Trade-offs between individual and ensemble forecasts of an emerging infectious disease

Rachel J. Oidtman^{1,2,3,*}, Elisa Omodei², Moritz U. G. Kraemer^{4,5,6}, Carlos A. Castañeda-Orjuela⁷, Erica Cruz-Rivera⁷, Sandra Misnaza-Castrillón⁷, Myriam Patricia Cifuentes⁸, Luz Emilse Rincon⁸, Viviana Cañon⁹, Pedro de Alarcon¹⁰, Guido España¹, John H. Huber¹, Sarah C. Hill^{4,11}, Christopher M. Barker¹², Michael A. Johansson¹³, Carrie A. Manore¹⁴, Robert C. Reiner, Jr.¹⁵, Isabel Rodriguez-Barraquer¹⁶, Amir S. Siraj¹, Enrique Frias-Martinez¹⁷, Manuel García-Herranz^{2,*}, and T. Alex Perkins^{1,*} ¹Department of Biological Sciences and Eck Institute for Global Health, University of Notre Dame, Notre Dame, IN, USA ² UNICEF, New York, NY, USA ³Department of Ecology and Evolution, University of Chicago, Chicago, IL, USA ⁴Department of Zoology. University of Oxford. Oxford. UK ⁵Boston Children's Hospital, Boston, MA, USA ⁶Harvard Medical School, Boston, MA, USA ⁷Instituto Nacional de Salud, Bogotá, Colombia ⁸Ministerio de Salud y Protección Social, Bogotá, Colombia ⁹ UNICEF, Bogotá, Colombia ¹⁰LUCA Telefonica Data Unit, Madrid, Spain ¹¹Department of Pathobiology and Population Sciences, The Royal Veterinary College, London, UK ¹²Department of Pathology, Microbiology, and Immunology, School of Veterinary Medicince, University of California, Davis, CA, USA ¹³Division of Vector-Borne Diseases, Centers for Disease Control and Prevention, San Juan, Puerto Rico ¹⁴Information Systems and Modeling (A-1), Los Alamos National Laboratory, Los Alamos, NM, USA ¹⁵Institute for Health Metrics and Evaluation, University of Washington, Seattle, WA, USA¹⁶Department of Medicine, University of California, San Francisco, CA, USA ¹⁷ Telefonica Research, Madrid, Spain *Corresponding authors: rjoidtman@gmail.com, mgarciaherranz@unicef.org, taperkins@nd.edu

June 30, 2021

Abstract

NOTE: This prephilipper and the probability of the spread, some of which can be addressed with probabilistic forecasts. The many uncertainties about the epidemiology of emerging pathogens can make

> it difficult to choose among model structures and assumptions, however. To assess the potential for uncertainties about emerging pathogens to affect forecasts of their spread, we evaluated the performance of a suite of 16 forecasting models in the context of the 2015-2016 Zika epidemic in Colombia. Each model featured a different combination of assumptions about the role of human mobility in driving transmission, spatiotemporal variation in transmission potential, and the number of times the virus was introduced. All models used the same core transmission model and the same iterative data assimilation algorithm to generate forecasts. By assessing forecast performance through time using logarithmic scoring with ensemble weighting, we found that which model assumptions had the most ensemble weight changed through time. In particular, spatially coupled models had higher ensemble weights in the early and late phases of the epidemic, whereas non-spatial models had higher ensemble weights at the peak of the epidemic. We compared forecast performance of the equally weighted ensemble model to each individual model and identified a trade-off whereby certain individual models outperformed the ensemble model early in the epidemic but the ensemble model outperformed all individual models on average. On balance, our results suggest that suites of models that span uncertainty across alternative assumptions are necessary to obtain robust forecasts in the context of emerging infectious diseases.

1 Introduction

Pathogen emergence, or the phenomenon of novel or established pathogens invad-2 ing a new host population, has been occurring more frequently in recent decades 3 [1]. In the last 40 years, more than 150 pathogens of humans have been identified 4 as emerging or re-emerging [2, 3]. In these situations, host populations are largely 5 susceptible, which can result in dynamics ranging from self-limiting outbreaks, as 6 with Lassa virus [4], to sustained pandemics, as with HIV [5], depending on the 7 pathogen's traits and the context in which it emerges. When emergence does oc-8 cur, mathematical models can be helpful for anticipating the future course of the 9 pathogen's spread [6, 7, 8]. 10

A necessary part of using models to forecast emerging pathogens is making deci-11 sions about how to handle the many uncertainties associated with these unfamiliar 12 microbes [8]. Given the biological and ecological diversity of emerging pathogens, 13 there is often considerable uncertainty about various aspects of their natural his-14 tories, such as their potential for superspreading [9], the role of human mobility in 15 their spatial spread [10, 11], drivers of spatiotemporal variation in their transmission 16 [6, 12], and even their modes of transmission [13]. In the case of MERS-CoV, for ex-17 ample, it took years to determine that the primary transmission route was spillover 18 from camels rather than sustained human-to-human transmission [14]. A lack of 19 definitive understanding about such basic aspects of natural history represents a 20 major challenge for forecasting emerging pathogens. 21

Inevitably, different forecasters make diverse choices about how to address un-22 known aspects of an emerging pathogen's natural history, as they do for numerous 23 model features. This diversity of approaches has itself been viewed as part of the 24 solution to the problem of model uncertainty, based on the idea that the biases of 25 different models might counteract one another to produce a reliable forecast when 26 viewed from the perspective of an ensemble of models [15]. This idea has support in 27 multi-model efforts to forecast seasonal transmission of endemic pathogens, such as 28 influenza and dengue viruses [16, 17, 18, 19, 20], with ensemble forecasts routinely 29 outperforming individual models. These successes with endemic pathogens have 30 motivated multi-model approaches in response to several emerging pathogens, in-31 cluding forecasting challenges for chikungunya [21] and Ebola [22], vaccine trial site 32 selection for Zika [23], and a multi-model decision-making framework for COVID-19 33 [15, 24].34

Although there has been increased attention to multi-model forecasting of emerg-35 ing pathogens in the last few years, these initiatives have involved significant effort to 36 coordinate forecasts among multiple modeling groups [25, 26]. Coordination across 37 multiple groups has clear potential to add value beyond what any single modeling 38 group can offer alone. At the same time, using multiple models to hedge against 39 uncertainties about a pathogen's natural history could potentially improve forecasts 40 from a single modeling group, too [16, 18]. This could, in turn, improve ensemble 41 forecasts based on contributions from multiple modeling groups. An ensemble-based 42 approach by one modeling group that contributes to forecasts of seasonal influenza 43 in the United States demonstrates the success that a single modeling group can 44 achieve with an ensemble-based approach [27], and that such an ensemble can con-45 tribute value to an ensemble of forecasts from multiple modeling groups [18]. Similar 46 approaches have not been widely adopted for forecasting emerging pathogens by a 47 single modeling group (although see [28]), despite the heightened uncertainty inher-48 ent to emerging pathogens. 49

Here, we evaluate the potential for an ensemble of models that span uncertain-50 ties in pathogen natural history, but share a common core structure, to accurately 51 forecast the dynamics of an emerging pathogen. We do so in the context of the 52 2015-2016 Zika epidemic in Colombia, which was well-characterized epidemiologi-53 cally (Fig. 1) [29] and involved potentially consequential uncertainties about: i) 54 the role of human mobility in facilitating spread across the country [30], ii) the 55 relationship between environmental conditions and transmission of this mosquito-56 borne virus [6, 12], and iii) the number of times the virus was introduced into the 57 country [31]. In this retrospective analysis, we used data assimilation to update 16 58 distinct models throughout the epidemic period and assessed forecast performance 59 of all models relative to an equally weighted ensemble model. This allowed us to 60 quantify the contribution of variants of each of the three aforementioned uncertain-61 ties to model performance during different phases of the epidemic. In doing so, we 62 sought to not only assess the performance of the ensemble model relative to indi-63 vidual models, but also to learn about features of individual models that may be 64 associated with improved forecast accuracy over the course of an epidemic. 65



Figure 1: Temporal and spatial variation of Zika incidence, temperature, and mosquito occurrence probability in Colombia. a. Weekly Zika incidence from August 9, 2015 to October 1, 2016 with all 31 mainland departments approximately ordered from south to north. b. Points indicate average temperature data and lines indicate temperature by department. c. Points indicate average mosquito occurrence probability and lines indicate mosquito occurrence probability by department. d-f. Mobility matrices under three different assumptions of mobility, with departments ordered south to north on y-axis and north south on x-axis. Tan indicates high rates of mobility, dark purple indicates low rates of mobility, white indicates no movements.

66 2 Results

67 2.1 General forecast performance

Before any data assimilation had occurred, our 16 models (See Table 1) initially 68 forecasted very low incidence across most departments over the 60-week period of 69 our analysis (Figs. 2 top row, S12). Even so, short-term forecasts over a four-week 70 horizon were consistent with the still-low observed incidence at that time (Figs. S3) 71 purple, S18). By the time twelve weeks of data had been assimilated into the mod-72 els, forecasts over the 60-week period of our analysis were considerably higher than 73 the initial forecasts and better aligned with the observed trajectory of the epidemic 74 (Figs. 2 second row, S13). Over those first twelve weeks, model parameters changed 75 modestly (Fig. S6) and correlations among parameters began to emerge (Figs. S7, 76 S8, S9, S10). We observed a more substantial change in the proportion of individ-77 ual stochastic realizations (where the n^{th} stochastic realization is the n^{th} "particle" 78 generated from some set of parameters $\vec{\theta}_{t,n}$ at time t) resulting in an epidemic, with 79 those particles resulting in no epidemic being filtered out almost entirely by week 80 12 (Fig. S1). Because each particle retained its stochastic realization of past in-81 cidence across successive data assimilation periods, stochastic realizations of past 82 incidence were inherited by particles much like parameter values. By week 24, many 83 of the models correctly recognized that they were at or near the epidemic's peak 84 and forecasted a downward trajectory for the remainder of the 60-week period of our 85 analysis (Figs. 2 third row, S27). The particle filtering algorithm replaced nearly 86 half of the original particles by that point (Fig. S_2), with the new particles con-87 sisting of stochastic realizations of past incidence selected through data assimilation 88 and updated every four weeks with forward simulations based on either original or 89 new parameter combinations. As the end of the 60-week period of our analysis was 90 approached, parameter correlations continued to strengthen (Figs. S7, S8, S9, S10), 91 our estimate of the reporting probability increased (Fig. S6), and only around 20% 92 of the original particles remained (Fig. S1). 93

Table 1: Different model assumptions regarding the role of human mobility in facilitating pathogen spread across the country, the relationship between environmental conditions and transmission of ZIKV, and the number of times the virus was introduced into Colombia. The suite of 16 models reflected factorial combinations of these three assumptions.

