nearly entirely. Both the SDK and the customized models performed well for BBD, leading to a tight spatial match between its predicted and observed distribution. The HWA customized model explained the lowest amount of spatial variation of the three species ($R^2_{om}$ = 0.45), likely due to the increased complexity of this species' spread mechanisms (i.e., climatic limitation), leading to an inability to capture the southernmost part of the range without overpredicting to the north. The customized model's fitted temperature threshold was quite low (~2°C lower than in the SDK), resulting in only a small effect on restricting pest distributions. The GM uncorrected GDK overpredicted spread, producing a distribution that included much of its invaded range, but lacking climatic limitation in northern areas. In contrast, the SDK no longer overpredicted spread in the north, and also did not predict disjointed jumps outside the observed distribution that the

customized model predicted ($R^2_{om}$ = 85%), but still did not capture the southernmost distribution.

*Forecasts*

Using the optimal set of SDK layers for each of the 63 species (see Appendix S5 for details of SDK corrections), our simulations project the distribution of pests at 2030 to remain highly aggregated in the northeastern United States, as it was in 2005, but pest species richness to increase (Fig. 3a, b). Northern Minnesota and Wisconsin, western Montana and northern Idaho, parts of New Mexico, and northern New England are predicted to have the largest increase in local establishments by 2030 (Fig. 4, regions B, C, E, F, K). Some smaller, more concentrated areas of increase are also predicted (Fig. 4, regions A, D, G-J). We predict very few new local establishments in the middle of the country. The areas at high risk correspond to high forested land and increasing human population densities. New local establishments are especially high in urban centers close to regions of high forested land area (Fig. 4, dashed lines). However, some less populated areas also see large increases in local establishments due to their high amount of forested land (Fig. 4, solid lines).

The combination of uncertainties in climate change, future population growth, and model parameters led to strong regional variability across future pest richness predictions (Fig. 5). However, the northeastern United States typically had the greatest number of relative establishments, indicating a consistent pattern of future spatial risk despite uncertainty. The simulations that produced the fewest novel local establishments were those from the high human population growth scenario, given decreased dispersal out of high population density sites. Scenarios with the highest future spread had increased pest growth rates and less preferential dispersal into high population areas, leading to more even dispersal patterns across space (see Appendix S6 for further discussion). The median range of predicted pest load was five species, but distinct regional differences were observed. The central portion of the United States had the lowest uncertainty (range of <5 pests), but was consistently predicted to have low numbers of future local establishments. The western United States had more moderate levels (~5–10 pests), while the eastern US had the highest levels (~10–20 pests). Some future hotspots were particularly variable (Regions C, D, H in Fig. 4). Additionally, many of the high uncertainty patches, which are particularly dense across the eastern US, had not been identified as hotspots (Fig. 4), indicating some potential for additional regions of high future pest load that warrant managerial surveillance.

