**Global correlates of emerging zoonoses: anthropogenic, environmental, and biodiversity risk factors**

Authorship:

Toph Allen, Kris A. Murray, Carlos Zambrana-Torrelio, Stephen S. Morse, Carlo Rondinini, Moreno Di Marco, Nathan Breit, Kevin J. Olival, Peter Daszak*

*Corresponding Author

**Abstract**

Zoonoses originating from wildlife represent a significant threat to global health, security and economic growth, and combating their emergence is a public health priority. However, our understanding of the mechanisms underlying their emergence remains rudimentary. Here, we update a global database of emerging infectious disease (EID) events, create a novel measure of reporting effort, and fit boosted regression tree models to analyse the demographic, environmental and biological correlates of their occurrence. After accounting for reporting effort, we show that zoonotic EID risk is elevated in forested tropical regions experiencing land-use changes and where wildlife biodiversity (mammal species richness) is high. We present a new global hotspot map of spatial variation in our zoonotic EID risk index, and partial dependence plots illustrating relationships between events and predictors. Our results may help to improve surveillance and long-term EID monitoring programs, and design field experiments to test underlying mechanisms of zoonotic disease emergence.

**Introduction**

Emerging infectious diseases (EIDs) are a significant and growing threat to global health, global economy and global security [1,2]. Analyses of their trends suggest that their frequency and their economic impact are on the rise [3,4], yet our understanding of the causes of disease emergence is incomplete. The majority of EIDs (and almost all recent pandemics) originate in animals, mostly wildlife, and their emergence often involves dynamic interactions among populations of wildlife, livestock, and people within rapidly-changing environments [5-7]. The mechanisms underlying this process are likely complex, and occur in contexts that are often characterized by a paucity of systematically collected data [8].

Global efforts to reduce the impacts of emerging diseases are largely focused on post-emergence outbreak control, quarantine, drug and vaccine development [3]. However, delays in detection of, or response to, newly emerged pathogens combined with increased global urbanization and connectivity have resulted in recent EIDs causing extensive mortality across cultural, political and national boundaries (e.g. HIV), and disproportionately high economic damages (e.g. SARS, H1N1). Efforts to identify the origins and causes of disease emergence at local scales, and regions where novel diseases may be more likely to emerge from, are valuable for focusing surveillance, prevention, and control programs earlier in the chain of emergence, containing EIDs closer to their source, and more effectively limiting their subsequent spread and socioeconomic impacts [8].

A previous analysis of global EID trends modeled the spatial variation of "EID events", representing records of disease occurrence of the first appearance of a pathogen in a human

population related to increased distribution (e.g., new geographic location, new host species), incidence, virulence, or other factors [4]. The EID events were divided into four groups, including wildlife-origin zoonoses [4]. To model the potential risk of disease emergence, these four groups were regressed as a function of human population density and growth, latitude, rainfall, and wildlife species richness. The results suggest that wildlife origin EIDs are more likely to occur in regions with higher human population density and greater wildlife diversity (mammal species richness) [8]. However, the study is limited in its mechanistic inference due, in part, to the lack of specificity of the predictors. For example, the effect of population density could represent anthropogenic environmental changes (human pressure on landscapes), human-animal contact rates, reporting biases, or a combination of these. Furthermore, a range of potential mechanisms may not be adequately represented by this predictor set; a lack of an effect of rainfall, for example, does not discount the potential for other climatic factors to play a role, and a lack of an effect of latitude could mean that it is simply a poor proxy for other more meaningful factors that nevertheless exhibit some latitudinal variation (e.g., temperature, habitat types, biodiversity, GDP). Improving the predictor set to better target underlying mechanisms could improve model performance and our ability to explain spatial variation in EID risk.

The current study aims to better analyse the mechanistic underpinnings of disease emergence for zoonotic EIDs of wildlife origin, while addressing some methodological limitations of Jones et al. [4]. We focus on EIDs of wildlife origin, which are responsible for nearly all recent pandemics (e.g. Ebola, MERS), constitute the majority of the high impact EIDs from the last few decades, and are a significantly growing proportion of all EIDs combined [4]. We updated the EID

database from [4], and employed a new modelling framework (boosted-regression trees, BRT) to capture high-dimensional interactions and generate response functions for individual variables. We selected a refined set of spatial predictors for their relevance to *a priori* hypotheses on plausible mechanisms underlying zoonotic EID emergence, including proxies for human activity, environmental factors, and the zoonotic pathogen pool from which novel diseases could emerge, all key features of conceptual models of zoonotic spillover [7-11]. We used an improved dataset of mammal species distributions[12], and included numerous datasets on measures of land use, land-use change and land cover. Further, all datasets with sufficient temporal coverage were matched to events in the EID database by decade, such that covariates more accurately reflect the prevailing conditions at the time of disease emergence. We also constructed a novel proxy of reporting effort to match the spatial resolution of the other predictors, where previous studies have relied on coarse, country level measures, and compared EID risk predictions with and without corrections for reporting effort. Finally, we accounted for spatial uncertainty in EID event data by random resampling to explicitly take into account the difficulties of accurately geocoding EID events.

**Results**

After factoring out reporting effort (weighted model), evergreen broadleaf trees (median 7.6% of predictive power), human population density (6.9%), Global Environmental Stratification (climate) (5.9%), and mammal species richness (an aspect of biodiversity) (5.6%) had the largest relative influence over the distribution of EID events (Figure 1). Across 1,000 iterations of the model, no variables consistently emerged as much stronger predictors than others but an average ranking of predictor importance could be derived. Of the top predictors, evergreen

4

broadleaf trees (representing tropical rainforests) exhibited an overall positive trend, human

population density an overall negative trend, the Global Environmental Stratification (climate)

an idiosyncratic trend towards warmer and wetter (i.e., more tropical) climates, and mammal

species richness showed an idiosyncratic trend, with higher risk values at both lower and higher

richness values (Figure 2). After mammal species richness, three variables involving agricultural

practices followed in importance: cultivated/managed vegetation (5.6%), pasture change

(5.2%), and areas dedicated to pasture (5.1%). In the unweighted model, not accounting for

reporting effort (see Supplementary Results 2), urban/built-up land was by far the strongest

predictor of observed events, explaining a median of 30.6% of the model's variation and
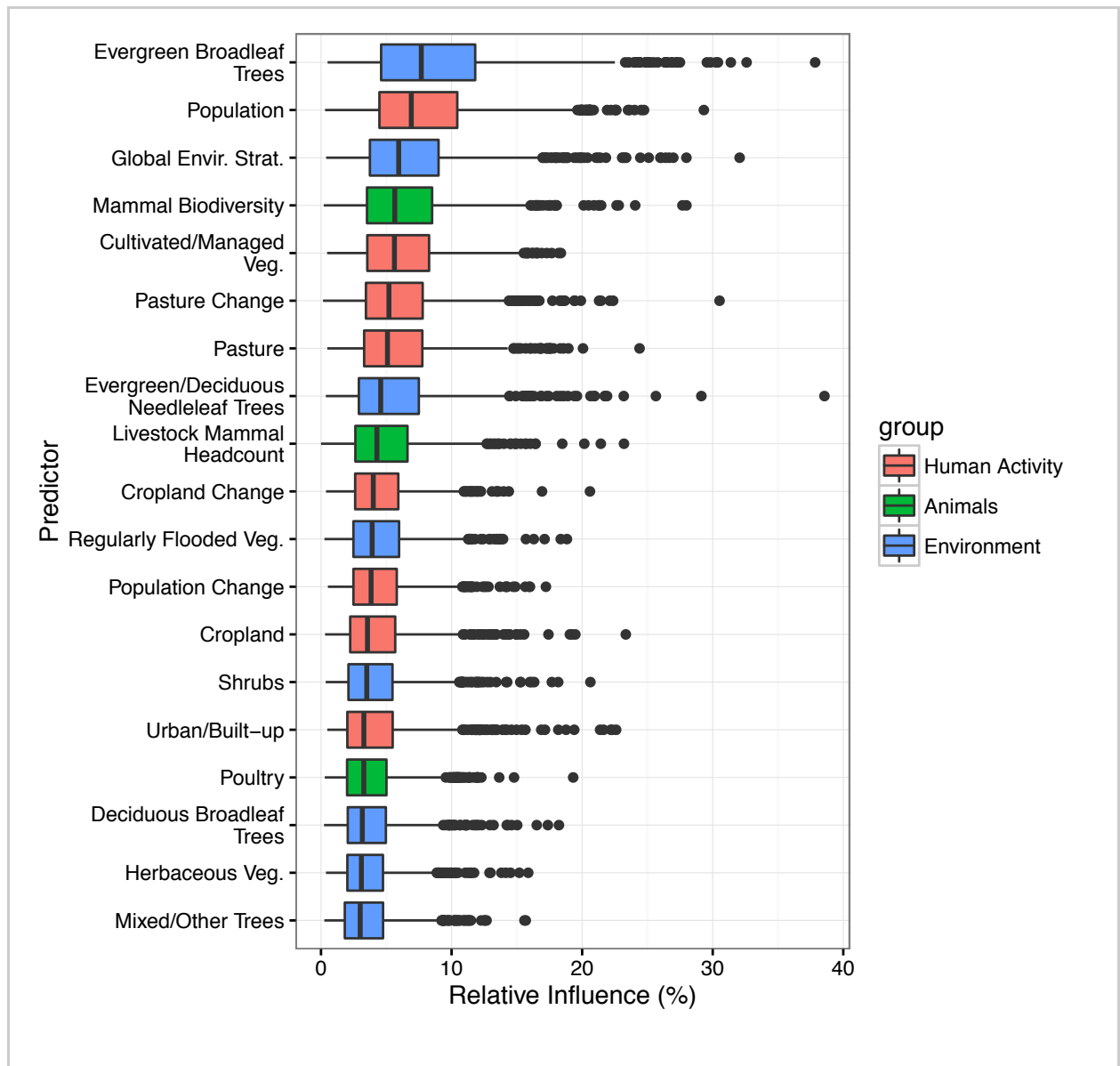
exhibiting a distinct positive trend.