Human mobility	Transmission potential	Number of ZIKV introductions
CDR-informed	Fixed R [6]	One
Gravity model	Dynamic R [12]	Two
Radiation model		
No human mobility		

⁹⁴ 2.2 Model-specific forecast performance

To quantify the forecast performance of individual models over time, we used logarithmic scoring (hereafter, log scoring) to compare forecasts of cumulative incidence four weeks into the future to observed values at departmental and national levels. We assessed log scores once the first case was reported nationally for spatially coupled models (i.e., models with explicit human mobility), and once the first case was



Figure 2: Observed incidence (navy points) with the median forecast for 16 models (black lines) with the equally weighted ensemble model (green band) for Antioquia, Norte de Santander, Cauca, and Amazonas at five points throughout the epidemic. Plotted departments reflect differences in population, epidemic size, and geographic regions of Colombia and are represented by each column. The vertical pink line indicates the point at which the forecast was made (also labeled on the right axis), with data to the left of the line assimilated into the model fit. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change. The green band reflects the 50% credible interval of the equally weighted ensemble model.

reported in each department for non-spatial models (i.e., models with no explicit 100 human mobility). Log scores were generally high for spatially coupled models early 101 in the epidemic, given that observed cases and forecasts were both low at that time 102 (Fig. S18, a-c). By week 12, as cases were reported in more departments, the accu-103 racy of forecasts from non-spatial models improved (Fig. S18 d onward). Forecast 104 performance around the peak of the epidemic differed considerably across models 105 and departments, with forecasts from non-spatial models being somewhat lower 106 than observed incidence and forecasts from spatially coupled models being some-107 what higher (Fig. S14, Fig. S18 f-j). Around the peak of the epidemic, forecasts 108 from spatially coupled models generally had higher log scores in departments with 109 lower incidence (e.g., Nariño). Later in the epidemic (weeks 40-56), some models 110 continued to forecast higher incidence than observed in some departments, despite 111 having passed the peak incidence of reported cases (Fig. S16). In particular, models 112 that used the dynamic instead of the static formulation of the reproduction number 113 (i.e., the temporal relationship between R and environmental drivers is dynamic 114 instead of static) were more susceptible to this behavior (note lower log scores in 115 "Rt" versus "R" models in Fig. S18 k-o), given that their forecasts were sensitive to 116 seasonal changes in temperature and mosquito occurrence. 117

Next, we used these log scores in an expectation-maximization (EM) optimiza-118 tion algorithm [32] to identify an optimal weighting of retrospective model-specific 119 forecasts into an ensemble forecast (Fig. S25-S29) in each forecasting period (Fig. 120 S17). To learn how model assumptions affected the inclusion of different models into 121 the optimally weighted ensemble for each forecasting period, we summed and then 122 normalized models' ensemble weights across each class of assumption (Fig. 3). Over 123 the course of the epidemic, changes in weighting for the assumptions about human 124 mobility and spatiotemporal variation in transmission, but not about the number 125 of virus introductions into the country, closely followed patterns in the trajectory 126 of the national epidemic. Spatially coupled models had most or all of the weight 127 in the early and late stages of the epidemic, while non-spatial models had most of 128 the weight around the peak of the epidemic (Fig. 3 b). Although the non-spatial 129 models somewhat under-predicted incidence in the middle stages of the epidemic, 130 this was often to a lesser extent than the spatially coupled models' over-predictions 131 of incidence (Fig. S_3). As a result, the EM algorithm achieved a balance between 132 the over- and under-predictions of these different models. 133

The maximum ensemble weight in any forecasting period was 0.802, held by 134 one model with a static R, two ZIKV introductions into the country, and CDR-135 informed human mobility 12 weeks after the first reported Zika case (Fig. S17). 136 Combined, the two models with static R and CDR-informed human mobility data 137 had the most instances of a non-zero ensemble weight (Fig. S17), occurring in 13 138 of 15 assimilation periods, with an average weight of 0.18. Around the peak of 139 the epidemic, non-spatial models had the highest ensemble weight, reflecting the 140 accuracy of short-term forecasts in some departments (e.g., Magdalena and Vaupés) 141 and their overall accuracy in nationally-aggregated forecasts (Fig. S11). Near the 142 end of the epidemic, the ensemble weight for models with a static R (Fig. 3 c) 143 increased as their forecasts more closely matched the downturn of incidence later 144 in the epidemic relative to models with dynamic R (Fig. S20). This was likely the 145 result of mosquito occurrence probability and temperature becoming more favorable 146 for transmission in many departments later in the epidemic (Fig. S21-S22), causing 147 the dynamic R models to forecast a late resurgence in Zika incidence. 148



Figure 3: Ensemble weight partitioned across assumptions about the role of human mobility in driving transmission, drivers of spatiotemporal variation in R, and the number of ZIKV introductions. a: Weekly Zika incidence aggregated to the national scale. b-d: Proportion of ensemble weight across assumption type colored by explicit assumption.

¹⁴⁹ 2.3 Target-oriented forecast performance

Short-term changes in incidence are an important target of infectious disease fore-150 casting, but there are other targets of potentially greater significance to public health 151 decision making. To explore these, we evaluated the ability of the 16 models—and an 152 evenly weighted ensemble—to forecast three targets at the department level: peak 153 incidence, week of peak incidence, and onset week, which we defined as the week 154 by which ten cases were first reported. We evaluated models based on log scores 155 of these targets. Summing log scores across departments to allow for comparisons 156 across different forecasting periods (Fig. 4), we found that, on average, the en-157 semble model outperformed every individual model for all three forecasting targets 158 (indicated by the ensemble model's location on the y-axis). Early in the epidemic, 159 spatially coupled models with a static R performed only slightly better (up to 1%) 160 than the equally weighted ensemble (Fig. 4). For the remainder of the epidemic, 161 the equally weighted ensemble model outperformed all individual models (Fig. 4). 162 Such small changes in forecast performance when averaging over space shows that 163 differences in forecast performance across space dominate relative to those across 164 time. 165

By summing log scores across forecasting periods to allow for comparisons across 166 departments (Fig. 5), we found that some individual models outperformed the 167 ensemble model in forecasting the peak incidence and the week of peak incidence. 168 In departments on the Caribbean Coast that experienced intermediate epidemic 169 sizes (e.g., Antioquia, Sucre, Atlántico), spatially coupled models with a static R170 outperformed the ensemble model at forecasting the peak week by about 10% (Fig. 171 5 A). At those same locations, the equally weighted ensemble performed better than 172 or similar to those same models at forecasting peak incidence and onset week (Fig. 5) 173 b-c). Over forecasting periods and departments, the non-spatial models consistently 174 had lower average forecast scores than the spatially coupled models (indicated by 175 their location on the y-axis in Figs. 4-5). This trend appeared because initial 176 forecasts from non-spatial models were not updated until the first case appeared in 177 each department, while initial forecasts from spatially coupled models were updated 178 when the first case appeared in the country. 179

180 3 Discussion

We assessed the potential for a suite of individual models that span a range of un-181 certainties, and ensembles of these models, to accurately forecast the dynamics of an 182 emerging pathogen. Results from the general forecast performance analysis demon-183 strated that once we began assimilating data into models, forecasts rapidly became 184 more accurate. Models were initialized with a wide range of parameter values [33], 185 with many initial parameter combinations producing unrealistic forecast trajectories, 186 but after only four assimilation periods (12 weeks), nearly 100% of those parameters 187 that produced zero infections were dropped. Similar to other retrospective forecast 188 analyses [16, 34], as more data were assimilated into the models over time, the model 189 fits and forecasts generally became more closely aligned with temporal trends in the 190 data. This was because the particle filter allowed model parameters to continually 191 adapt to noisy data [35]. There were still some exceptions where the particle filter 192 could not fully compensate for shortcomings of the transmission model, such as the 193 drastic underestimates of incidence in departments with sub-optimal conditions for 194 transmission (e.g., static R model in Risaralda in Fig. S20). At the same time, 195



Figure 4: Model-specific forecast scores relative to equally weighted ensemble model for each assimilation period and forecasting target. *a.* Timing of peak week (within two weeks). *b.* Incidence at peak week. *c.* Onset week. Forecast scores are averaged over department. Models are ordered on the y-axis by average forecast score for each forecasting target. Model names on the y-axis are abbreviated such that "R" or "Rt" indicates assumption about spatiotemporal variation, "1" or "2" indicates number of introduction events, and "CDRs", "gravity", "radiation" or "nonspatial" indicates the human mobility assumption. In the heat plot, blue indicates individual model performed better than the ensemble model in a given department, red indicates individual model performed roughly the same as the ensemble model.



Departments ordered by overall incidence (high to low)

Figure 5: Model-specific forecast scores relative to equally weighted ensemble model for each department and forecasting target. *a.* Timing of peak week (within two weeks). *b.* Incidence at peak week. *c.* Onset week, or the week by which ten cumulative cases occurred. Forecast scores are averaged over department. Models are ordered on the y-axis by average forecast score for each forecasting target, with model names abbreviated in the same manner as Fig. 4. Departments are ordered on the x-axis from high to low for overall incidence. In the heat plot, blue indicates individual model performed better than the ensemble model in a given department, red indicates individual model performed roughly the same as the ensemble model.

the broader suite of models buffered against shortcomings of any single transmissionmodel.