DISCUSSION

*The performance of general versus single-species models*

The customized models performed 17% better than the GDK, averaged across our three case studies, with
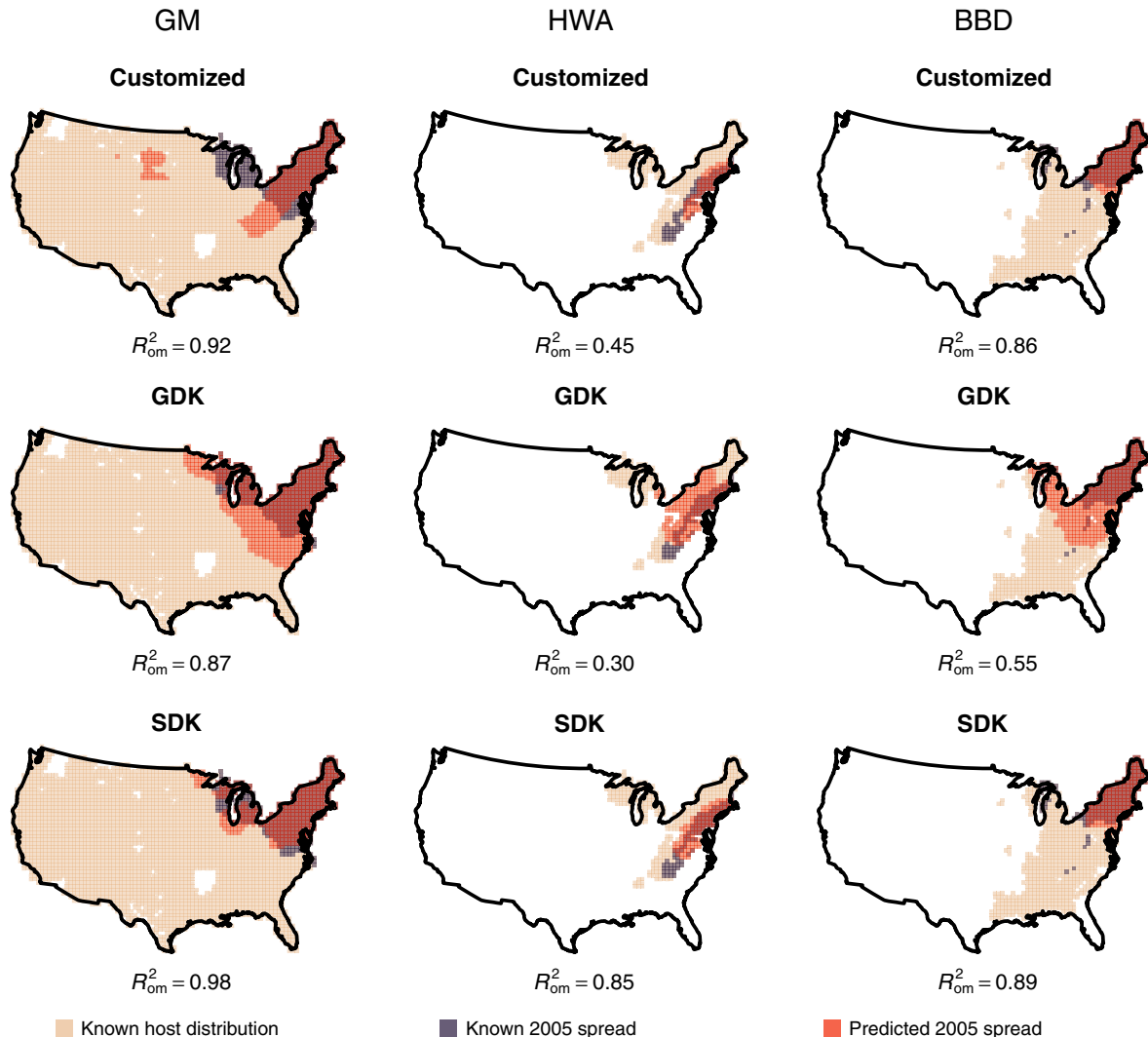
Fig. 2. A comparison of the predictive ability of customized (top row), uncorrected generalized dispersal kernel (GDK, middle row), and semi-generalized dispersal kernel (SDK) models for gypsy moth (GM, left column), hemlock woolly adelgid (HWA, middle column), and beech bark disease (BBD, right column). Host presence is indicated in beige, predicted distributions after a forecast (5-yr) are shown in red, and observed distributions are shown in blue. Areas of overlap between predicted and observed distributions produce a darker red color due to the overlap of the red and blue colors.