Figure 1: The relative influence of predictors on EID event occurrence probability. The box plots show the spread of relative influence across 1000 replicate model runs to account for uncertainty in EID event location (see above). BRTs do not provide p-values or coefficients, but rank variables by their relative influence in explaining variation in the outcome [13].
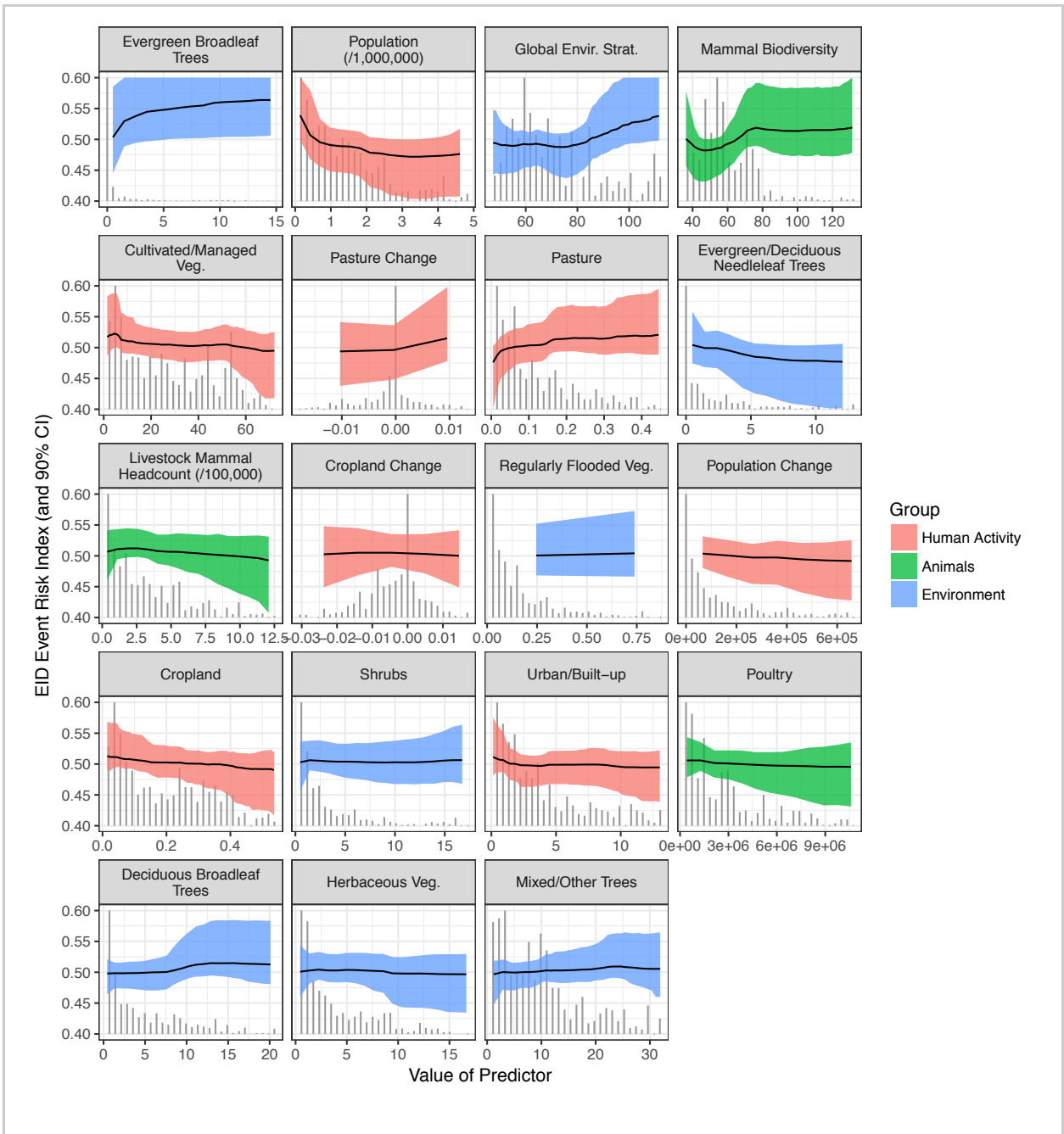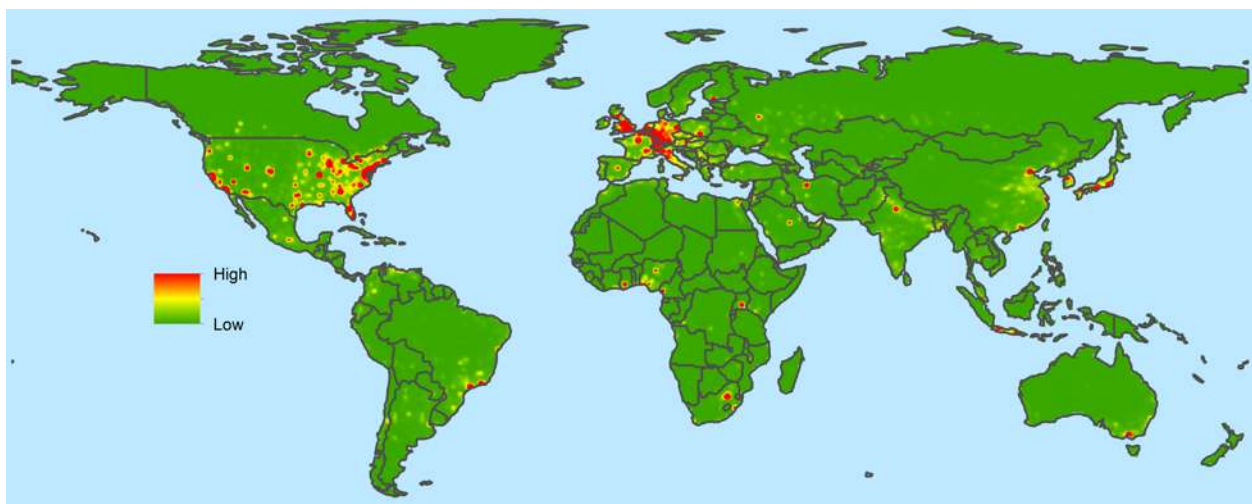
Figure 2: Partial dependence plots for the zoonotic EID events for all predictors in the boosted regression tree (BRT) model (ordered by relative influence). X-axes show the range from the 10<sup>th</sup> to 90<sup>th</sup> percentiles of sampled values of predictors (e.g. number of mammal species per

grid square (mammalian richness) or percentage of land cover type and its change per grid square), the distribution of which are displayed in the histograms (grey bars). Y-axes show the EID event risk index. Black lines show the median and coloured areas show the 90% confidence intervals, computed using a bootstrap resampling regime incorporating uncertainty in EID event locations. Our sampling regime fixes the outcome, which indexes EID event risk, between 0 and 1, with a mean at 0.5. Y-axes are centred around the mean and scaled to 0.1 above and below. Partial dependence plots display the response for an individual variable in the model while holding all other variables constant[13,14]. They allow a visualization of what are mostly non-linear relationships between drivers and the EID event risk index (in this case, after reporting effort is factored out.). See Supplementary Results 2 for results of the model unweighted by reporting effort.