In the model-specific forecast performance analysis, we identified clear temporal 198 trends related to when models with a static R versus a dynamic R should be in-199 cluded in an optimally weighted ensemble. In contrast, there were no clear temporal 200 trends in weighting regarding the assumption about the number of times the virus 201 was introduced into the country, potentially reflecting that, even with multiple in-202 troductions, most transmission may have been linked to a single introduction [31]. 203 Models with a dynamic R had higher weights in the ensemble at the peak of the 204 epidemic, while models with a static R had higher weights at the beginning and 205 end of the epidemic. This was likely due to temporal shifts in temperature and 206 mosquito occurrence probabilities dominating forecasts of transmission potential for 207 the models with a dynamic R. For example, in the latter parts of the epidemic 208 when reported cases were declining, mosquito conditions and temperature became 209 more suitable for transmission in many departments. This caused models with a dy-210 namic R to forecast a resurgence in ZIKV transmission in those departments, while 211 models with a static R forecasted a downturn in incidence that was more similar 212 to the observed dynamics. This finding that susceptible depletion may have been 213 more influential than temporal variation in environmental conditions for the Zika 214 epidemic is consistent with recent findings for SARS-CoV-2 [36]. 215

Through the model-specific forecast performance analysis, we also found that 216 spatially coupled models had higher ensemble weights in the early and late stages of 217 the epidemic, while non-spatial models had higher weights around the peak of the 218 epidemic. The importance of including spatially coupled models in the optimally 219 weighted ensemble early in the epidemic supports the general notion that human 220 mobility may be particularly predictive of pathogen spread early in an epidemic [7, 221 30, 37, 38. In part, temporal shifts in weighting around the peak of the epidemic 222 were due to more accurate nationally-aggregated forecasts from the non-spatial mod-223 els. This result was consistent with a previous modeling analysis of the invasion of 224 chikungunya virus in Colombia, which showed that models fitted independently to 225 sub-national time series recreated national-level patterns well when aggregated [39]. 226 A shift in ensemble weights toward non-spatial models around the peak of the epi-227 demic was also due to less accurate department-level forecasts from the spatially 228 coupled models. At that point in the epidemic, prevalence was at its highest, which 229 means that we would expect local epidemics to be more endogenously driven and 230 less sensitive to pathogen introductions across departments. 231

In the target-oriented forecast performance analysis, we found that the equally 232 weighted ensemble generally outperformed individual models, with a few key ex-233 ceptions. In the months leading up to the peak of the epidemic, spatially coupled 234 models with a static R had slightly, but consistently, higher forecast scores with 235 respect to peak week and onset week. Like the model-specific analysis results, this 236 result illustrates the importance of human mobility in facilitating the spread of an 237 emerging pathogen across a landscape [30]. Individual models outperforming the 238 equally weighted ensemble model in the early phase of the epidemic is not wholly 239 surprising given that non-spatial models were represented equally in that ensemble 240 throughout the epidemic. Non-spatial models may be realistic when locations have 241 self-sustaining epidemics, but they are not appropriate for capturing early-phase 242 growth and its dependence on importations [40]. Another instance when individ-243 ual models had higher forecast scores than the equally weighted ensemble was with 244 respect to peak week for spatially coupled models with a static R in departments 245

along the Caribbean Coast. Compared to dynamic R models, the static R models 246 more accurately forecasted peak week in these departments (e.g., Magdalena, Cesar. 247 Sucre), as they did not forecast a late-stage resurgence in transmission. The equal 248 weighting of the dynamic R models in the ensemble therefore led to overall lower 249 peak week forecast scores for the ensemble relative to static R models. Still, our 250 results indicating that an equally weighted ensemble mostly outperformed individ-251 ual models adds to the growing literature highlighting the importance of ensemble 252 methods in epidemiological forecasting [16, 17, 27, 41, 42]. 253

We considered both equally and optimally weighted ensembles and found that 254 the equally weighted ensemble had a lower root mean square error than the op-255 timally weighted ensemble (RMSE=0.640 and 0.705, respectively)—therefore pro-256 viding slightly more accurate forecasts of the observed data (Fig. S23). With the 257 optimally weighted ensemble, which we updated at each data assimilation period, 258 we found that model weights changed over the course of the epidemic Fig. S17). Al-259 though this is intuitive given the changing nature of an emerging epidemic through 260 time [8], it may be problematic in practice. It is almost as if the ensemble weights 261 require their own forecast. On the one hand, promising new advances in ensem-262 ble modeling [27, 41]—such as adaptive stacking for seasonal influenza forecasting 263 [43]—are being used to address this issue of identifying optimal, adaptive weights 264 without training to historical data. On the other hand, in an emerging pathogen 265 context, establishing optimal model weights by way of model fitting and forecast 266 generation is often reliant on available incidence data (rather than historical data) 267 that is highly variable [44], given the delayed nature of data reporting [45]. In this 268 context, our results demonstrate that it is preferable to use an equally weighted en-269 semble to buffer against uncertainty in optimal ensemble weights. As is also being 270 demonstrated in forecasts of COVID-19, equally weighted ensembles can provide 271 accurate forecasts [26, 46, 44] and may be a better reflection of the considerable 272 structural uncertainty inherent to models of emerging pathogen transmission [24]. 273

A few limitations of our study should be noted. First, while an equally weighted 274 ensemble approach allowed us to consider contributions of several alternative model 275 assumptions, there was high uncertainty associated with these forecasts (sometimes 276 spanning orders of magnitude, see Fig. S24). Potential end-users of these types 277 of forecasts could consider high levels of uncertainty to be problematic for decision-278 making [47], though if the uncertainty does not affect the choice of a control measure, 279 then the uncertainty may not be as relevant [48]. In the future, ensemble approaches 280 aimed at increasing precision and reducing uncertainty [49, 27] could be used in con-281 junction with equally weighted ensembles. Second, we considered alternative models 282 across only three assumptions. With ZIKV transmission, there are additional struc-283 tural uncertainties that could be considered, such as the role of sexual transmission 284 [50]. In real-time applications of our or other Zika forecasting models, it could be 285 worthwhile to explore these types of ZIKV-specific structural uncertainties. Relat-286 edly, the static and dynamic R had minor differences in their formulations, such 287 that the static R also included a socioeconomic index. In future work, it could be 288 interesting to explore if the inclusion of this time-independent variable affected the 289 dynamic R. Third, in this analysis we did not explicitly consider delays in reporting 290 that likely would have occurred had these forecasts been generated in real time [51]. 291 In that context, temporally aggregating data to a wider interval (e.g., at 2-week 292 intervals rather than 1-week intervals) could potentially help mitigate the effects of 293 reporting delays to some extent. Fourth, we assumed that the reporting probability 294 was constant through time. Although this is a standard assumption [52] given the 295

lack of data to inform a time-varying relationship for this mechanistic element [53],
it would be interesting to include and test a reporting dynamics model (e.g., the reporting probability scales with incidence [54]) as an additional component included
in our ensemble framework. Fifth, we conducted this analysis at the departmental
level instead of this municipality level, which could obfuscate meaningful differences
across regions of a single department [29]. In future work, it would be useful to test
and assess our forecasting algorithm and outputs at different spatial scales [39].

As the world is reminded of on a daily basis with COVID-19, pathogen emer-303 gence is an ongoing phenomenon that will continue to pose threats in the future [55]. 304 A better understanding of an emerging pathogen's natural history could help to re-305 duce pathogen-specific structural uncertainties, but these insights may not always 306 occur in time to inform model development for real-time forecasting [8]. Our results 307 highlight important trade-offs between individual and ensemble models in this con-308 text. Specifically, we demonstrated that an equally weighted ensemble forecast was 309 almost always more accurate than individual models. Instances in which individ-310 ual models were better than the ensemble, or greatly improved the ensemble, also 311 provided insight. For example, incorporating human mobility into models improved 312 forecasts in the early and late phases of an epidemic, which underscores the impor-313 tance of making aggregated mobility data available early in an epidemic [56]. The 314 range of outcomes resulting from alternative modeling assumptions in model-specific 315 forecasts demonstrates why it will continue to be important to address structural 316 uncertainties in forecasting models in the future. 317

³¹⁸ 4 Materials and methods

319 4.1 Data

We used passive mandatory surveillance data for reported cases of Zika, from the 320 National Surveillance System (Sivigila) at the first administrative level (31 mainland 321 departments) in Colombia. To span the beginning, peak, and tail of the epidemic in 322 Colombia, we focused on the 60-week period between August 9, 2015 and October 323 1, 2016. We used the version of these data collated by Siraj et al. [29], as well as 324 modeled values of weekly average temperature and estimates of department-level 325 population from that data set. For some models, we worked with monthly estimates 326 of mosquito occurrence probability (i.e., dynamic R models) obtained from Bogoch 327 et al. [57], and for others we worked with time-averaged estimates (i.e., static R 328 models) from Kraemer et al. [58]. 329

For models that relied on cell phone data to describe human mobility, we used 330 anonymized and aggregated call detail records (CDRs). Every time a user receives 331 or makes a call, a CDR including the time, date, ID, and the tower (BTS) providing 332 the service is generated. The positions of the BTSs are georeferenced and so the 333 aggregated mobility between towers can be tracked in time. We used this information 334 to derive daily mobility matrices at the municipality level in Colombia from February 335 2015 to August 2015. Mobility matrices captured the number of individuals that 336 moved in each given day from one municipality to another (i.e., that appeared 337 in BTSs of different municipalities). The change for each day was captured by 338 comparing the last known municipality to the current one. No individual information 339 or records were available. 340

As these data did not align with the time frame of the epidemic, and to calculate a mobility matrix at a department level, we computed a representative mobility

matrix by summing all available CDRs within the municipalities of each department 343 and normalizing them to sum to one relative to the sum of CDRs originating from 344 that department. In five departments (Amazonas, Cudinamarca, Guainía, Vaupés, 345 Vichada), the proportion of CDRs linking callers within the same department was 346 below 60%. Given that this implied an unrealistically low proportion of time spent 347 within an individual's department of residence, we interpreted those values as id-348 iosyncrasies of the data and not representative of human mobility [59]. Thus, for 349 those five departments, we replaced the proportion of within-department CDRs with 350 the mean proportion of within-department CDRs from all other departments. We 351 then re-normalized the number of CDRs originating from each department in our 352 mobility matrix to sum to one. 353

354 4.2 Summary of models

To produce weekly forecasts of ZIKV transmission across Colombia, we sought to use 355 a computationally efficient model with the flexibility to include relevant epidemiolog-356 ical and ecological mechanisms. We used a previously described semi-mechanistic, 357 discrete-time, and stochastic model [60] that had been previously adapted and used 358 to model mosquito-borne pathogen transmission [61, 62]. Using this model, we 359 were able to account for the extended generation interval of ZIKV using overlapping 360 pathogen generations across up to five weeks of the generation interval distribu-361 tion of ZIKV [62]. Furthermore, we could specify this model to be either spatially 362 connected or non-spatial—a key assumption that we considered in our analysis. 363

We considered a suite of 16 models that spanned all combinations of four as-364 sumptions about human mobility across Colombia's 31 mainland departments, two 365 assumptions about the relationship between environmental conditions and the re-366 production number (R), and two assumptions about how many times Zika virus 367 was introduced to Colombia (Table 1). Twelve of 16 models allowed for spatial con-368 nectivity across departments [60], while four models were non-spatial. There were 369 up to two steps in the transmission process: transmission across departments (for 370 spatially connected models) and local transmission within departments. 371

Across departments, we simulated the movement of individuals using a spatial connectivity matrix (**H**), the d^{th} column of which corresponds to the proportion of time spent by residents of department d in all departments \vec{d} . Using this matrix, we redistributed infections in department d in week t ($I_{d,t}$) across \vec{d} as a multinomial random variable,

$$I'_{\vec{d}\,t} \sim \text{multinomial}(I_{d,t}, \mathbf{H}_{\vec{d},d}),\tag{1}$$

where the first and second arguments represent the number of trials and the probabilities of the outcomes, respectively. By taking this Lagrangian approach to modeling human mobility, transmission across departments can occur either by infected visitors transmitting to local susceptibles or susceptible visitors becoming infected by local infecteds. The relative occurrence of these events depends on the prevalence of infection, susceptibility, local transmission potential, and mobility patterns of a given pair of departments.