74% spatial variation explained. While customization showed a marked benefit, notably, the GDK was still able to capture a respectable 57% of spatial variation in spread, for these three species. Moreover, GDK's predictive ability may be substantially higher for most other species: We note that HWA and BBD were much more poorly fit by the uncorrected GDK than the majority of species, while GM was fit better than average, and that the magnitude of over or underprediction in the uncorrected GDK appears to predict forecasting ability (Appendix S7). Thus, we expect the GDK's average predictive ability to be between BBD (55%) and GM (87%). While researchers have reasonably focused on customized single-species models for prediction (Liebhold

et al. 1992, Koch and Smith 2008, Kovacs et al. 2011), the GDK yielded useful predictions even without any modification, and will be useful in situations where customized models cannot be built, e.g., in the case of novel invaders.

### Comparing predictors in GDK versus customized models

We found different suites of predictors to be important for single-species spread, and that fewer predictors were important compared to the general predictors for all species in the GDK.

The differences in predictors between the GDK and customized models could have arisen due to two
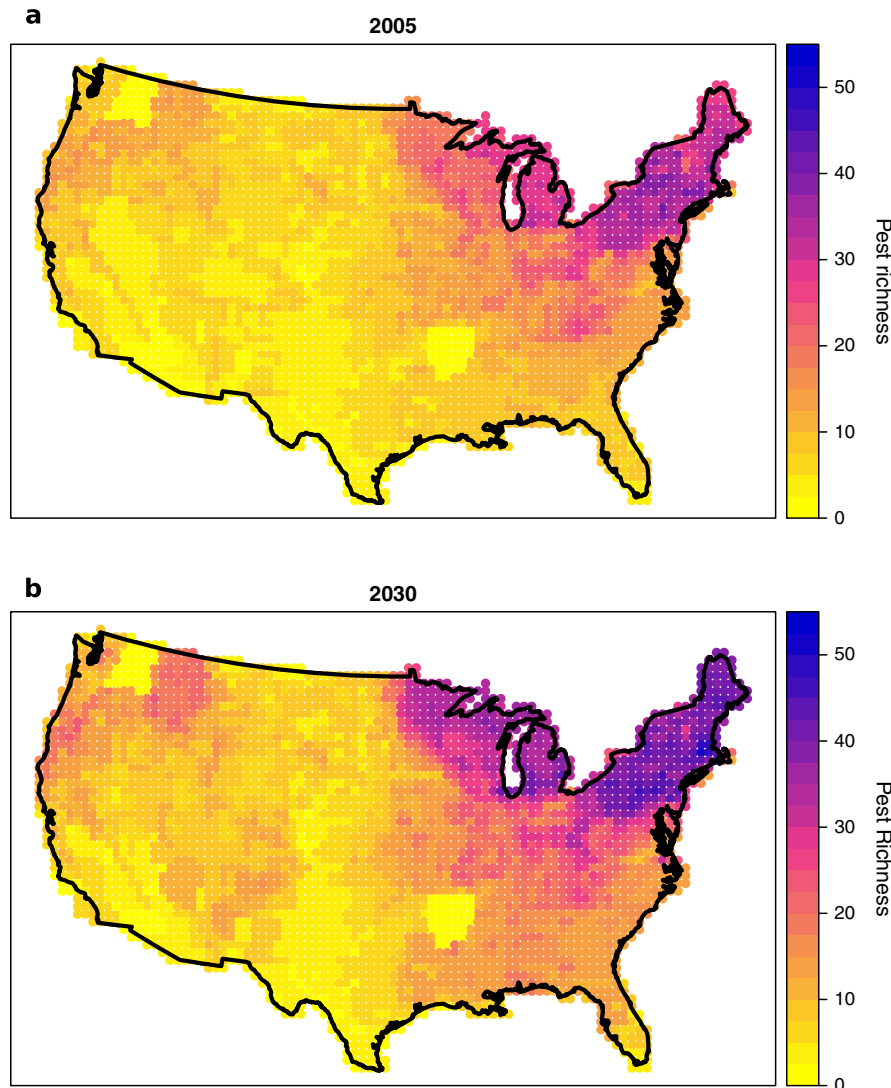
**a**