Relative to the observed risk index for EID events, the model's estimated risk index correcting for reporting bias (Figure 3) is more concentrated in tropical regions. Areas of higher suitability for EID occurrence are fairly evenly distributed across the continents, with no major land mass free from areas predicted to be suitable for EIDs. In particular, areas of high population outside the tropics, such as cities in Europe, the United States, Asia and Latin America remain among areas at the high end of the risk index. Tropical regions in North America, Asia, central Africa, and regions of South America have more extensive areas of predicted EID occurrence.

Our model validation statistics were computed both for the weighted model — with a background, or absence, sample weighted by reporting effort, effectively computing statistics on the residuals of that variable — and our unweighted model, using a background sample

uniform across land area. The weighted bootstrap model reported a median of 31.6% of deviance explained across the 1000 replicate models (empirical 90% CI 15.9% to 50.5%), whereas the unweighted model explained a median 50.2% of deviance (empirical 90% CI 35.8% to 67.2%). Our weighted model's cross-validation statistics, computed over 100 runs of ten-fold cross-validation, varied depending on the weighting of the null validation sample. With validation absences weighted by reporting effort, the weighted model had a median AUC of 0.64, with an empirical 90% confidence interval ranging from 0.54 to 0.69 (out of possible values between 0 and 1, with 0.5 indicating performance no better than random). The median True Skill Statistic (TSS) was 0.23 with an empirical 90% CI of 0.14 to 0.33 (out of a range of -1 to 1). These indicate low to moderate predictive performance [15-17]. Evaluated against an unweighted null, the weighted model had a median AUC of 0.78 (90% CI [0.75, 0.81]) and a median TSS of 0.43 (90% CI [0.37, 0.50]). The unweighted model evaluated against to an unweighted null, had a median AUC of 0.77 (90% CI [0.73, 0.81]) and a median TSS of 0.44 (90% CI [0.37, 0.50]).
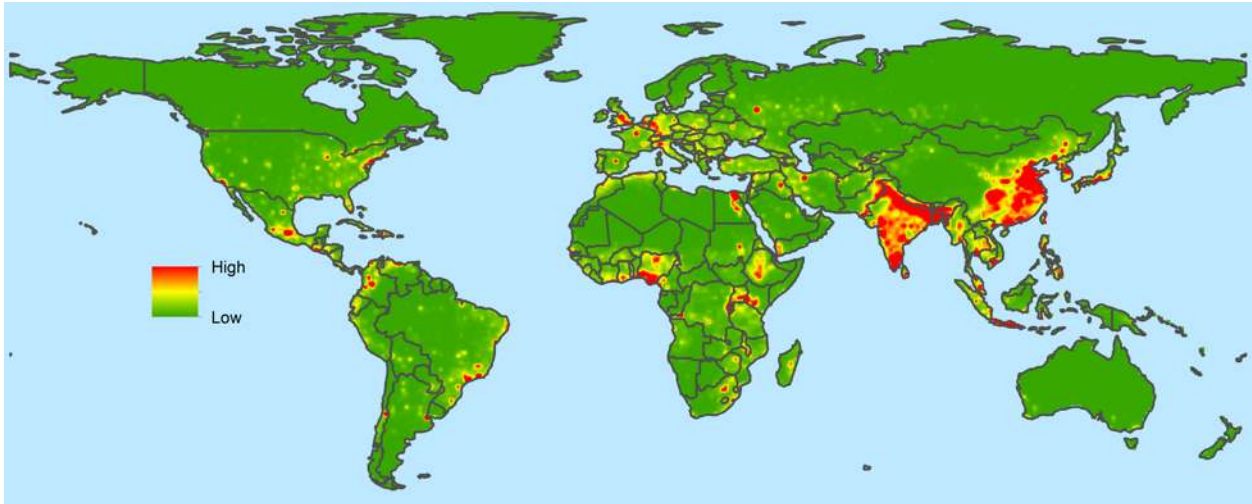
Figure 3: Heat maps of predicted relative risk distribution of zoonotic EID events. The top panel shows the predicted distribution of new events being observed (weighted model output with current reporting effort); the bottom panel shows the estimated risk of event locations after factoring out reporting bias (weighted model output reweighted by population). See Figure 4 for raw weighted model output. Maps were created in ArcGIS 10.2.2 using standard deviation scaling, with the colour palette scaled to 2.5 standard deviations above and below the mean [18].
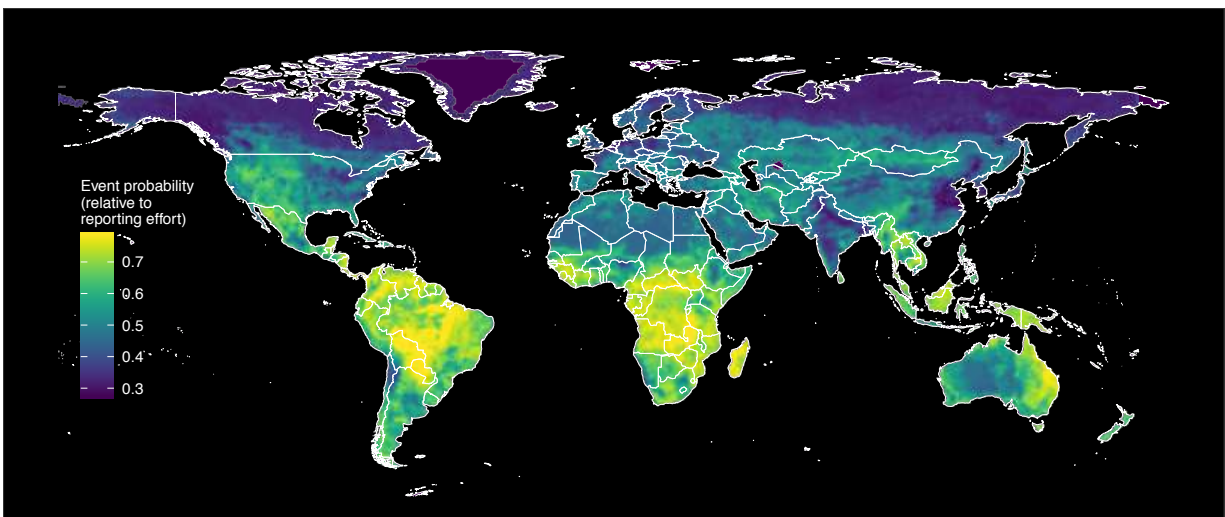
Figure 4. Heat map of weighted model response, i.e. EID risk relative to reporting effort. Value

indicates the binomial probability that a grid cell sampled at that location will contain an EID

event as opposed to a background sample, when drawing equal numbers of absence and

background samples weighted by reporting effort (see Materials & Methods). This layer was

weighted by reporting effort to produce the "observed" EID risk index map (Figure 3, upper

panel) and by population to produce the risk index map with bias factored out (Figure 3,

bottom panel).

## Discussion

We developed a spatial model to describe the global spatial patterns of zoonotic EIDs. The

model factored out (referred to as the weighted model) clear effects of reporting effort, which

otherwise biases our ability to interpret EID event observations. Our model ranked risk factors

according to their predictive power, capturing both their main effects and potential interactions

with other variables, and we derived the directionality and shape of their relationships to EID

events for graphical interpretation.  Our results suggest that the risk of disease emergence is

elevated in tropical forest regions, high in mammal biodiversity, and experiencing

anthropogenic land use changes related to agricultural practices [19-21].