Within each department, we defined a variable representing the effective number of infections that could have generated new infections in week $t(I''_{d,t})$ as

$$I_{d,t}'' = \sum_{j=1}^{5} \omega_j^{GI} I_{d,t-j}',$$
(2)

where ω_j^{GI} is the probability that the generation interval is j weeks [63]. The relationship between $I''_{d,t}$ and the expected number of new local infections in week t+1 $(I_{d,t+1})$ follows

$$I_{d,t+1} = \beta_{d,t} \frac{I_{d,t}''}{N_d} S_{d,t},$$
(3)

where $\beta_{d,t}$ is the transmission coefficient, N_d is the total population, and $S_{d,t}$ is the total susceptible population prior to local transmission in week t. We accounted for the role of stochasticity in transmission by using the stochastic analogue of Eqn. 3, such that

$$I_{d,t+1} \sim \text{negative binomial}\left(\beta_{d,t} \frac{I_{d,t}''}{N_d} S_{d,t}, I_{d,t}''\right)$$
 (4)

where the first and second arguments are the mean and dispersion parameters, respectively [60].

To allow for comparison of the model's simulations of infections $(I_{d,t})$ with empirical data on reported cases $(y_{d,t})$, we applied a reporting probability (ρ) to simulated infections to obtain simulated cases $(C_{d,t})$, such that $C_{d,t} \sim \text{binomial}(I_{d,t}, \rho)$. Using this, we defined the contribution to the overall log-likelihood of the model and its parameters from a given department d and week t as

$$\boldsymbol{\ell}_{d,t}(\boldsymbol{\theta}_{t}) = \ln\left(\text{negative binomial}(y_{d,t}+1 \,|\, \boldsymbol{\phi}, C_{d,t}+1)\right),\tag{5}$$

where ϕ is a dispersion parameter that we estimated to account for variability in case reporting beyond that captured by ρ . Shifting $y_{d,t}$ and $C_{d,t}$ by one in eq. (5) was intended to safeguard against $\ell_{d,t}$ being undefined in situations where $C_{d,t} = 0$.

403 4.2.1 Assumptions about human mobility

We allowed for spatial coupling across departments in 12 of 16 models. In these 404 models, we informed **H** in three alternative ways: i) with mobility data extracted 405 from mobile phone CDRs, ii) with a gravity model, or iii) with a radiation model 406 (Fig. 1d-f). For the gravity model, we used parameters previously fitted to CDRs 407 from Spain and validated in West Africa [11]. For the radiation model, we calculated 408 human mobility fluxes according to the standard formulation of this model [64], 409 which depends only on the geographic distribution of population. In four of 16 410 models, we assumed that departments were spatially uncoupled (Table 1), such that 411 each department was modeled individually with its own set of parameters. In those 412 models, each department's epidemic was seeded independently with its own set of 413 imported infections. Further details about the spatially uncoupled models can be 414 found in the Supplemental Text. 415

416 4.2.2 Assumptions about environmental drivers of transmission

We parameterized the transmission coefficient, $\beta_{d,t}$, based on a description of the reproduction number, $R_{d,t}$, appropriate to environmental drivers for department dand time t. We considered two alternative formulations of $R_{d,t}$ that were informed by data that were available prior to the first reported case of Zika in Colombia. Specifically, both of these alternative formulations used different outputs from previous modeling efforts [6, 12] and because of this they contain slightly different components. Both formulations were defined such that

$$\beta_{d,t} = kR_{d,t} \tag{6}$$

where k is a scalar that we estimated over the course of the epidemic to account for the unknown magnitude of ZIKV transmission in Colombia. In addition to the summary below, further details about these formulations of $R_{d,t}$ are provided in the Supplementary Methods.

⁴²⁸ The formulation of $\beta_{d,t}$ that we refer to as "dynamic" is defined at each time ⁴²⁹ t in response to average temperature at that time $(T_{d,t})$ and mosquito occurrence ⁴³⁰ probability at that time $(OP_{d,t})$. This relationship can be expressed generically as

$$\beta_{d,t} = k \dot{R}_{d,t} (T_{d,t}, OP_{d,t} | c, \psi, \alpha, v), \tag{7}$$

where c, ψ, α , and v are parameters governing the relationship among $T_{d,t}, OP_{d,t}$, 431 and $R_{d,t}$. We informed the component of $R_{d,t}$ related to mosquito density with 432 monthly estimates of $OP_{d,t}$, which derive from geostatistical modeling of Aedes ae-433 gypti occurrence records globally [57]. Other components of $R_{d,t}$, which include 434 several temperature-dependent transmission parameters, were informed by labora-435 tory estimates [12]. Given that this formulation of $R_{d,t}$ was not validated against 436 field data prior to the Zika epidemic in Colombia, we estimated values of c, ψ, α , 437 and v over the course of the epidemic. 438

The formulation of $\beta_{d,t}$ that we refer to as "static" is defined as a time-averaged value that is constant across all times t. Temporal variation in $T_{d,t}$ is still accounted for, but its time-varying effect on $R_{d,t}$ is averaged out over all times \vec{t} to result in a temporally constant R_d . Mosquito occurrence probability is also incorporated through a temporally constant value (OP_d) [58]. The relationship among these variables can be expressed generically as

$$\beta_{d,t} = k\bar{R}_d(T_{d,\vec{t}}, OP_d, PPP_d), \tag{8}$$

where PPP_d is purchasing power parity in department d (a feature not included in the dynamic model) [65]. This input is an economic index that was intended to serve as a proxy for spatial variation in conditions that could affect exposure to mosquito biting, such as housing quality or air conditioning use [6]. Given that this formulation of \bar{R}_d was informed by data from outbreaks of Zika and chikungunya prior to the Zika epidemic in Colombia, we did not estimate its underlying parameters over the course of the epidemic in Colombia.

452 4.2.3 Assumptions about introduction events

Although many ZIKV infections were likely imported into Colombia throughout the 453 epidemic, we assumed that ZIKV introductions into either one or two departments 454 drove the establishment of ZIKV in Colombia [31]. Under the two different sce-455 narios, there was either one introduction event into one department or there were 456 two independent introduction events into two randomly drawn departments. For 457 each parameter set, the initial number of imported infections was seeded randomly 458 between one and five in a single week, the timing of which was estimated as a pa-459 rameter. Following the initial introduction(s), we assumed that ZIKV transmission 460 was driven by a combination of movement of infected people among departments 461 and local transmission within departments, as specified by each model. 462

463 4.3 Data assimilation and forecasting

⁴⁶⁴ For each particle, we produced a single forecast to "initialize" the model prior to ⁴⁶⁵ the first reported case in Colombia. Beginning with the time of the first reported

case in Colombia, we then assimilated new data, updated parameter estimates, 466 and generated forecasts every four weeks, consistent with the four-week frequency 467 used by Johansson et al. in an evaluation of dengue forecasts [16]. We specified 468 20,000 initial parameter sets $(\vec{\theta}_{1,n})$, indexed by n, by drawing independent samples 469 from prior distributions of each parameter [66] (see Supplemental Methods). Each 470 parameter set was used to generate a corresponding particle: a stochastic realization 471 of the state variables $(I_{d,1,n} \text{ and } C_{d,1,n})$. At each assimilation period, we normalized 472 log-likelihoods summed across departments over the preceding four weeks to generate 473 particle weights. 474

$$\omega(t,n) = \frac{\sum_{d} \sum_{\tau=t-3}^{t} \ell_{d,\tau}(\vec{\theta}_{t,n})}{\sum_{n} \sum_{d} \sum_{\tau=t-3}^{t} \ell_{d,\tau}(\vec{\theta}_{t,n})}.$$
(9)

Proportional to these particle weights $(\omega(t, n))$, we sampled 18,000 sets of corre-475 sponding parameters $(\vec{\theta}_t^{\text{resampled}})$ and state variables $(\{I_{d,t}^{\text{resampled}}, C_{d,t}^{\text{resampled}}\})$ from time t with replacement and used them at the next data assimilation step four 476 477 weeks later, where boldface indicates a set of parameters or state variables, respec-478 tively, over all n. In doing so, information including the initial prior assumptions 479 $(\vec{\theta}_{1,n})$ and the likelihoods at each four-week period was assimilated into the model 480 sequentially over time. Given that particle filtering algorithms are susceptible to 481 particle drift—or the convergence of particles onto very few states through iter-482 ative re-sampling [33]—we also generated 2,000 new parameter sets at each data 483 assimilation step. To do so, we drew random samples of model parameters from 484 a multivariate normal distribution with parameter means and covariances fitted to 485 the resampled 18,000 parameter sets $(\vec{\theta}_t^{\text{resampled}})$. Whereas the 18,000 resampled 486 parameter sets already included simulated values of state variables $I_{d,t,n}$ and $C_{d,t,n}$ 487 through time t, the 2,000 new parameter sets did not and so we informed initial conditions of $I_{d,t}^{\text{new}}$ with draws from $I_{d,t}^{\text{resampled}}$ for those parameter sets at the time 488 489 they were created. Together, the 18,000 resampled parameter sets $(\vec{\theta}_t^{\text{resampled}})$ and the 2,000 new parameter sets $(\vec{\theta}_t^{\text{new}})$ constituted the set of parameter sets used as 490 491 input for the next data assimilation step $(\vec{\theta}_{t+4} = \{\vec{\theta}_t^{\text{resampled}}, \vec{\theta}_t^{\text{new}}\})$. We also used 492 this new set of parameters as the basis for forecasts made at time t, which simply 493 involved simulating forward a single realization of the model for each parameter set. 494

495 4.4 Evaluating forecast performance

At each of the 15 time points at which we performed data assimilation through the 496 60-week forecasting period, we created an ensemble forecast that evenly weighted 497 contributions from each of the 16 models [46]. To populate this ensemble, we spec-498 ified 20,000 total samples, with 1,250 samples from each model. We assessed the 499 model-specific performance of individual and ensemble forecasts using log scores, 500 which are forecast scoring rules that assess both the precision and accuracy of fore-501 casts [67]. For a specific forecasting target, z, and model, m, the log score is defined 502 as $\log f_m(z^*|\mathbf{x})$, where $f_m(z|\mathbf{x})$ is the predicted density conditioned on the data, \mathbf{x} , 503 and z^* is the empirical value of the target Z [16]. 504

We computed log scores for departmental and national incidence over each fourweek assimilation period. Following [17], we used an expectation-maximization algorithm to generate ensemble weights for each model in each assimilation period. For each model, we computed 32 log scores (i.e., one for each department and one

nationally). To compute the ensemble weight for a given model feature, such as the static R assumption, we summed the weights of all models with that feature.