**2005**



**b**

**2030**



Fig. 3. Forecasted pest species richness from (a) 2005 (fitting year of SDK) to (b) 2030 generated by extending simulated spread patterns for each species from the best-fitting SDK parameters.

processes: (1) species specificity in spread mechanisms or (2) differences in power across models to detect predictors of spread. While there is likely some species specificity in the predictors of spread, we also believe that these differences reflect noisiness of single-species data (i.e., power to detect predictors of spread), given that the SDK was more predictive than customized models, and used the same predictors as the GDK. Arguably, the general model could "borrow" power from numerous species, where spread processes are partially consistent across species. The incorporation of multispecies' information has also allowed for recent advances in analogous fields, such as species distribution modeling, improving spatial predictions of individual species occurrences (Fithian et al. 2015, Leung et al. 2019).

*Semi-generalized Dispersal Kernel approach*

We explored the value of using the GDK as a basic structure upon which to build models, adding context-specific information where it was known. This yielded a 17% average improvement compared to the fully customized model (and 34% improvement compared to GDK). Both GM and BBD customized models were already highly predictive, and SDK yielded modest improvements (3% and 6%, respectively). However, for HWA, by including the three additional corrections to the GDK, SDK yielded a 40% increase in spatial variation explained (from 45% customized model to 85% SDK). Well-documented biases exist in HWA's spread pattern to support the incorporation of niche limitations (Paradis et al. 2008, Morin et al. 2009). Thus, we
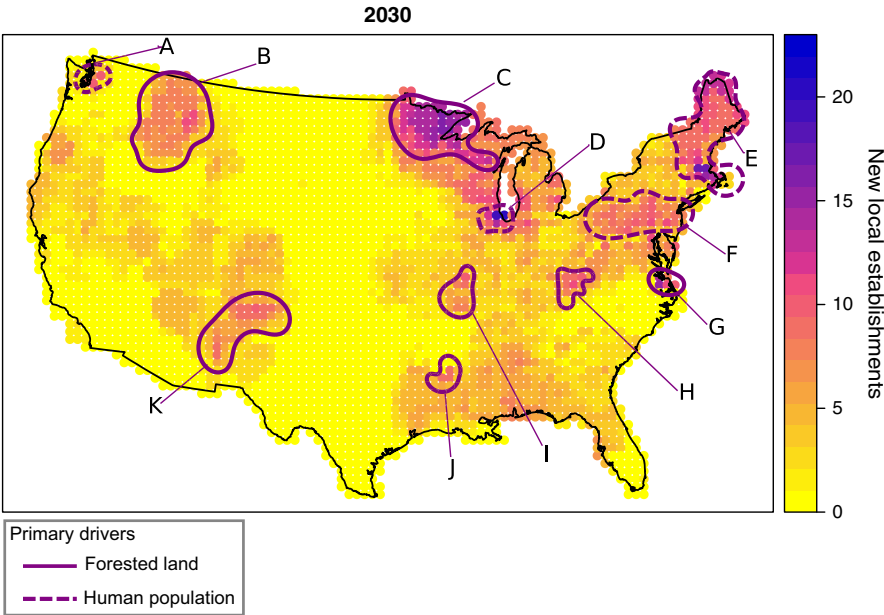
**2030**