The link between mammal biodiversity and zoonotic disease emergence has been identified

previously [4] and hypothesized widely [8,22]. Areas with tropical forest and high mammalian

biodiversity were elevated on our EID risk index (henceforth "EID risk"), although uncertainty of

the estimates was high). It may be that these variables represent the same mechanism, as

tropical forests are generally areas of high biodiversity [23], and the apparent causal effect may

be attenuated by the presence of both in the model. This trend is consistent with existing hypotheses that suggest greater host biodiversity increases the 'depth' of the pathogen pool from which novel pathogens may emerge, which in turn increases the potential for novel zoonotic pathogens to emerge [24]. There is a large literature on the relationship between biodiversity and infectious disease risk in people, with some studies suggesting that high host biodiversity decreases risk or that biodiversity loss may increase risk (i.e. the dilution effect) [25], while others refute the generalizability of this [26,27] or suggest disease richness or prevalence increases with increasing wildlife species richness [15]. Our findings look at the global scale and a large group of pathogens, and so do not speak directly to this debate: although the dominant trend is an increase in risk of disease emergence with higher mammalian richness, this neither rules out nor substantiates the possibility of a dilution effect for specific diseases. Rather, it is consistent with previous suggestions that the relationship between biodiversity and disease risk is complex, context-specific and idiosyncratic [26].

When not accounting for reporting effort ("unweighted"), our model showed urban land as having a very strong positive association with EID events. However, this can be interpreted as an effect of reporting bias, since (1) urban land was also strongly associated with our measure of reporting effort, and (2) fitting our weighted model, relative to reporting effort, attenuated this effect. Similarly, although population density was not found to be an important predictor in the unweighted model (median relative influence 2.2%), weighting the model by reporting effort drove up its importance (median rel. inf. 6.9%), such that EID risk was inversely related to population density. Population density was also included in the reporting effort model, but was not as strong a predictor (rel. inf. 3.6%) as urban land (rel. inf. 45.2%). Theoretically, population

has a baseline multiplicative effect on human disease events [28] — of which EID events are a subclass — and their detection is modulated by reporting effort. Reporting effort appears to be associated with urbanization, but reporting effort and urbanization are also both products of human population. We did not attempt to fully disentangle these factors, instead using our measure of reporting effort to present a map of emerging infectious disease hotspots with bias "factored out" (described below in Materials & Methods).

Our reporting effort measure was created by matching placenames in a subset of the biomedical literature. The BRT model of reporting effort model suggested that the distribution of this effort was strongly and positively related to urban areas. This could be because our extraction of placenames biases the outcome toward urban areas, or it may accurately represent the true distribution of reporting towards urban areas, or a combination of the two. In either case, our reporting effort dataset is likely to be a large improvement over previous studies that have used country-level data to control heterogeneous reporting effort in better-than-country-level spatial analyses of disease risks used in similar previous studies[4,28] (detailed fully in Supplementary Methods: PubCrawler).

The work presented here builds on previous research [4] in a number of important ways to advance our understanding of wildlife-origin zoonotic disease emergence. Firstly, our model building approach explores the explanatory value of a large collection of globally-gridded data on environmental, demographic, and host diversity variables, including newly developed models of mammal distributions and richness patterns. This has allowed us to close the gap between predictors and *a priori* mechanistic hypotheses specifically relevant to zoonotic disease emergence from wildlife reservoirs. Secondly, we adopted a machine-learning

modelling approach (boosted regression trees) suited to the analysis of complex ecological data[13], and used various resampling regimes to measure and visualize multiple sources of uncertainty (model uncertainty, spatial uncertainty of EID events, temporal uncertainty of covariates matching with events) and predictive performance. Thirdly, we have attempted to improve how the model accounts for uneven global distribution of surveillance and research on disease event detection (i.e. report effort). This includes an algorithm-based approach to more realistically map reporting effort and shows the significant implications that a finer-scale, sub-national resolution variable for reporting effort can have for a model. Finally, we were able to temporally match predictors to events.

Despite using a more flexible modelling framework, there are limitations to our approach. When differentiating between EID events and a uniformly-weighted background sample, our weighted and unweighted models an AUC of 0.78 and 0.77, and a TSS of 0.43 and 0.41 respectively, indicating moderate predictive performance. However, against a background sample weighted by reporting effort, our weighted model had an AUC of 0.61 and a TSS of 0.18, indicating low–moderate performance. These statistics indicate much unexplained variation. While broad changes in zoonotic EID relative risk are evident in the partial dependence plots, in areas of elevated risk confidence intervals are generally wide enough that quantitative relationships remain uncertain.

Wherever possible, we tried to define and incorporate uncertainty into our model (for example, correcting for uncertainty in location by sampling EID events from within known areas of occurrence, and correcting for literature-level biases by weighting background samples by our measure of observation effort). Multiple factors contribute to this uncertainty. Firstly, analyses

were conducted using gridded data at 1º WGS84 resolution (c. 100 km at the equator), the same resolution used previously [4]. Our choice of resolution for predictor datasets was constrained by data availability, since all were downscaled to the lowest common spatial resolution. Secondly, confidence intervals are widest in regions for each variable where fewer grid cells were sampled. Since our weighted model sampled fewer grid cells proportional with reporting effort, these represent areas where more reporting effort — including ground-truthing studies — may increase confidence. Thirdly, another limitation shared with [4] is the underlying accuracy and suitability of EID event data, which were drawn from a review of published literature. Individual studies, though, carry their own biases, inaccuracies, and different approaches to collecting and documenting data, and this alone adds an unknown amount of imprecision and potential bias to our outcome dataset. Finally, our goal of creating a single model, to look for common trends in emerging wildlife-origin zoonotic diseases, likely imposes limitations on the specificity of trends we can examine. In reality, different classes of diseases (e.g., viruses versus bacteria) and indeed individual diseases have their own unique ecology, with different drivers and sets of conditions being more or less important in shaping the emergence process [29]. Because of these limitations, we refrain from making specific (e.g. city-by-city) interpretations of the model's output, rather noting broad trends in geographic regions and environment types of intererest.

Wide confidence intervals in areas of elevated EID risk suggest areas for future study, and underscore the need for targeted long-term disease surveillance and monitoring in these areas. Collection of more accurate spatiotemporal data on events surrounding disease emergence, including initial emergence events, using a combination of large scale field research (e.g.

USAID's PREDICT project [30]) and digital disease detection tools [31] would help alleviate this issue in the future by generating more consistent data on a larger scale, potentially automatically [32]. Additionally, we propose to launch an editable version of the EID event database to crowd-source data using the EID research community and ultimately the public and improve the overall spatial resolution and number of events [33]. These datasets will aid efforts to better define the point at which a disease becomes 'emerging', allowing the programmatic definition and examination of different definitions of emergence (e.g. first appearance vs. increasing incidence, etc.) in testable form [34].

Future work may be able to enhance the predictive power of this approach by focusing on even tighter classes of disease, taxonomic groups of pathogens, or transmission modes, and building models to forecast changes in risk distribution or to examine more specific mechanistic hypotheses. Efforts to examine the commonalities of disease emergence may benefit from incorporating disease specific models in a hierarchical approach, allowing certain parameters to vary across diseases or disease classes, while pooling other parameters.

Despite shortcomings, our improvements to the earlier model allowed us to find quantitative support for previously only hypothesized factors that increase the risk of EID events. Our findings therefore have broad implications for surveillance, monitoring, control and research on emerging infectious diseases. Like Jones et al., [4], we find that EID events are observed predominantly in developed countries, where surveillance is strongest, but that our predicted risk is higher in tropical, developing countries.

Our spatial mapping has direct relevance to ongoing surveillance and pathogen discovery efforts [35]. It shows that the global distribution of zoonotic EID risk (and the presence of EID 'hotspots') is concentrated in tropical regions where wildlife biodiversity is high and land use change occurring. These regions are likely to be the most cost effective for surveillance programs targeting wildlife, livestock or people for novel zoonoses, and for pandemic prevention programs that build capacity and infrastructure to pre-empt and control outbreaks [30]. Further honing the EID risk index within regions and countries might also inform the planning of large land use change programs such as logging and mining concessions, dam-building, and road development [36]. These activities carry an intrinsic risk of disease emergence by increasing human or livestock contact with wildlife in new regions or by disrupting disease dynamics in reservoir hosts[24,37], and have been repeatedly linked to outbreaks of novel EIDs.