We assessed target-oriented forecast performance using log scores for three fore-511 casting targets: timing of peak week (within two weeks), incidence at peak week, 512 and onset week, which we defined as the week by which ten cumulative cases oc-513 curred. These choices were motivated by forecasting assessments for influenza and 514 dengue [16, 17, 18, 68] and deemed applicable to public health objectives for fore-515 casting an emerging pathogen such as Zika. For peak week and onset week, we used 516 modified log scores [18], such that predictions within two weeks of the correct week 517 were considered accurate. We evaluated a total of 7,680 log scores, reflecting three 518 targets for each of 16 models in each of 31 departments plus at the national level 519 and at each of 15 time points at which data assimilation occurred. 520

As log scores only yield a relative measure of model performance, we used fore-521 casting scores [18] as a way to retrospectively compare forecast performance for 522 different forecasting targets. Whereas log scores are preferable for comparing per-523 formance across models on the same data, forecasting scores are preferable for com-524 paring forecast performance across data composed of different locations and times. 525 A forecasting score is defined simply as the exponential of the average log score, 526 where the latter reflects an average over one or more indices, such as models, time 527 points, targets, or locations. 528

529 Data availability

The mobile phone data set used in this study is proprietary and subject to strict 530 privacy regulations. The access to this data set was granted after reaching a non-531 disclosure agreement with the proprietor, who anonymized and aggregated the orig-532 inal data before giving access to the authors. The data could be available on request 533 after negotiation of a non-disclosure agreement. The contact person is Enrique Frías-534 Martínez (enrique.friasmartinez@telefonica.com). Epidemiological, meteorological, 535 and demographic data are available from Siraj et al. [29] and additionally available 536 on https://github.com/roidtman/eid_ensemble_forecasting. 537

⁵³⁸ Author contributions

RJO, EO, MUGK, CMB, MAJ, CAM, RCR, IR-B, MG-H, and TAP conceptual-539 ized the study; RJO, EO, MUGK, CAC-O, EC-R, AM-C, PC, LER, VC, PA, GE, 540 JHH, SCH, ASS, EF-M, and MG-H provided and / or processed data; RJO, EO, 541 MUGK, CA-O, EC-R, SM-C, PC, LER, VC, PA, EF-M, MG-H, and TAP partic-542 ipated in biweekly meetings to provide feedback on research: RJO, EO, MUGK. 543 MG-H, and TAP developed the model and wrote the first draft of the manuscript; 544 RJO, EO, MUGK, JHH, and SCH analyzed data; EO, MG-H, and TAP supervised 545 the research; all authors reviewed the manuscript. 546

547 Acknowledgements

The authors would like to thank Clara Palau Montava for help with managing the early stages of this project. The authors would additionally like to thank the UNICEF-Colombia Representative, Aida Oliver Arostegui, INS Director, Martha

- $_{\tt 551}$ $\,$ Lucia Ospina Martinez, and the past and present Ministers of the Ministry of Health,
- ⁵⁵² Juan Pablo Uribe Restrepo and Fernado Ruiz Gomez.

RJO acknowledges support from an Eck Institute for Global Health Fellowship, 553 GLOBES grant, Arthur J. Schmitt Fellowship, and the UNICEF Office of Innova-554 tion. MUGK is supported by The Branco Weiss Fellowship - Society in Science, 555 administered by the ETH Zurich and acknowledges funding from the Oxford Martin 556 School and the European Union Horizon 2020 project MOOD (#874850). SCH is 557 supported by the Wellcome Trust (220414/Z/20/Z). This research was funded in 558 whole, or in part, by the Wellcome Trust [Grant number 220414/Z/20/Z]. For the 559 purpose of open access, the author has applied a CC BY public copyright licence to 560 any Author Accepted Manuscript version arising from this submission 561

References

- Kate E. Jones et al. "Global trends in emerging infectious diseases". In: Nature 451.7181 (Feb. 2008), pp. 990-993. ISSN: 1476-4687. DOI: 10.1038/ nature06536. URL: https://doi.org/10.1038/nature06536.
- Katherine F. Smith et al. "Global rise in human infectious disease outbreaks". In: Journal of The Royal Society Interface 11.101 (2014), p. 20140950. DOI: 10.1098/rsif.2014.0950. URL: https://royalsocietypublishing.org/ doi/abs/10.1098/rsif.2014.0950.
- Juliet Bedford et al. "A new twenty-first century science for effective epidemic response". In: *Nature* 575.7781 (Nov. 2019), pp. 130–136. DOI: 10.1038/ s41586-019-1717-y. URL: https://doi.org/10.1038/s41586-019-1717y.
- [4] Giovanni Lo Iacono et al. "Using modelling to disentangle the relative contributions of zoonotic and anthroponotic transmission: the case of Lassa fever". In: *PLOS Neglected Tropical Diseases* 9.1 (Jan. 2015), pp. 1–13. DOI: 10.1371/journal.pntd.0003398. URL: https://doi.org/10.1371/journal.pntd.0003398.
- [5] Thomas C Quinn. "Global burden of the HIV pandemic". In: *The Lancet* 348.9020 (1996), pp. 99-106. ISSN: 0140-6736. DOI: https://doi.org/10.1016/S0140-6736(96)01029-X. URL: http://www.sciencedirect.com/science/article/pii/S014067369601029X.
- T. Alex Perkins et al. "Model-based projections of Zika virus infections in childbearing women in the Americas". In: *Nature Microbiology* 1.9 (2016), p. 16126. DOI: 10.1038/nmicrobiol.2016.126. URL: https://doi.org/10.1038/nmicrobiol.2016.126.
- [7] Moritz U. G. Kraemer et al. "Spread of yellow fever virus outbreak in Angola and the Democratic Republic of the Congo 2015-2016: a modelling study". In: *The Lancet Infectious Diseases* 17.3 (Mar. 2017), pp. 330–338. DOI: 10. 1016/S1473-3099(16)30513-8. URL: https://doi.org/10.1016/S1473-3099(16)30513-8.
- C. Jessica E. Metcalf and Justin Lessler. "Opportunities and challenges in modeling emerging infectious diseases". In: Science 357.6347 (2017), pp. 149– 152. DOI: 10.1126/science.aam8335. URL: https://science.sciencemag. org/content/357/6347/149.

- J. O. Lloyd-Smith et al. "Superspreading and the effect of individual variation on disease emergence". In: *Nature* 438.7066 (Nov. 2005), pp. 355–359. ISSN: 1476-4687. DOI: 10.1038/nature04153. URL: https://doi.org/10.1038/ nature04153.
- [10] Amy Wesolowski et al. "Impact of human mobility on the emergence of dengue epidemics in Pakistan". In: *Proceedings of the National Academy of Sciences* 112.38 (2015), pp. 11887–11892. DOI: 10.1073/pnas.1504964112. URL: https://www.pnas.org/content/112/38/11887.
- [11] M. U. G. Kraemer et al. "Utilizing general human movement models to predict the spread of emerging infectious diseases in resource poor settings". In: *Scientific Reports* 9.1 (2019), p. 5151. ISSN: 2045-2322. DOI: 10.1038/s41598-019-41192-3.
- [12] Erin A. Mordecai et al. "Detecting the impact of temperature on transmission of Zika, dengue, and chikungunya using mechanistic models". In: *PLOS Ne*glected Tropical Diseases 11.4 (Apr. 2017), pp. 1–18. DOI: 10.1371/journal. pntd.0005568. URL: https://doi.org/10.1371/journal.pntd.0005568.
- Birgit Nikolay et al. "Transmission of Nipah Virus 14 Years of Investigations in Bangladesh". In: New England Journal of Medicine 380.19 (2019). PMID: 31067370, pp. 1804–1814. DOI: 10.1056/NEJMoa1805376. URL: https://doi. org/10.1056/NEJMoa1805376.
- [14] Gytis Dudas et al. "MERS-CoV spillover at the camel-human interface". In: eLife 7 (Jan. 2018), e31257. ISSN: 2050-084X. DOI: 10.7554/eLife.31257. URL: https://doi.org/10.7554/eLife.31257.
- [15] Katriona Shea et al. "Harnessing multiple models for outbreak management". In: Science 368.6491 (2020), pp. 577-579. DOI: 10.1126/science.abb9934. URL: https://science.sciencemag.org/content/368/6491/577.
- [16] Michael A. Johansson et al. "An open challenge to advance probabilistic forecasting for dengue epidemics". In: *Proceedings of the National Academy of Sciences* 116.48 (2019), pp. 24268-24274. DOI: 10.1073/pnas.1909865116. URL: https://www.pnas.org/content/116/48/24268.
- [17] Nicholas G. Reich et al. "Accuracy of real-time multi-model ensemble forecasts for seasonal influenza in the U.S." In: *PLOS Computational Biology* 15.11 (Nov. 2019), pp. 1–19. DOI: 10.1371/journal.pcbi.1007486. URL: https://doi.org/10.1371/journal.pcbi.1007486.
- [18] Nicholas G. Reich et al. "A collaborative multiyear, multimodel assessment of seasonal influenza forecasting in the United States". In: Proceedings of the National Academy of Sciences 116.8 (2019), pp. 3146–3154. DOI: 10.1073/ pnas.1812594116. URL: https://www.pnas.org/content/116/8/3146.
- [19] Craig J. McGowan et al. "Collaborative efforts to forecast seasonal influenza in the United States, 2015–2016". In: *Scientific Reports* 9.1 (Jan. 2019), p. 683.
 DOI: 10.1038/s41598-018-36361-9. URL: https://doi.org/10.1038/ s41598-018-36361-9.
- [20] Leah R. Johnson et al. "Phenomenological forecasting of disease incidence using heteroskedastic Gaussian processes: A dengue case study". In: *The Annals* of Applied Statistics 12.1 (2018), pp. 27–66. DOI: 10.1214/17-AOAS1090. URL: https://doi.org/10.1214/17-AOAS1090.