FIG. 4. Newly predicted local establishments (for existing United States invasive forest pests) between years 2005 and 2030, created by subtracting Fig. 3b from Fig. 3a. Areas of particular interest are labelled and dominant mechanisms promoting new invasions are denoted with dashed vs. solid lines. A, Seattle, Washington region; B, northern Idaho and western Montana (includes Kootenai, Nez Perce-Clearwater, and Flathead National Forests); C, northern Minnesota and Wisconsin (includes Kabetogama state forest); D, Chicago, Illinois region; E, northern New England (Maine, New Hampshire, Vermont, and Massachusetts), where blue represents the Boston, Massachusetts region; F, Pennsylvania and New Jersey; G, Chesapeake, Virginia region; H, Huntington, West Virginia region; I, Saint Louis, Missouri region; J, Monroe, Louisiana region (includes Upper Ouachita National Wildlife Refuge); K, Carson and Gila National Forests, New Mexico.
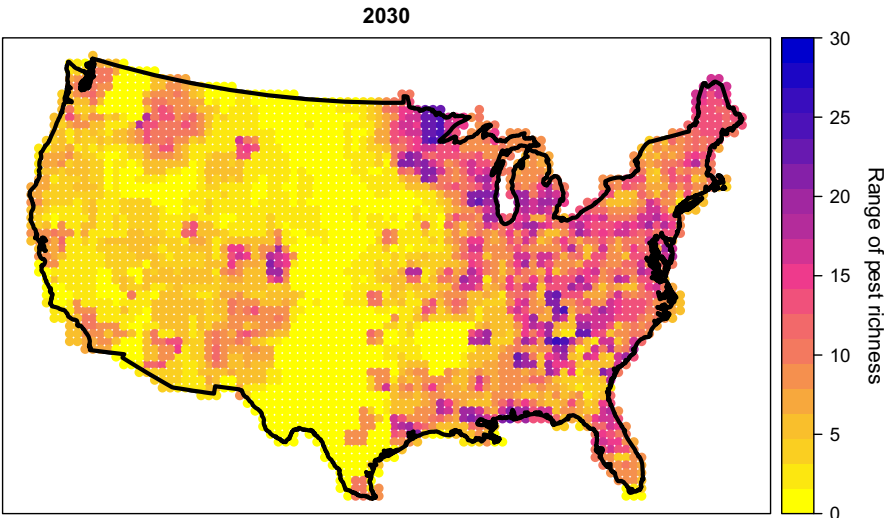
**2030**



FIG. 5. Projections of combined SDK uncertainty at 2030 (range of predicted pest richness at each site) arising from two climate change scenarios (rcp2.6 and rcp8.5 BIOCLIM scenarios), two alternative population growth scenarios (SSP3: "Regional Rivalry", SSP: "Fossil-fueled Development"), and from a sensitivity analysis of model parameters (see Appendix S6).

recommend integrating such context-specific information into an SDK, following fitting protocols described in this manuscript.

Given that spread models have previously been parameterized using these same time-series data, it was

reasonable to expect that customized models would be highly predictive across species. However, the customized model for HWA was only moderately predictive, providing a cautionary tale that even models using the best available data may not produce highly predictive

forecasts. It may be that the quality of data is too poor to build single-species models, in some cases. On the other hand, it appears that a snapshot of a species distribution, a known initial invasion location, and when necessary, a known niche limitation, synthesized with a semi-generalized model (SDK), can outperform the best-fitting customized model.

### Forecasts of future invasion risk

Our simulations suggest future invasions to be even more aggregated in space. Urban centers, areas of high forest cover and tree density appear to be large-scale attractors of invasive propagules from all sources (sensu Colunga-Garcia et al. 2010, Gaertner et al. 2017). Surprisingly, in the GDK, these attractors do not send out as many propagules as other sites, leading to fewer surrounding invasions than if they were also major sources of propagule pressure. Instead, invaders arrive at these sites (sometimes up to 15–20 new pests in the next 25 yr in areas like Chicago, Illinois and Boston, Massachusetts), and remain there, perhaps due to a lack of favorable conditions elsewhere.

While new establishments driven by population density such as Boston, Massachusetts and Chicago, Illinois are relatively unsurprising "hotspots" of future pest load, those driven by forested land, such as Chesapeake, Virginia, western Montana and northern Idaho, and the national forests of New Mexico, are less obvious. Many of these regions coincide with National and State Forests and National Wildlife Refuges, highlighting the role of forested land in the dispersal model. Moreover, areas such as Northern Minnesota and Wisconsin are projected to be the largest "hotspot," possibly reflecting historically low numbers of establishments and hence a lack of saturation in comparison to the northeastern United States or the Midwest. These results indicate a high risk of spread into Canada in the Great Lakes region, and may support a management regime that prioritizes limiting propagule entry to these hubs, though an explicit analysis of the consequences of this prioritization requires further study.