Similarly, the partial dependence plots allow a deeper understanding of the largely non-linear relationships between EID drivers and disease emergence that can be used to design field experiments to test specific and generalizable hypotheses on the drivers of zoonotic disease emergence. These should include field sites along land use gradients within EID hotspot countries where controlled sampling protocols are used to identify how wildlife biodiversity, known and unknown pathogen diversity (e.g., using viral family level degenerate primers for PCR[38]), and human contact with wildlife varies across a landscape. Such an approach will provide a way to identify the fine-scale rules that govern disease emergence and provide a richer understanding of what drives EID risk on-the-ground, a critical extension of this modelling approach.

**Materials & Methods**

All data and code used to generate the models are available on GitHub (doi: 10.5281/zenodo.400978)[39], as is the code used to generate the reporting effort layer (doi: 10.5281/zenodo.400977)[40].

<u>Data sources</u>

*Response variable (zoonotic EID events)*

We followed the definition of an emerging infectious disease and an EID event used in [4] —

specifically, events documented in the scientific literature denoting the first emergence of

pathogen in a human population where that pathogen was classified as "emerging" due to

recent spillover from an animal reservoir, a significant increase in its incidence or geographic

distribution in the human population, a marked change in its pathogenicity or virulence, or

other factors. In this study we focus only on EID events of wildlife origin ('wildlife zoonoses')

because these represent the majority of EID events in the most recent decade studied, are

increasing significantly as a proportion of all EIDs after correcting for reporting bias, include

most of the highest impact EIDs of recent decades (e.g. Ebola viruses, Nipah virus) and almost

all recent pandemics (e.g. pandemic influenza viruses, SARS). Data on EID events were derived

from an updated version of the database originally used by [4], which contained EID events

ranging from 1940 to 2004 (n = 335 total, n = 145 for wildlife zoonoses (43.3% of all EIDs)). We

updated the database to include EID events for wildlife zoonoses through 2008 (n = 224),

following the methodology in [4] so as to include only diseases reported in the peer-reviewed

literature, where there is evidence that a disease is emerging for one of the reasons laid out

above. Additionally, we only included the first emergence of a new disease-causing agent, such

that the MERS Coronavirus was included, but not reports of new strains of Ebola virus. For each

EID event, data were derived from the literature, if available, for date, location (see below),

pathogen genus and species, zoonotic origin and type, and associated or hypothesized drivers,

following [4]. Location data for initial EID emergence events were variable in their geographic

specificity, ranging from precise coordinates to broader regions (e.g., municipalities, counties,

districts) or entire continents depending on details reported in the primary literature. A spatial

polygon was created for each event that represented the most precise municipal region the EID

event was known to have occurred in. All EID event polygons, regardless of precision, were

included in our bootstrap resampling framework; removing those with geographic uncertainty

(e.g. those with only country level resolution) may artificially inflate the certainty of our model;

our resampling scheme limits their impact to appropriate levels. Events with precise

coordinates were also assigned a polygon for consistency of data format, but rather than using

a municipal boundary, the event was assigned a 5 km circular buffer zone. EID polygons were

subsampled for model fitting as described below. Because our model matches EID events with

decadal population and land use data (described below), we restricted our analyses to decades

for which covariate data exist, excluding events before 1970 and leaving n = 147 records for

analysis (66% of wildlife zoonosis events).

*Explanatory variables*

We compiled spatial data layers for 20 predictors in four broad categories to decompose which

factors are associated with zoonotic disease emergence. These reflected the most frequently

hypothesized drivers of zoonotic disease emergence and included (Table 1): human

presence/activity, animals/hosts, the environment, and reporting effort. Explanatory variables

came from a variety of data sources, and all were rescaled or transformed to a spatial grid of 1º

resolution (WGS84, c. 110 km at the equator) prior to their use in models. Full details of

sources, original resolutions and rescaling are presented in Tables 1 and 2.

"Human Activity" data were compiled and eight predictors derived based on the following

rationale: **1)** Population density likely influences EID risk in two discrete ways. Firstly, as EID

events are defined as diseases emerging in the human population, their frequency—before the

effects of other predictors—is assumed to be proportional to population density, with the other

predictors modifying the per-person risk of EID events. To represent this, we treated human

population as a baseline multiplicative factor in our models [41]. Secondly, population density

may affect transmission dynamics such that EID events in areas of denser population may be

more likely to produce outbreaks large enough to be detected [42]. We used the Global Rural-

Urban Mapping Project [43] human population dataset, which provides gridded estimates of

human population every five years for 1970–2000. **2)** Population change acts as a proxy for

changing demands on ecosystems leading to environmental perturbation, which has been

hypothesized to drive disease emergence [24]. We created a measure for population change by

calculating the inter-decadal difference of human population per grid cell. **3)** Land-use type

represents largely anthropogenic influence on the landscape (as opposed to 'land cover' below)

and has been hypothesized to play a role in disease emergence and spatial distribution [24,44-47].

We used the HYDE dataset which estimates the percentage of land-use types in each grid cell of

a global dataset every ten years for 1900–2000 [48] to derive predictors representing percentage

of land used for cropland and percentage used for pasture. We also include the layers for Urban

Land and Managed/Cultivated Vegetation from the EarthEnv dataset, described below under

"Environment", in this category, as they index human impact on the environment. **4)** <u>Land-use change</u> has been hypothesized as a key driver for disease emergence by perturbing ecosystems and bringing humans into close proximity with wildlife [5,7,8,24,29]. We created metrics of change for pasture and cropland by calculating the between-decade difference in values for each grid cell for cropland and pasture.

For datasets with multiple temporal layers (human population, cropland, and pasture), we included the intersection of available dates in different datasets (decades 1970–2000) and calculated inter-decadal change layers by differencing consecutive decades. All presence and absence samples drawn for each event (see below) were matched to the nearest decadal layers (years ending in 5 were rounded up) and the change layer for the decade they fell in.

"Animal/host" data were represented by two predictors: **1)** <u>Mammalian biodiversity.</u> The diversity and prevalence in a host population of potentially zoonotic pathogens in an area is hypothesised to be a key factor in the risk of novel pathogen emergence [8,24,49]. However, spatial data on global pathogen diversity do not currently exist, and it is estimated that we have identified less than 1% of mammalian viral diversity [38]. Consistent with previous studies, we therefore assume that the number of available pathogens in an area is proportional to the diversity (species richness) of wildlife species [4,5,38,50]. The overwhelming majority of emerging zoonoses have mammalian hosts [51], and global biogeographic patterns of human infectious diseases is highly correlated with global patterns of mammalian diversity [32]. We therefore used mammal biodiversity (species richness) measured as number of mammal species per grid cell as a proxy for pathogen species richness. To do this, we used the most up to date mammal species distribution maps available, derived from species distribution ranges filtered according to

species-specific habitat preferences [12]. These habitat suitability models reflected species preferences for land cover types, their altitudinal limits, their tolerance to human presence, and their relationship with water bodies. The full-resolution mammal biodiversity data (representing all 5,291 terrestrial mammal species)[12] was rescaled to the study grid by summing the number of species' distributions that overlapped each grid cell; **2)** Domestic animal density. A number of past EID events with wildlife origin have emerged through farmed or domestic animal intermediate or amplifier hosts (e.g. Hendra and Nipah virus, SARS). Additionally, there is growing evidence that the global trend of intensification of livestock production increases the emergence risk of novel wildlife-origin zoonoses, e.g. Nipah virus in Malaysia [52], influenza viruses and others [6]. We used the Gridded Livestock of the World (GLW) dataset [53], which contains data for poultry, goat, buffalo, cattle, sheep and pig headcounts. We summed mammals to a single predictor (livestock mammal headcount) and retained poultry as a discrete predictor.