- [21] Sara Y. Del Valle et al. "Summary results of the 2014-2015 DARPA Chikungunya challenge". In: *BMC Infectious Diseases* 18.1 (May 2018), p. 245. DOI: 10.1186/s12879-018-3124-7. URL: https://doi.org/10.1186/s12879-018-3124-7.
- [22] Cécile Viboud et al. "The RAPIDD ebola forecasting challenge: Synthesis and lessons learnt". In: *Epidemics* 22 (2018). The RAPIDD Ebola Forecasting Challenge, pp. 13–21. DOI: https://doi.org/10.1016/j.epidem.2017.08.002. URL: http://www.sciencedirect.com/science/article/pii/S1755436517301275.
- [23] ZIKAVAT Collaboration et al. "Preliminary results of models to predict areas in the Americas with increased likelihood of Zika virus transmission in 2017". In: *bioRxiv* (2017). DOI: 10.1101/187591. URL: https://www.biorxiv.org/ content/early/2017/09/29/187591.
- [24] Katriona Shea et al. "COVID-19 reopening strategies at the county level in the face of uncertainty: Multiple Models for Outbreak Decision Support". In: medRxiv (2020). DOI: 10.1101/2020.11.03.20225409. URL: https://www.medrxiv.org/content/early/2020/11/05/2020.11.03.20225409.
- [25] Dylan B. George et al. "Technology to advance infectious disease forecasting for outbreak management". In: Nature Communications 10.1 (Sept. 2019), p. 3932. ISSN: 2041-1723. DOI: 10.1038/s41467-019-11901-7. URL: https://doi.org/10.1038/s41467-019-11901-7.
- [26] Evan L Ray et al. "Ensemble Forecasts of Coronavirus Disease 2019 (COVID-19) in the U.S." In: medRxiv (2020). DOI: 10.1101/2020.08.19.20177493.
 URL: https://www.medrxiv.org/content/early/2020/08/22/2020.08. 19.20177493.
- [27] Logan C. Brooks et al. "Nonmechanistic forecasts of seasonal influenza with iterative one-week-ahead distributions". In: *PLOS Computational Biology* 14.6 (June 2018), pp. 1–29. DOI: 10.1371/journal.pcbi.1006134. URL: https://doi.org/10.1371/journal.pcbi.1006134.
- [28] G. Chowell et al. "Real-time forecasting of epidemic trajectories using computational dynamic ensembles". In: *Epidemics* 30 (2020), p. 100379. ISSN: 1755-4365. DOI: https://doi.org/10.1016/j.epidem.2019.100379. URL: http: //www.sciencedirect.com/science/article/pii/S1755436519301112.
- [29] Amir S. Siraj et al. "Spatiotemporal incidence of Zika and associated environmental drivers for the 2015-2016 epidemic in Colombia". In: Scientific Data 5.1 (2018), p. 180073. DOI: 10.1038/sdata.2018.73. URL: https://doi.org/10.1038/sdata.2018.73.
- [30] Moritz U. G. Kraemer et al. "The effect of human mobility and control measures on the COVID-19 epidemic in China". In: Science 368.6490 (2020), pp. 493-497. DOI: 10.1126/science.abb4218. URL: https://science.sciencemag.org/content/368/6490/493.
- [31] Allison Black et al. "Genomic epidemiology supports multiple introductions and cryptic transmission of Zika virus in Colombia". In: *BMC Infectious Diseases* 19.1 (2019), p. 963. DOI: 10.1186/s12879-019-4566-2. URL: https: //doi.org/10.1186/s12879-019-4566-2.
- [32] Roni Rosenfeld. The EM Algorithm. 1997. URL: http://www.cs.cmu.edu/ afs/cs.cmu.edu/academic/class/11761-s97/WWW/tex/EM.ps.

- [33] Michael C. Dietze. *Ecological Forecasting*. Princeton University Press, 2017. ISBN: 9780691160573. URL: http://www.jstor.org/stable/j.ctvc7796h.
- [34] Nicholas B. DeFelice et al. "Ensemble forecast of human West Nile virus cases and mosquito infection rates". In: *Nature Communications* 8.1 (Feb. 2017), p. 14592. ISSN: 2041-1723. DOI: 10.1038/ncomms14592. URL: https://doi. org/10.1038/ncomms14592.
- [35] Wan Yang, Alicia Karspeck, and Jeffrey Shaman. "Comparison of Filtering Methods for the Modeling and Retrospective Forecasting of Influenza Epidemics". In: *PLOS Computational Biology* 10.4 (Apr. 2014), pp. 1–15. DOI: 10.1371/journal.pcbi.1003583. URL: https://doi.org/10.1371/journal.pcbi.1003583.
- [36] Rachel E. Baker et al. "Susceptible supply limits the role of climate in the early SARS-CoV-2 pandemic". In: Science 369.6501 (2020), pp. 315-319. DOI: 10.1126/science.abc2535. URL: https://science.sciencemag.org/content/369/6501/315.
- [37] S. Cauchemez et al. "Local and regional spread of chikungunya fever in the Americas". In: *Euro surveillance* 19.28 (2014), pp. 20854–20854. ISSN: 1560-7917. DOI: 10.2807/1560-7917.es2014.19.28.20854. URL: https://pubmed.ncbi.nlm.nih.gov/25060573.
- [38] Michael A. Johansson et al. "Nowcasting the Spread of Chikungunya Virus in the Americas". In: *PLOS ONE* 9.8 (Aug. 2014), pp. 1–8. DOI: 10.1371/ journal.pone.0104915. URL: https://doi.org/10.1371/journal.pone. 0104915.
- [39] Sean M. Moore et al. "Local and regional dynamics of chikungunya virus transmission in Colombia: the role of mismatched spatial heterogeneity". In: *BMC Medicine* 16.1 (Aug. 2018), p. 152. ISSN: 1741-7015. DOI: 10.1186/s12916-018-1127-2. URL: https://doi.org/10.1186/s12916-018-1127-2.
- [40] Shengjie Lai et al. "Seasonal and interannual risks of dengue introduction from South-East Asia into China, 2005-2015". In: *PLOS Neglected Tropical Diseases* 12.11 (Nov. 2018), pp. 1–16. DOI: 10.1371/journal.pntd.0006743.
 URL: https://doi.org/10.1371/journal.pntd.0006743.
- [41] Tom Lindström, Michael Tildesley, and Colleen Webb. "A Bayesian ensemble approach for epidemiological projections". In: *PLOS Computational Biology* 11.4 (Apr. 2015), pp. 1–30. DOI: 10.1371/journal.pcbi.1004187. URL: https://doi.org/10.1371/journal.pcbi.1004187.
- [42] Teresa K. Yamana, Sasikiran Kandula, and Jeffrey Shaman. "Superensemble forecasts of dengue outbreaks". In: Journal of The Royal Society Interface 13.123 (2016), p. 20160410. DOI: 10.1098/rsif.2016.0410. URL: https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2016.0410.
- [43] Thomas McAndrew and Nicholas G. Reich. Adaptively stacking ensembles for influenza forecasting with incomplete data. 2020. arXiv: 1908.01675 [stat.AP].
- [44] Evan L. Ray et al. "Challenges in training ensembles to forecast COVID-19 cases and deaths in the United States". In: International Institute of Forecasters (2021). URL: https://forecasters.org/blog/2021/04/09/ challenges-in-training-ensembles-to-forecast-covid-19-casesand-deaths-in-the-united-states/.

- [45] T. Alex Perkins et al. "Estimating unobserved SARS-CoV-2 infections in the United States". In: Proceedings of the National Academy of Sciences 117.36 (2020), pp. 22597-22602. ISSN: 0027-8424. DOI: 10.1073/pnas.2005476117. eprint: https://www.pnas.org/content/117/36/22597.full.pdf.
- [46] Thomas McAndrew et al. "Aggregating predictions from experts: A review of statistical methods, experiments, and applications". In: WIREs Computational Statistics (), e1514. DOI: 10.1002/wics.1514. URL: https:// onlinelibrary.wiley.com/doi/abs/10.1002/wics.1514.
- [47] Korryn Bodner, Marie-Josée Fortin, and Péter K. Molnár. "Making predictive modelling ART: accurate, reliable, and transparent". In: *Ecosphere* 11.6 (2020), e03160. DOI: 10.1002/ecs2.3160. URL: https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1002/ecs2.3160.
- [48] Shou-Li Li et al. "Essential information: Uncertainty and optimal control of Ebola outbreaks". In: *Proceedings of the National Academy of Sciences* 114.22 (2017), pp. 5659-5664. ISSN: 0027-8424. DOI: 10.1073/pnas.1617482114. eprint: https://www.pnas.org/content/114/22/5659.full.pdf. URL: https://www.pnas.org/content/114/22/5659.
- [49] Sen Pei and Jeffrey Shaman. "Counteracting structural errors in ensemble forecast of influenza outbreaks". In: *Nature Communications* 8.1 (Oct. 2017), p. 925. ISSN: 2041-1723. DOI: 10.1038/s41467-017-01033-1. URL: https://doi.org/10.1038/s41467-017-01033-1.
- [50] Antoine Allard et al. "The risk of sustained sexual transmission of Zika is underestimated". In: *PLOS Pathogens* 13.9 (Sept. 2017), pp. 1–12. DOI: 10. 1371/journal.ppat.1006633. URL: https://doi.org/10.1371/journal. ppat.1006633.
- [51] Axel Bonačić Marinović et al. "Quantifying Reporting Timeliness to Improve Outbreak Control". In: *Emerging Infectious Disease journal* 21.2 (2015), p. 209.
 DOI: 10.3201/eid2102.130504. URL: https://doi.org/10.3201/eid2102. 130504.
- [52] Matt J. Keeling and Pejman Rohani. Modeling Infectious Diseases in Humans and Animals. Princeton University Press, 2008. ISBN: 9780691116174. URL: http://www.jstor.org/stable/j.ctvcm4gk0.
- [53] LT Figueiredo, SM Cavalcante, and MC Simões. "Dengue serologic survey of schoolchildren in Rio de Janeiro, Brazil, in 1986 and 1987". In: Bull Pan Am Health Organ. 24.2 (1990), pp. 217–225. URL: https://iris.paho.org/ bitstream/handle/10665.2/27164/ev24n2p217.pdf?sequence=1.
- [54] Jue Tao Lim et al. "Time varying methods to infer extremes in dengue transmission dynamics". In: *PLOS Computational Biology* 16.10 (Oct. 2020), pp. 1–19. DOI: 10.1371/journal.pcbi.1008279. URL: https://doi.org/10.1371/journal.pcbi.1008279.
- [55] Honglei Sun et al. "Prevalent Eurasian avian-like H1N1 swine influenza virus with 2009 pandemic viral genes facilitating human infection". In: Proceedings of the National Academy of Sciences (2020). ISSN: 0027-8424. DOI: 10.1073/ pnas.1921186117. eprint: https://www.pnas.org/content/early/2020/ 06/23/1921186117.full.pdf. URL: https://www.pnas.org/content/ early/2020/06/23/1921186117.