We note that while the intercept and niche corrections can only be employed once a species begins spreading and has a substantial enough distribution for these limitations to be fit, the starting point-corrected GDK can be used as a first pass to predict invasion risk of new invaders from likely points of entry, as the only information it requires is an estimated initial introduction location. If the pest does successfully establish, an SDK combining additional corrections based on model fit can more closely hone in on its future trajectory. For species with well-known niche limitations such as HWA, niche limitations can be similarly incorporated by maximizing SDK fit to the observed distribution, once distributional information is available, as we have done here.

### Caveats and limitations

As detailed above, even with our best current models, there is substantial uncertainty in future pest distributions, given available data. Intuitively, such uncertainty will commonly occur, and we argue that invasive species models should be validated using temporal data withheld from fitting, where possible.

Our model was based on current and historical conditions. However, climate change could alter environmental suitability either due to its direct influence on the invading species or indirectly via effects on hosts and other species (Hellmann et al. 2008, see Appendix S5 for additional species with climatic limitations). However, we note that much of the Northeast, Midwest, and central United States is predicted to have colder minimum winter temperatures with climate change, even if mean temperatures are predicted to increase (Appendix S3), which will lead to more complex future spread dynamics for temperature-limited species. According to the SDK, HWA will be even more constrained with climate change (Appendix S3). Spread could also be affected by conditions becoming hotter or more humid with climate change, potentially affecting GM (Tobin et al. 2014, though these limitations might improve GM forecasts by only 2% based on our analyses). Fortunately, the SDK can easily parameterize any type of spatial limitation for any pest (though these can only be validated using time-series information), and can thus incorporate future knowledge of pest distributional thresholds.

The validation set used in this analysis was not a random selection of species. Instead, it included the three species with time-series data, for which comparative analyses of general vs. customized models could be conducted. It was notably useful from an applied perspective, as they represent some of the most damaging invasive forest pests (Aukema et al. 2011). Emerald ash borer has caused more damage than these species, but was not included in this analysis due to its short invasion history. Further, while fine-scale spatiotemporal GM data are available from pheromone trapping, we applied the GDK to new detections at the 50 × 50 km grid scale in 5-yr time steps, which represent a much coarser spatiotemporal dispersal pattern. To account for the finer-resolution dispersal, a second sub-model could be developed for small-scale dispersal and integrated into country-scale model (although such data do not presently exist for species other than GM).

### Conclusion

While customizing models for each species based on their ecological context yielded 17% higher predictive power compared to the fully generalized GDK, combining both into the SDK yielded the most powerful approach, outperforming the customized model by an additional 17% of spatial variation explained. These results show that the spread process has a substantial

component that is generalizable, and that this generality can be effectively synthesized with context-specific information. The SDK is a strong predictive tool to examine the future distributions of these pests, which we predict are becoming increasingly aggregated at urban centers and are beginning to invade less populated areas with high numbers of trees. These forecasts can aid in estimating future damages due to invasive forest pests, and will be helpful in optimizing future management by highlighting areas of high future pest risk.

## Literature Cited

Aslan, B., and G. Zech. 2005. New test for the multivariate two-sample problem based on the concept of minimum energy. Journal of Statistical Computation and Simulation 75:109–119.

Aukema, J. E., et al. 2011. Economic impacts of non-native forest insects in the continental United States. PLoS ONE 6: e24587.

Aylor, D. E. 1990. The role of intermittent wind in the dispersal of fungal pathogens. Annual Review of Phytopathology 28:73–92.