We analysed eight predictors from two datasets representing "Environmental" variables: **1)** Climate. Climatic factors have been repeatedly hypothesized as important in the global biogeography of human infectious diseases, including EIDs [32,54,55]. Climate may influence disease distribution through enhanced suitability for vectors of wildlife origin zoonoses (e.g. West Nile virus), more rapid vector reproduction rates and biting rates, changes in the efficiency or rates of pathogen transmission among hosts and vectors, and changes in the ability of pathogens to persist in the environment, among other factors [56,57]. Climate was represented by a single layer in our study, the Global Environmental Stratification [58], which uses a quantitative model to stratify the Earth's surface into zones of similar climate on a single

scalar measure, where higher values equate to warmer, wetter (more tropical) regions; **2)** <u>Land</u> <u>cover type</u>: Land cover type is associated with the distribution of terrestrial mammals [12] and other taxa [59], potentially exposing humans present to different assemblages of viral species. It is also likely that the types of contact between wildlife and people vary with land cover type. For land cover, we used the EarthEnv dataset [60], which divides the Earth's surface into twelve classes. These include different classes of natural ecosystems, urban land and cultivated vegetation (grouped with "Human Activity" above). We excluded barren areas, open water and snow/ice due to a lack of biologically plausible mechanisms for disease emergence. EarthEnv represents each class as a percentage per grid cell.

*Reporting effort*

The distribution of reported EID events is likely strongly influenced by an inconsistent spatial distribution of detection and reporting of disease outbreaks. Previous studies have used proxies of reporting effort such as the interpolated locations of known sampling sites ('sampling effort') [61]; frequency of countries of residence for all authors of all articles in the Journal of Infectious Disease ('reporting effort') [4]; and PubMed searches for keywords for each country ('reporting bias') [28]. Other studies have used occurrence records for a similar class of observations as a surrogate for background sampling effort; for example, in ecology, modelling the distribution of a particular species and utilizing occurrence records from multiple other species to represent background samples [62].

We adapted these approaches by deriving an index for reporting effort based on the spatial distribution of toponyms (place names) in peer-reviewed biomedical literature. We wrote a

Python package, *PubCrawler* (see Supplementary Methods: PubCrawler for full details), to

search the full text of each of the 1,266,085 (as of April 2016) articles in the PubMed Central

Open-Access Subset (PMCOAS) (http://www.ncbi.nlm.nih.gov/pmc/tools/ftp/) for toponyms

from the GeoNames database (http://geonames.org/), which includes data on population (if

appropriate), country, and geographical coordinates for each toponym. *PubCrawler* uses a set

of heuristics, based on textual and geographic features of the identified toponyms, to minimize

the number of false positives and select amongst ambiguous matches. We selected articles

matching terms from the Human Disease Ontology[63] and exported extracted toponyms. After

excluding a further round of potentially spurious matches, place name matches were assigned a

weight, normalized by article, and then summed to the study grid. To impute missing data

(resulting in a number of zero-value grid cells) and smooth noise in the raw output, we fit a

Poisson boosted regression tree model (using human population, accessibility, urbanized land,

DALY rates, health expenditure, and GDP as predictors), and used this to represent reporting

effort in our model. This approach produced a layer that adequately represented the underlying

data whilst achieving a similar coverage of grid cells to other layers.

Table 1: List of predictor layers included in the model.

| Variable | Unit per grid cell | Type | Source Dataset | Processing | Temporal Resolution |
|---|---|---|---|---|---|
| Human population | Population | Human activity | GRUMP | Rescaled | Decadal |
| Population change | Change in population | Human activity | GRUMP (calculated) | Calculated from rescaled layers | Decadal |
| Cropland | Proportion | Human activity | HYDE | Rescaled | Decadal |
| Cropland change | Change in proportion | Human activity | HYDE (calculated) | Calculated from rescaled layers | Decadal |
| Pasture | Proportion | Human activity | HYDE | Rescaled | Decadal |
| Pasture change | Change in proportion | Human activity | HYDE (calculated) | Calculated from rescaled layers | Decadal |
| Urban land | Percentage | Human activity | EarthEnv | Rescaled | Decadal |
| Managed/cultivated vegetation | Percentage | Human activity | EarthEnv | Rescaled | Static |
| Mammalian species richness | Count of species | Animals/hosts | Global Mammal Assessment | Reprojected, rescaled | Static |
| Domestic mammal headcount | Count of animals | Animals/hosts | GLW | Rescaled, summed buffalo, cattle, goat, pig, sheep headcounts | Static |
| Poultry headcount | Count of animals | Animals/hosts | GLW | Rescaled | Static |
| Global environmental stratification | Global environmental stratification | Environment | GEnS | Rescaled | Static |
| Evergreen/Deciduous Needleleaf Trees | Percentage | Environment | EarthEnv | Rescaled | Static |
| Evergreen Broadleaf Trees | Percentage | Environment | EarthEnv | Rescaled | Static |
| Deciduous Broadleaf Trees | Percentage | Environment | EarthEnv | Rescaled | Static |
| Mixed/Other Trees | Percentage | Environment | EarthEnv | Rescaled | Static |
| Shrubs | Percentage | Environment | EarthEnv | Rescaled | Static |
| Herbaceous Vegetation | Percentage | Environment | EarthEnv | Rescaled | Static |
| Regularly Flooded | Percentage | Environment | EarthEnv | Rescaled | Static |

| Vegetation Reporting Effort | Weighted number of mentions in publications | Observation bias | (Internal) | (See methods) | Static |
|---|---|---|---|---|---|

Table 2: Original resolutions and extents of source datasets.

| Source Dataset | Spatial Resolution | Temporal Resolution and Extent |
|---|---|---|
| GRUMP (Global Rural Urban Mapping Project)[43] | 0º5' | 5 years, 1970–2000 |
| HYDE (History Database of the Global Environment)[48] | 0º5' | 10 years, 1900–2000 |
| GMA (Global Mammal Assessment)[12] | 300m | N/A |
| GLW (Gridded Livestock of the World)[53] | 0.05º | N/A |
| GEnS (Global Environmental Stratification)[58] | 0º0'30" | N/A |
| EarthEnv[60] | 0º0'30" | N/A |

Statistical Framework

1. *Modelling algorithm:* We used boosted regression trees (BRT) to model EID occurrence [13,14,54] and to determine how conditions varied between locations where EID events have been observed compared to areas where they have not. BRTs handle non-linear relationships and higher order interactions among many variables more robustly than many other modelling methods, and are robust to monotonic transformations of data [13,14]. They fit potentially complex, non-linear relationships by aggregating the predictions of multiple simpler models, and are trained iteratively on random partitions of the data [13,14]. In addition, predictive accuracy of BRTs, as determined by common validation methodologies (e.g. Area Under the Curve of the Receiver-Operator Characteristic (AUC of the ROC), True Skill Statistic (TSS)), frequently exceeds conventional linear methods [13]. Unlike conventional models, they do not produce confidence intervals or p-values.

2. *Data and model fitting*: We used various resampling techniques to incorporate our measure of reporting effort[62,64], estimate the predictive power of our models, account for spatial uncertainty in EID events[17], and generate empirical confidence intervals for effects representing both sampling uncertainty and spatial uncertainty[65]. Each time an event was sampled, one presence point and one absence point were drawn (artificially fixing overall prevalence at 0.5)[17]; the presence point from the grid cells overlapped by that event's polygon and the absence point from all grid cells, both weighted by

reporting effort (the effect of weighting *presence* points by reporting effort made little difference for points with small, precisely-specified occurrence polygons, and for events with high uncertainty, acted as a prior specifying that, absent other knowledge, the event was more likely detected where reporting effort was higher).

All replicate BRT models were fit using the R packages *dismo* and *gbm*[13]. The function gbm.step() was called, with the parameters tree.complexity = 3 (governing interaction depth), learning.rate = 0.0035 (setting the "shrinkage" applied to individual trees), and n.trees = 35 (governing the initial number of trees fit, as well as the "step size" or number added at each step of the stagewise fitting process).[13] These values were selected through an iterative process, starting with the default parameters, adding tree complexity, and tuning the shrinkage and step size parameters to achieve successful gradient descent consistently across resampling runs, following [13] and [65]. With the final parameters, the BRTs composing the bootstrap model fit a mean of 1005 trees.

The sampling regimes were as follows:

1) A bootstrap resampling regime was used to fit 1000 replicate models. For each model, 147 events were drawn randomly with replacement from the set the 147 EID events of interest; for each selected event, one presence and one absence value were drawn as described above. The fitted models were used to generate Relative Influence boxplots and Partial Dependence plots with empirical 90% confidence intervals. The mean of the predictions of these models were used to generate all maps.