- [56] Caroline O. Buckee et al. "Aggregated mobility data could help fight COVID-19". In: Science 368.6487 (2020), pp. 145-146. DOI: 10.1126/science. abb8021. URL: https://science.sciencemag.org/content/368/6487/ 145.2.
- [57] Isaac I. Bogoch et al. "Potential for Zika virus introduction and transmission in resource-limited countries in Africa and the Asia-Pacific region: a modelling study". eng. In: *The Lancet. Infectious diseases* 16.11 (Nov. 2016), pp. 1237– 1245. DOI: 10.1016/S1473-3099(16)30270-5. URL: https://doi.org/10. 1016/S1473-3099(16)30270-5.
- [58] Moritz UG Kraemer et al. "The global distribution of the arbovirus vectors *Aedes aegypti* and *Ae. albopictus*". In: *eLife* 4 (June 2015), e08347. ISSN: 2050- 084X. DOI: 10.7554/eLife.08347. URL: https://doi.org/10.7554/eLife. 08347.
- [59] Amy Wesolowski et al. "Heterogeneous mobile phone ownership and usage patterns in Kenya". In: PLOS ONE 7.4 (Apr. 2012), pp. 1–6. DOI: 10.1371/ journal.pone.0035319. URL: https://doi.org/10.1371/journal.pone. 0035319.
- [60] Yingcun Xia, Ottar N. Bjørnstad, and Bryan T. Grenfell. "Measles metapopulation dynamics: A gravity model for epidemiological coupling and dynamics". English (US). In: American Naturalist 164.2 (Aug. 2004), pp. 267–281. ISSN: 0003-0147. DOI: 10.1086/422341.
- [61] Rachel J. Oidtman et al. "Inter-annual variation in seasonal dengue epidemics driven by multiple interacting factors in Guangzhou, China". In: *Nature Communications* 10.1 (Mar. 2019). DOI: 10.1038/s41467-019-09035-x. URL: https://doi.org/10.1038/s41467-019-09035-x.
- [62] T.A. Perkins et al. "Estimating drivers of autochthonous transmission of chikungunya virus in its invasion of the Americas". In: *PLOS Currents Outbreaks* (2017). DOI: 10.1371/currents.outbreaks.a4c7b6ac10e0420b1788c9767946d1fc.
- [63] Amir S. Siraj et al. "Temperature modulates dengue virus epidemic growth rates through its effects on reproduction numbers and generation intervals". In: *PLOS Neglected Tropical Diseases* 11.7 (July 2017), pp. 1–19. DOI: 10. 1371/journal.pntd.0005797. URL: https://doi.org/10.1371/journal.pntd.0005797.
- [64] Filippo Simini et al. "A universal model for mobility and migration patterns".
 In: Nature 484.7392 (2012), pp. 96–100. ISSN: 1476-4687. DOI: 10.1038/ nature10856. URL: https://doi.org/10.1038/nature10856.
- [65] William D. Nordhaus. "Geography and macroeconomics: New data and new findings". In: Proceedings of the National Academy of Sciences 103.10 (2006), pp. 3510-3517. ISSN: 0027-8424. DOI: 10.1073/pnas.0509842103. eprint: https://www.pnas.org/content/103/10/3510.full.pdf. URL: https://www.pnas.org/content/103/10/3510.
- [66] M. Sanjeev Arulampalam et al. "A tutorial on particle filters for online nonlinear/nongaussian Bayesian tracking". In: *IEEE Transactions on Signal Processing* 50.4 (Feb. 2002), pp. 174–188.

- [67] Tilmann Gneiting, Fadoua Balabdaoui, and Adrian Raftery. "Probabilistic forecasts, calibration and sharpness". In: Journal of the Royal Statistical Society: Series B 69 (2007), pp. 243–268. URL: https://doi.org/10.1371/ journal.pcbi.1004187.
- Sen Pei et al. "Forecasting the spatial transmission of influenza in the United States". In: Proceedings of the National Academy of Sciences 115.11 (2018), pp. 2752-2757. DOI: 10.1073/pnas.1708856115. URL: https://www.pnas.org/content/115/11/2752.

¹ Supplemental methods

² Priors on parameters common to all models

In each model that we considered, we iteratively estimated the reporting probability (ρ) , dispersion parameter of the negative binomial distribution (ϕ) , R multiplier (k), and the timing of the first infection (ι) . When possible, we leveraged previous estimates of parameters to inform prior distributions for model parameters. In some instances, we used dengue-specific parameter estimates as priors for Zika-specific parameters [1].

For the reporting probability (ρ) , we assumed a mean of 0.2 and a variance of 0.05. Although this mean and variance are not directly informed by an empirical study of Zika reporting, they are in line with what we would expect for dengue [2, 3] and Zika [4]. We assumed that ρ was a beta random variable and, using the method of moments, we specified a prior distribution for ρ such that

$$\rho \sim \text{beta}(0.44, 1.76)$$
(10)

¹⁴ which resulted in mean and variance consistent with our prior assumptions.

The dispersion parameter of the negative binomial distribution accounts for variability in case reporting beyond that captured by ρ . Lower values of the dispersion parameter indicate overdispersion, such that variability in cases cannot be explained by a single rate of case incidence, as would be generated by a Poisson distribution with rate $\rho I_{d,t}$. Given the likelihood of variation in the reporting probability over the course of the epidemic [5] and across departments [6], we specified a uniform prior for ϕ ,

$$\phi \sim \text{uniform}(0, 1), \tag{11}$$

which resulted in a level of overdispersion in reporting equal to at least a geometric a distribution ($\phi = 1$) but potentially greater ($\phi < 1$).

To relate the transmission coefficient $(\beta_{d,t})$ to environmentally driven descriptions of the reproduction number $(R_{d,t})$, we used a multiplier (k) that applied to both the static and dynamic models of R. We specified a gamma prior distribution,

$$k \sim \text{gamma}(2.25, 0.75),$$
 (12)

27 the parameters of which were chosen by moment matching to result in a mean of 28 three and a variance of two.

To seed the Zika epidemic in Colombia, we assumed that undetected transmission could have been occurring at any time throughout the first 34 weeks of 2015. For reference, the first case was not reported until August 9, 2015 (week 35). Thus, we specified a uniform prior for the date of the initial introduction (ι_1) between weeks 1 and 34 of 2015. We assumed that the location of the first introduction (l_1) could have been any of the departments in Colombia with equal prior probability.

³⁵ Assumptions about human mobility

³⁶ Spatially coupled models

For the spatially coupled models, we assumed that transmission was coupled across departments by human mobility. In these models, we informed the spatial connectivity matrix in three ways: i) aggregated mobility patterns extracted from mobile phone call detail records (CDRs), ii) a gravity model, or iii) a radiation model. In

the gravity and radiation models, $T_{i,j}$ is defined as the total number of trips from 41 department i to department j. This takes the form $T_{i,j} = c \frac{N_i^{\alpha} N_j^{\beta}}{d_{i,j}^{\alpha}}$ in the gravity 42 model and $T_{i,j} = T_i \frac{N_i N_j}{(N_i + s_{i,j})(N_i + N_j + s_{i,j})}$ in the radiation model, where N_i and N_j are 43 population sizes at the origin i and destination j, $d_{i,j}^{\gamma}$ is the distance between i and 44 j, $s_{i,j}$ is the total population within radius $d_{i,j}$ from i, and T_i is the total number 45 of individuals who make a trip. The parameters c, α, β , and γ were fitted to CDRs 46 from Spain and validated in West Africa [7]. All three mobility models were row-47 normalized to correspond to the proportion of time spent by residents of department 48 i in department j. 49

⁵⁰ Spatially uncoupled models

For the spatially uncoupled models, we assumed that each department's epidemic 51 occurred independently of all other departments. We used the same prior distribu-52 tions as described above for ρ , ϕ , k for each department. Under this assumption, we 53 did not include a parameter for the location of the initial introduction into Colom-54 bia, as we instead were concerned with the initial introduction into that department. 55 It was still necessary to specify the timing of that introduction and, for models that 56 considered it, the timing of a second introduction. Following the rationale that 57 undetected transmission could have occurred for up to 34 weeks prior to the first 58 reported case in a given department, we assumed an even prior for the timing of the 59 introduction(s) into a given department. 60

⁶¹ Assumptions about environmental drivers of transmission

62 Dynamic model

- ⁶³ The environmentally driven component of $\beta_{d,t}$ for the dynamic model, $\tilde{R}_{d,t}(T_{d,t}, OP_{d,t}|c, \psi, \alpha, v)$,
- was defined as the product of two functions: one that depends on $T_{d,t}$ and one that depends on $OP_{d,t}$.

The function of $OP_{d,t}$ was defined as $c(-\log(1 - OP_{d,t}))$, which converts occurrence probability into an expectation of mosquito abundance [8]. The constant cscales this expectation to account for uncertainty in its magnitude, which we estimated and assumed to have a gamma-distributed prior with a shape parameter of 16 and a rate parameter of 0.07. This choice of prior resulted in average $\tilde{R}_{d,t}$ being equal to one when evaluated at the mean values of the prior for c.

The function of $T_{d,t}$ was based on a temperature-dependent description of the 72 basic reproduction number by Mordecai et al. [9]. Specifically, we used the version 73 of that model based on Briere functions for parameters that were not otherwise ac-74 counted for in estimates of mosquito occurrence probability [10]. Those parameters 75 include the mosquito biting rate (a), mosquito-to-human transmission probability 76 (b), human-to-mosquito transmission probability (c), average adult mosquito lifes-77 pan (lf), and parasite development rate (PDR). Those functions combine to form 78 the temperature-dependent component of $R_{d,t}$, 79

$$a(T_{d,t})b(T_{d,t})c(T_{d,t})e^{-1/(lf(T_{d,t})PDR(T_{d,t}))}PDR(T_{d,t}).$$
(13)

We did not include other parameters from Mordecai et al. [9] related to mosquitoes, such as immature development rate and egg-to-adult survival rate, as the effects of

those parameters were accounted for in estimates of $OP_{d,t}$ from Kraemer et al. [10].

To reduce the number of parameters that we needed to estimate, we worked with 83 a simplified description of the temperature curves produced by eq. (13) (Fig. S5). 84 To do so, we took random draws from the posterior distribution of parameters from 85 Mordecai et al. [9], computed functions of temperature according to eq. (13), and 86 fitted parameters of a skew normal distribution to those curves by least squares. 87 The skew normal distribution has three parameters—location (ψ), scale (α), shape 88 (v)—that together control the mean, variance, and skew of the distribution, which is 89 sufficient to approximate the uncertainty in posterior predictions of the temperature 90 curves described by eq. (13). We then took the mean and covariance across those 91 estimates of ψ , α , and v to describe their variation with a multivariate normal 92 distribution, which was the prior distribution we used for those parameters at the 93 first step of our particle filter. 94

95 Static model

The environmentally driven component of β_d for the dynamic model, $R_d(T_{d,\vec{t}}, OP_d, PPP_d)$, 96 was defined as the product of three functions: one that depends on $T_{d,\vec{t}}$, one that 97 depends on OP_d , and one that depends on PPP_d . We used values of R_d at the 5 98 km x 5 km grid cell level as calculated by Perkins et al. [8]. To aggregate them to 99 the department level, we took a population-weighted mean of R_d across 5 km x 5 100 km grid cells within each department. Although we made no other modifications to 101 the calculation of R_d , we summarize the methodology used by Perkins et al. [8] in 102 the interest of comparability with the dynamic model. 103

The function of OP_d was defined as $-\log(1 - OP_d))$, which was the same procedure used to convert occurrence probability into an expectation of mosquito abundance as used in the dynamic model. For the static model, however, we followed Perkins et al. [8] and relied on a description of occurrence probability that was not defined on a time-varying basis [10].