Bigsby, K. M., P. C. Tobin, and E. O. Sills. 2011. Anthropogenic drivers of gypsy moth spread. Biological Invasions 13:2077–2090.

Bradie, J., and B. Leung. 2015. Pathway-level models to predict non-indigenous species establishment using propagule pressure, environmental tolerance and trait data. Journal of Applied Ecology 52:100–109.

Colunga-Garcia, M., R. A. Haack, R. A. Magarey, and M. L. Margosian. 2010. Modeling spatial establishment patterns of exotic forest insects in urban areas in relation to tree cover and propagule pressure. Journal of Economic Entomology 103:108–118.

Evans, M. R., et al. 2013. Do simple models lead to generality in ecology? Trends in Ecology & Evolution 28:578–583.

Fick, S. E., and R. J. Hijmans. 2017. WorldClim 2: new 1-km spatial resolution climate surfaces for global land areas. International Journal of Climatology 37:4302–4315.

Fithian, W., J. Elith, T. Hastie, and D. A. Keith. 2015. Bias correction in species distribution models: pooling survey and collection data for multiple species. Methods in Ecology and Evolution 6:424–438.

Gaertner, M., J. R. Wilson, M. W. Cadotte, J. S. MacIvor, R. D. Zenni, and D. M. Richardson. 2017. Non-native species in urban environments: patterns, processes, impacts and challenges. Biological Invasions 19:3461–3469.

Gilbert, M., J. C. Grégoire, J. F. Freise, and W. Heitland. 2004. Long-distance dispersal and human population density allow the prediction of invasive patterns in the horse chestnut leafminer Cameraria ohridella. Journal of Animal Ecology 73:459–468.

Grayson, K. L., and D. M. Johnson. 2018. Novel insights on population and range edge dynamics using an unparalleled spatiotemporal record of species invasion. Journal of Animal Ecology 87:581–593.

Haack, R. A. 2006. Exotic bark-and wood-boring Coleoptera in the United States: recent establishments and interceptions. Canadian Journal of Forest Research 36:269–288.

Hansen, B. B. 2007. Flexible, optimal matching for observational studies. R News 7:18–24.

Hauer, M. E. 2019. Population projections for US counties by age, sex, and race controlled to shared socioeconomic pathway. Scientific data 6:190005.

Hellmann, J. J., J. E. Byers, B. G. Bierwagen, and J. S. Dukes. 2008. Five potential consequences of climate change for invasive species. Conservation Biology 22:534–543.

Houston, D. R. 1994. Major new tree disease epidemics: beech bark disease. Annual Review of Phytopathology 32:75–87.

Hudgins, E. J., A. M. Liebhold, and B. Leung. 2017. Predicting the spread of all invasive forest pests in the United States. Ecology Letters 20:426–435.

Koch, F. H., and W. D. Smith. 2008. Spatio-temporal analysis of *Xyleborus glabratus* (Coleoptera: Circulionidae: Scolytinae) invasion in eastern US forests. Environmental Entomology 37:442–452.

Kot, M., M. A. Lewis, and P. Van Den Driessche. 1996. Dispersal data and the spread of invading organisms. Ecology 77:2027–2042.

Kovacs, K. F., R. J. Mercader, R. G. Haight, N. W. Siegert, D. G. McCullough, and A. M. Liebhold. 2011. The influence of satellite populations of emerald ash borer on projected economic costs in US communities, 2010–2020. Journal of Environmental Management 92:2170–2181.

Leung, B., E. J. Hudgins, A. Potapova, and M. Ruiz-Jaen. 2019. A new baseline for countrywide α-diversity and species distributions: illustration using > 6000 plant species in Panama. Ecological Applications 29(3):e01866.

Liebhold, A., V. Mastro, and P. W. Schaefer. 1989. Learning from the legacy of Leopold Trouvelot. Bulletin of the Entomological Society of America 35:20–21.