2) To compute validation statistics (described below), we conducted 100 rounds of 10-fold cross-validation[17,65]. In each round, a single presence and absence sample were drawn for each event, which were assigned randomly to ten groups. Each group in turn was held out, and a model was trained on the remaining groups' samples. The model's predictions for the presence and absences samples of the held-out group were used to construct confusion matrices, and calculate the AUC and TSS. This process was repeated 100 times, and the median, 0.05 and 0.95 quantiles for all scores were reported. The entire process was conducted for each AUC and TSS reported.

3. *Factoring reporting bias out:* We assumed that the distribution of *observed* EID events was conditional on the distribution of reporting effort across the globe following [62].

We fit our main model weighting by reporting effort. The models produce a response relative to this. We multiplied this response by the value of reporting effort in each grid cell to map an index of observed EID event risk. We produced an estimate of the risk index factoring out reporting *bias* as follows:

We assumed that the optimal distribution of reporting effort for human disease events in a location is proportional to the distribution of the human population. In reality, other unmeasured factors likely affect this. If we take this assumption, we can define *reporting bias* as proportional to the ratio of reporting effort to the human population.

$$reporting\ bias \propto \frac{reporting\ effort}{population}$$

When bias is known, it is possible to estimate the true distribution of a phenomenon by "factoring bias out" [62]. In ecological studies, this generally means dividing by the measured "survey effort", assuming that the optimal distribution of search effort is uniform across the landscape.

$$true\ risk\ index\ \propto \frac{observed\ risk\ index}{reporting\ bias}$$

We posit that, in the case of human disease events, uniform search effort across a landscape is also suboptimal, and that it is safer to assume optimal reporting effort distribution would be proportional to the human population. In this case, we remove "bias" by factoring **out** measured effort and factoring *in* assumed optimal effort, and obtain a hypothetical map of the true event risk index, thus:

$$true\ risk\ index\ \propto\ observed\ risk\ index \times \frac{human\ population}{reporting\ effort}$$

*Validation and model performance:* We used multiple tools for model validation and performance. For our bootstrap model, we calculated deviance explained using the gbm.step() function [13] and also derived median and empirical 90% confidence intervals by taking the 0.05, 0.5, and 0.95 quantiles of those values for the replicate models. Since this model is fit relative to reporting effort, percentage deviance explained is calculated relative to that variable. For the ten-fold cross-validation runs, we calculated the area under the receiver-operating characteristic curve (AUC), a threshold-independent measure of model predictive performance that is commonly used as a validation metric in species distribution modelling [66]. The AUC can be interpreted as "the probability that the model will rank a randomly chosen presence site higher than a randomly chosen absence site" [67], or more accurately in our application, a measure of a model's performance to discriminate EID events from random points [62]. Because the use of AUC has been criticized for its lack of sensitivity to absolute predicted probability and its inclusion of *a priori* untenable prediction thresholds [15], we also calculated the True Skill Statistic (TSS)[17].

Because all test statistics and figures from our main model are relative to the reporting effort measure, we also ran "unweighted" models. We expected these would score yield higher cross-validation scores, since we expected that reporting effort would be correlated both with some important predictor variables and the outcome, and weighting background samples uniformly rather than according to this variable would present a clearer contrast. To avoid bias from land area in the WGS84 grid cells, we additionally weighted our "unweighted models" by land area per grid cell. The figures from these models are presented fully in SI.

**Acknowledgements**

**References**

1       Heymann, D. L. *et al.* Global health security: the wider lessons from the west African Ebola virus disease epidemic. *Lancet* **385**, 1884-1901 (2015).
2       Morens, D. M. & Fauci, A. S. Emerging Infectious Diseases in 2012: 20 Years after the Institute of Medicine Report. *Mbio* **3** (2012).
3       Pike, J., Bogich, T. L., Elwood, S., Finnoff, D. C. & Daszak, P. Economic optimization of a global stategy to reduce the pandemic threat. *Proceedings of the National Academy of Sciences, USA* **111**, 18519-18523 (2014).
4       Jones, K. E. *et al.* Global trends in emerging infectious diseases. *Nature* **451**, 990-993, doi:10.1038/nature06536 (2008).
5       Wolfe, N. D., Dunavan, C. P. & Diamond, J. Origins of major human infectious diseases. *Nature* **447**, 279-283 (2007).
6       Jones, B. A. *et al.* Zoonosis emergence linked to agricultural intensification and environmental change. *Proceedings of the National Academy of Sciences* **110**, 8399-8404 (2013).
7       Karesh, W. B. *et al.* Zoonoses 1 Ecology of zoonoses: natural and unnatural histories. *Lancet* **380**, 1936-1945 (2012).
8       Morse, S. Factors in the Emergence of Infectious Diseases. *Emerging infectious diseases* **1**, 7-15 (1995).
9       Coker, R. *et al.* Towards a conceptual framework to support one-health research for policy on emerging zoonoses. *The Lancet infectious diseases* **11**, 326-331 (2011).
10      Woolhouse, M., Scott, F., Hudson, Z., Howey, R. & Chase-Topping, M. Human viruses: discovery and emergence. *Philosophical Transactions of the Royal Society B-Biological Sciences* **367**, 2864-2871 (2012).
11      Brierley, L., Vonhof, M. J., Olival, K. J., Daszak, P. & Jones, K. E. Quantifying Global Drivers of Zoonotic Bat Viruses: A Process-Based Perspective. *The American Naturalist* **187**, E53-E64, doi:doi:10.1086/684391 (2016).
12      Rondinini, C. *et al.* Global habitat suitability models of terrestrial mammals. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **366**, 2633-2641, doi:10.1098/rstb.2011.0113 (2011).
13      Elith, J., Leathwick, J. R. & Hastie, T. A working guide to boosted regression trees. *The Journal of animal ecology* **77**, 802-813 (2008).
14      De'ath, G. Boosted trees for ecological modeling and prediction. *Ecology* **88**, 243-251 (2007).

15    Lobo, J. M., Jiménez-Valverde, A. & Real, R. AUC: a misleading measure of the performance of predictive distribution models. *Global Ecology and Biogeography* **17**, 145-151, doi:10.1111/j.1466-8238.2007.00358.x (2008).

16    Allouche, O., Tsoar, A. & Kadmon, R. Assessing the accuracy of species distribution models: prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology* **43**, 1223-1232, doi:10.1111/j.1365-2664.2006.01214.x (2006).

17    Barbet-Massin, M., Jiguet, F., Albert, C. H. & Thuiller, W. Selecting pseudo-absences for species distribution models: how, where and how many? *Methods in Ecology and Evolution* **3**, 327-338, doi:10.1111/j.2041-210X.2011.00172.x (2012).

18    ArcGIS v. 10.2.2 (ESRI, Redlands, CA, 2016).

19    Weiss, R. A. & McMichael, A. J. Social and environmental risk factors in the emergence of infectious diseases. *Nature Medicine* **10**, S70-S76 (2004).

20    McFarlane, R., Sleigh, A. & McMichael, A. Land-Use Change and Emerging Infectious Disease on an Island Continent. *International Journal of Environmental Research and Public Health* **10**, 2699-2719 (2013).

21    Patz, J. A. *et al.* Unhealthy landscapes: Policy recommendations on land use change and infectious disease emergence. *Environmental Health Perspectives* **112**, 1092-1098 (2004).

22    Keesing, F. *et al.* Impacts of biodiversity on the emergence and transmission of infectious diseases. *Nature* **468**, 647-652 (2010).

23    Myers, N., Mittermeier, R. A., Mittermeier, C. G., da Fonseca, G. A. B. & Kent, J. Biodiversity hotspots for conservation priorities. *Nature* **403**, 853-858, doi:http://www.nature.com/nature/journal/v403/n6772/suppinfo/403853a0_S1.html (2000).

24    Murray, K. A. & Daszak, P. Human ecology in pathogenic landscapes: two hypotheses on how land use change drives viral emergence. *Current Opinion in Virology* **3**, 79-83, doi:10.1016/j.coviro.2013.01.006 (2013).

25    Schmidt, K. A. & Ostfeld, R. S. Biodiversity and the dilution effect in disease ecology. *Ecology* **82**, 609-619 (2001).

26    Salkeld, D. J., Padgett, K. A. & Jones, J. H. A meta-analysis suggesting that the relationship between biodiversity and risk of zoonotic pathogen transmission is idiosyncratic. *Ecology Letters* **16**, 679-686, doi:10.1111/ele.12101 (2013).