Liebhold, A. M., J. A. Halverson, and G. A. Elmes. 1992. Gypsy moth invasion in North America: a quantitative analysis. Journal of Biogeography 19:513–520.

Lodge, D. M., et al. 2006. Biological invasions: recommendations for US policy and management. Ecological Applications 16:2035–2054.

Lovett, G. M., M. Weiss, A. M. Liebhold, T. P. Holmes, B. Leung, K. F. Lambert, et al. 2016. Nonnative forest insects and pathogens in the United States: Impacts and policy options. Ecological Applications 26:1437–1455.

McCullough, D. G., and R. J. Mercader. 2012. Evaluation of potential strategies to SLow Ash Mortality (SLAM) caused by emerald ash borer (*Agrilus planipennis*): SLAM in an urban forest. International Journal of Pest Management 58:9–23.

Morin, R. S., A. M. Liebhold, P. C. Tobin, K. W. Gottschalk, and E. Luzader. 2007. Spread of beech bark disease in the eastern United States and its relationship to regional forest composition. Canadian Journal of Forest Research 37:726–736.

Morin, R. S., A. M. Liebhold, and K. W. Gottschalk. 2009. Anisotropic spread of hemlock woolly adelgid in the eastern United States. Biological Invasions 11:2341–2350.

Paradis, A., J. Elkinton, K. Hayhoe, and J. Buonaccorsi. 2008. Role of winter temperature and climate change on the survival and future range expansion of the hemlock woolly adelgid

(*Adelges tsugae*) in eastern North America. Mitigation and Adaptation Strategies for Global Change 13:541–554.

Seebens, H., et al. 2015. Global trade will accelerate plant invasions in emerging economies under climate change. Global Change Biology 21:4128–4140.

Sharov, A. A., D. Leonard, A. M. Liebhold, E. A. Roberts, and W. Dickerson. 2002. "Slow the spread": a national program to contain the gypsy moth. Journal of Forestry Research 100:30–36.

Shigesada, N., K. Kawasaki, and Y. Takeda. 1995. Modeling stratified diffusion in biological invasions. American Naturalist 146:229–251.

Skellam, J. G. 1951. Random dispersal in theoretical populations. Biometrika 38:196–218.

Skuhravá, M., V. Skuhravý, and G. Csóka. 2007. The invasive spread of the gall midge *Obolodiplosis robiniae* in Europe. Cecidology 22:84–90.

Taylor, R. A. J., L. S. Bauer, T. M. Poland, and K. N. Windell. 2010. Flight performance of *Agrilus planipennis* (Coleoptera: Buprestidae) on a flight mill and in free flight. Journal of Insect Behavior 23:128–148.

Tobin, P. C., D. R. Gray, and A. M. Liebhold. 2014. Supraoptimal temperatures influence the range dynamics of a non-native insect. Diversity and Distributions 20:813–823.

Vilà, M., and P. E. Hulme, editors. 2017. Impact of biological invasions on ecosystem services. Volume 12. Springer, Cham, Switzerland.

Ward, S. F., S. Fei, and A. M. Liebhold. 2019. Spatial patterns of discovery points and invasion hotspots of non-native forest pests. Global Ecology and Biogeography. https://doi.org/10.1111/geb.12988.

Ward, J. S., M. E. Montgomery, C. A. S.-J. Cheah, B. P. Onken, and R. S. Cowles. 2004. Eastern hemlock forests: guidelines to minimize the impacts of hemlock woolly adelgid. NA-TP-03-04. USDA Forest Service, Northeastern Area State and Private Forestry, Morgantown, West Virginia, USA.

## Supporting Information

Additional supporting information may be found online at: http://onlinelibrary.wiley.com/doi/10.1002/eap.1988/full

## Data Availability

Model predictions, associated data, and code are available on Zenodo: https://doi.org/10.5281/zenodo.3343027