27    Randolph , S. E. & Dobson , A. D. M. Pangloss revisited: a critique of the dilution effect and the biodiversity-buffers-disease paradigm. *Parasitology* **FirstView**, 1-17, doi:doi:10.1017/S0031182012000200 (2012).

28    Yang, K. *et al.* Global Distribution of Outbreaks of Water-Associated Infectious Diseases. *PLoS neglected tropical diseases* **6**, e1483-e1483 (2012).

29    Loh, E. H. *et al.* Targeting Transmission Pathways for Emerging Zoonotic Disease Surveillance and Control. *Vector Borne Zoonotic Dis* **15**, 432-437, doi:10.1089/vbz.2013.1563 (2015).

30    Morse, S. S. *et al.* Prediction and prevention of the next pandemic zoonosis. *Lancet* **380**, 1956-1965 (2012).

31    Olson, S. H. *et al.* Drivers of Emerging Infectious Disease Events as a Framework for Digital Detection. *Emerg Infect Dis* **21**, 1285-1292, doi:10.3201/eid2108.141156 (2015).

32    Murray, K. A. *et al.* Global biogeography of human infectious diseases. *Proceedings of the National Academy of Sciences* **112**, 12746-12751, doi:10.1073/pnas.1507442112 (2015).

33    Gold, Z. *et al.* The Emerging Infectious Disease Repository (EIDR): A novel resource for investigating the emergence of infectious disease. *Ecohealth* (In Review).

34    Funk, S., Bogich, T. L., Jones, K. E., Kilpatrick, A. M. & Daszak, P. Quantifying trends in disease impact to produce a consistent and reproducible definition of an emerging infectious disease. *PLoS ONE* **8**, e69951 (2013).

35    Carroll, D. *et al.* The Global Virome Project. *Science* (In Review).

36    Laurance, W. F. *et al.* A global strategy for road building. *Nature* **513**, 229-+, doi:10.1038/nature13717 (2014).

37    Loh, E. H., Murray, K. A., Nava, A., Aguirre, A. A. & Daszak, P. in *Tropical Conservation: Perspectives on Local and Global Priorities* (eds A. Alonso Aguirre & Raman Sukumar) Ch. 6, 79-88 (Oxford. University Press, 2016).

38    Anthony, S. J. *et al.* A Strategy To Estimate Unknown Viral Diversity in Mammals. *mBio* **4**, e00598-00513-e00598-00513 (2013).

39    ecohealthalliance/hotspots2: "Global Correlates" paper (2016).

40    ecohealthalliance/pubcrawler: "Global correlates" paper (2016).

41    Moffett, A., Shackelford, N. & Sarkar, S. Malaria in Africa: vector species' niche models and relative risk maps. *PloS one* **2**, e824-e824 (2007).

42    McCallum, H. How should pathogen transmission be modelled? *Trends in Ecology & Evolution* **16**, 295-300, doi:10.1016/S0169-5347(01)02144-9 (2001).

43    Socioeconomic Data and Applications Center (sedac). *Global Rural-Urban Mapping Project (GRUMP), v1*, <http://sedac.ciesin.columbia.edu/data/collection/grump-v1> (2015).

44    Ostfeld, R. S. & Keesing, F. Biodiversity series: the function of biodiversity in the ecology of vector-borne zoonotic diseases. *Canadian Journal of Zoology* **78**, 2061-2078 (2000).

45    Keesing, F. *et al.* Impacts of biodiversity on the emergence and transmission of infectious diseases. *Nature* **468**, 647-652 (2010).

46    Ostfeld, R. S. & Keesing, F. Effects of Host Diversity on Infectious Disease. *Annual Review of Ecology, Evolution, and Systematics* **43**, 157-182, doi:doi:10.1146/annurev-ecolsys-102710-145022 (2012).

47    Bogich, T. L. *et al.* Preventing Pandemics Via International Development: A Systems Approach. *PLoS Medicine* **9**, doi:10.1371/journal.pmed.1001354 (2012).

48    Klein Goldewijk, K., Beusen, A., Van Drecht, G. & De Vos, M. The HYDE 3.1 spatially explicit database of human-induced global land-use change over the past 12,000 years. *Global Ecology and Biogeography* **20**, 73-86, doi:10.1111/j.1466-8238.2010.00587.x (2011).

49    Lloyd-Smith, J. O. *et al.* Epidemic dynamics at the human-animal interface. *Science* **326**, 1362-1367, doi:10.1126/science.1177345 (2009).

50    Dunn, R. R., Davies, T. J., Harris, N. C. & Gavin, M. C. Global drivers of human pathogen richness and prevalence. *Proceedings of the Royal Society B-Biological Sciences* **277**, 2587-2595, doi:10.1098/rspb.2010.0340 (2010).

51    Woolhouse, M. E. J. & Gowtage-Sequeria, S. Host Range and Emerging and Reemerging Pathogens. *Emerging infectious diseases* **11**, 1842-1847 (2005).

52    Pulliam, J. R. C. *et al.* Agricultural intensification, priming for persistence and the emergence of Nipah virus: a lethal bat-borne zoonosis. *Journal of The Royal Society Interface* **9**, 89-101 (2011).

53    Robinson, T. P. *et al.* Mapping the global distribution of livestock. *PloS one* **9**, e96084-e96084, doi:10.1371/journal.pone.0096084 (2014).

54    Hay, S. I. *et al.* Global mapping of infectious disease. *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **368**, 20120250-20120250, doi:10.1098/rstb.2012.0250 (2013).

55    Guernier, V., Hochberg, M. E. & Guégan, J.-F. Ecology Drives the Worldwide Distribution of Human Diseases. *PLoS Biology* **2**, e141, doi:10.1371/journal.pbio.0020141 (2004).

56    Rohr, J. R. *et al.* Frontiers in climate change-disease research. *Trends in Ecology & Evolution* **26**, 270-277, doi:10.1016/j.tree.2011.03.002 (2011).

57    Kilpatrick, A. M. & Randolph, S. E. Zoonoses 2 Drivers, dynamics, and control of emerging vector-borne zoonotic diseases. *Lancet* **380**, 1946-1955 (2012).

58    Metzger, M. J. *et al.* A high-resolution bioclimate map of the world: a unifying framework for global biodiversity research and monitoring. *Global Ecology and Biogeography* **22**, 630-638, doi:10.1111/geb.12022 (2013).

59    Jenkins, C. N., Pimm, S. L. & Joppa, L. N. Global patterns of terrestrial vertebrate diversity and conservation. *Proceedings of the National Academy of Sciences of the United States of America* **110**, E2602-2610 (2013).

60    Tuanmu, M.-N. & Jetz, W. A global 1-km consensus land-cover product for biodiversity and ecosystem modelling. *Global Ecology and Biogeography* **23**, n/a-n/a, doi:10.1111/geb.12182 (2014).

61    Hopkins, M. E. & Nunn, C. L. A global gap analysis of infectious agents in wild primates. *Diversity and Distributions* **13**, 561-572, doi:10.1111/j.1472-4642.2007.00364.x (2007).

62    Phillips, S. J. *et al.* Sample selection bias and presence-only distribution models: implications for background and pseudo-absence data. *Ecological applications : a publication of the Ecological Society of America* **19**, 181-197 (2009).

63    Kibbe, W. A. *et al.* Disease Ontology 2015 update: an expanded and updated database of human diseases for linking biomedical knowledge through disease data. *Nucleic Acids Res* **43**, D1071-1078, doi:10.1093/nar/gku1011 (2015).

64    Dorazio, R. M. Accounting for imperfect detection and survey bias in statistical analysis of presence-only data. *Global Ecology and Biogeography* **23**, 1472-1484, doi:10.1111/geb.12216 (2014).

65    Leathwick, J. R., Elith, J., Francis, M. P., Hastie, T. & Taylor, P. Variation in demersal fish species richness in the oceans surrounding New Zealand : an analysis using boosted regression trees.  **321**, 267-281 (2006).

66    Liu, C., White, M. & Newell, G. Measuring and comparing the accuracy of species distribution models with presence-absence data. *Ecography* **34**, 232-243, doi:10.1111/j.1600-0587.2010.06354.x (2011).

67    Fawcett, T. An introduction to ROC analysis. *Pattern Recognition Letters* **27**, 861-874, doi:10.1016/j.patrec.2005.10.010 (2006).