Rather than estimate a scaling parameter like c in the dynamic model, we relied on a scaling parameter defined as a function of PPP_d by Perkins et al. [8]. This function took the form of a monotonically decreasing, cubic spline function estimated with a shape-constrained additive model. The data that informed this estimate of both the relationship with PPP_d and the magnitude of \bar{R}_d were outbreak sizes of 12 chikungunya outbreaks and one Zika outbreak [8] that occurred prior to the ZIKV invasion of the Americas.

The function of $T_{d,\vec{t}}$ was based on a temperature-dependent description of the basic reproduction number that includes temperature-dependent descriptions of mosquito mortality (μ) [11] and the extrinsic incubation period (n) [12]. Those functions contribute to an expression for monthly contributions to $R_{d,t}$,

$$\frac{bca^2 e^{-\mu(T_{d,t})n(T_{d,t})}}{\mu(T_{d,t})r},$$
(14)

along with constants that represent the mosquito-to-human transmission probability (b), human-to-mosquito transmission probability (c), mosquito biting rate (a), and rate of recovery from human infection (r). For each department d, we took the average of the six largest values of eq. (14) as the contribution of $T_{d,\vec{t}}$ to \bar{R}_d .

124 Assumptions about introduction events

125 One-introduction models

Each of the one-introduction models assumed infections were seeded at one point in time (ι_1) in one location (l_1) , as specified in the general model parameters section of the Supplementary Methods.

129 Two-introduction models

Each of the two-introduction models made similar assumptions about ι_1 and l_1 for the first introduction. For these models, we additionally specified the timing (ι_2) and location (l_2) of the second introduction events, for which we assumed even priors. In reality, there were likely more than only one or two introductions throughout the epidemic, but genomic epidemiology suggests the majority of local transmission resulted from only one or two introductions [13].

136 Supplemental tables

Symbol	Definition	Class
ρ	Reporting probability	Estimated parameter
ϕ	Overdispersion in reporting	Estimated parameter
k	$R_{d,t}$ multiplier	Estimated parameter
ι_1	Timing of first infection	Estimated parameter
ι_2	Timing of second infection	Estimated parameter
l_1	Location of first infection	Estimated parameter
l_2	Location of second infection	Estimated parameter
c	Dynamic R scalar	Estimated parameter
ψ	Location for skew normal	Estimated parameter
α	Scale for skew normal	Estimated parameter
v	Shape for skew normal	Estimated parameter
θ	Parameter set	Set of estimated parameters
$I_{d,t}$	Simulated infections	State variable
$I'_{d.t}$	Redistributed infections	State variable
$I_{d,t}^{\prime\prime}$	Effective number of infections	State variable
$C_{d,t}$	Simulated cases	State variable
$y_{d,t}$	Observed data	Data
d	Department index	Data
t	Time; $t = 0,, T$	Data
n	Particle index; $n = 1,, N$	Data
η	Timing of first reported case	Data

Table S1: Mathematical symbols for parameters, state variables, and data with their respective meanings.

¹³⁷ Supplemental figures



Figure S1: Proportion of forecast trajectories predicting zero infections. Forecast trajectories are specific to each particle and this example is from the model using CDR-derived mobility data, static R, and one importation event, corresponding to Fig. S2.



Figure S2: Proportion of particles that remain from the original particle ensemble present in the retained ensemble at each assimilation period. This example is from the model using CDR-derived mobility data, static R, and one importation event, corresponding to Fig. S1.



Figure S3: Cumulative observed versus forecasted incidence at 4-week ahead intervals for each individual model for Antioquia, Norte de Satander, Cauca, and Amazonas at five points throughout the epidemic. *a-p.* Each plot represents a different model, with model features labels on rows and columns. Plotted departments reflect differences in population, epidemic size, and geographic regions of Colombia and are denoted by point type. Point shape denotes department. Point color indicates time at which forecasts were generated and is visually denoted in inset plot and color bar. 1-, 2-, 3-, and 4-week ahead forecasts and observed incidence were aggregated for ease of comparison. Points are the median values and lines are the 50% credible interval. 1:1 line is in grey.



Figure S4: Relative difference between the prior estimates of the parameters and the posterior estimates at the final time point in the epidemic. Parameters include the R multiplier (k), reporting rate (ρ) , overdispersion parameter (ϕ) , R_t scalar (c), and the location (ψ) , shape (α) , and scale (ν) parameters for the skew normal distribution. Blue indicates posterior estimates were higher than prior estimates, red indicates posterior estimates were lower than prior estimates, and grey indicates no difference. Areas in white indicate the corresponding model does not use that parameter. To calculate the relative difference, we subtracted the prior estimates of parameters from the posterior estimates of parameters and divided the difference by the prior to standardize over different parameter magnitudes. For comparison purposes, we left out the initial timing and initial location of ZIKV introduction parameters.



Figure S5: Estimates of R at each assimilation week across temperatures for average mosquito occurrence probability in Colombia. Blue bands indicate 95% credible interval, with think navy line indicating the median estimate of R. Horizontal red line indicates R = 1.



Figure S6: Posterior distributions of model parameters for the spatially coupled models with a static R at five different points in time. The Rmultiplier (k), reporting rate (ρ) , the overdispersion parameter (ϕ) , and the timing of the first importations (ι_1) were model parameters represented in each model, although we are showing the posterior distribution only for a subset of the models. Violin plots are colored by time at which the forecasts were made, and correspond to time points in Fig. S3 and Fig. S23, S10



Figure S7: R_d multiplier versus the overdispersion parameter in the negative binomial distribution. The overdispersion parameter in the negative binomial distribution represents the variability in the reporting probability. Pearson's correlation coefficient for the two parameters within a parameter set is listed in the top right corner. This example is from the model using CDR-derived mobility data, static R, and one importation event, corresponding to Fig. S1.



Figure S8: R_d multiplier versus the reporting probability. Pearson's correlation coefficient for the two parameters within a parameter set is listed in the top right corner. This example is from the model using CDR-derived mobility data, static R, and one importation event, corresponding to Fig. S1.



Figure S9: The negative binomial overdispersion parameter versus the reporting probability for one example model at each assimilation period. The overdispersion parameter in the negative binomial distribution represents the variability in the reporting probability. Pearson's correlation coefficient for the two parameters within a parameter set is listed in the top right corner. This example is from the model using CDR-derived mobility data, static R, and one importation event, corresponding to Fig. S1.



Figure S10: Correlation of model parameters across within a parameter set through time for the spatially coupled models with a static R at five different points in time. Parameters include, the R multiplier (k), reporting rate (ρ) , the overdispersion parameter (ϕ) , and the timing of the first importations (ι_1) as they were represented in each model, although we are showing the correlations only for a subset of the models. Time points shown here correspond to those time points in Fig. S3 and Fig. S23, S6.



Figure S11: Cumulative observed versus forecasted incidence at 1-, 2-, 3-, 4- week ahead intervals for each individual model aggregated nationally at each point throughout the epidemic. a-p. Each plot represents a different model, with model feature labels on rows and columns. Point color indicates time at which forecasts were generated. Points are the median values and lines are the 50% credible interval. 1:1 line is in grey.



Figure S12: Observed incidence with initial median forecast for each of the 16 models with no weeks of data yet assimilated forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change. With the forecasts generally being at zero cases in the majority of the departments, we see that models were unlikely to forecast an epidemic to occur when no data was yet to be assimilated.



Figure S13: Observed incidence with median forecast for each of the 16 models with 12 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S14: Observed incidence with median forecast for each of the 16 models with 24 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S15: Observed incidence with median forecast for each of the 16 models with 36 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S16: Observed incidence with median forecast for each of the 16 models with 48 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S17: Ensemble weights of each model at each assimilation period. Ensemble weights were calculated using the expectation-maximization algorithm on short-term forecast performances, where forecast performance was assessed against 4-wk ahead incidence in a given department and nationally.

Figure S18: Log scores of each model and department at each assimilation period. Log scores were assessed using cumulative 4-wk ahead forecasts for each department and aggregated nationally. These log scores were then used in the expectation-maximization algorithm. L Least accurate Most accurate 4-wk ahead forecast log scores Майолај ви Санов Тобита Паналики Полагиа Калалики Мала Калалики Калалики Калалики Калалики Калалики Калалики Калалики Санова С Санова Санова Санова С Санова С С Санова С С С С С С С С С С С 16 36 56 Œ c Wational Mational Mational Mational Mational Material Computer Scientification (Computer Material Mate 12 32 52 τ Мазіолаі Мазіолаі Галіпаліек Галіпаліек Мент Мент Мент Мент Малакса Мент Малакса Соскобо Калано Соскобо Соскобо Калано Соскобо Соскоб Соскобо Соскоб Соскобо Соскоб Соскоб Соскоб Соскоб Соскобо Соскоб Соскобо Соскобо Соскоб Соскобо Соскобо Соскобо Соскобо Соскобо Соскобо Соскобо Соскобо Соскобо Соскоб Соскоб Соскобо Соскоб Соскоб Соскоб Соскоб Соскоб Соскоб Соскобо Соскоб С Соскоб С Соскоб С Соскоб Соскоб Соскоб С Соскоб С Соскоб С С 28 48 œ ε c C 24 4 C σ 4 20 ര _____ _____ пп
 R-1--madiation

 RP-1-madiation

 RP-1-madiation

 RP-2-madiation

 RP-2-audiation

 RP-2-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-1-audiation

 RP-2-audiation

 RP-2-audiation

 RP-1-audiation

 Rt-2-granky Rt-1-granky Rt-2-granky Rt-2-granky Rt-2-comparial Rt-2-comparial Rt-2-comparial Rt-2-comparial Rt-1-comparial Rt-1-comparial Rt-1-comparial Rt-2-gravity Rt-1-gravity R-2-gravity R-1-gravity R-1-gravity -2-rocrepatial Rt-2-radiation Rt-1-radiation R-2-radiation R-1-radiation Rt-2-CDRs Rt-1-CDRs R-2-CDRs R-1-CDRs



Number of weeks since first case reported in Colombia

Figure S19: Forecasts for a spatially coupled (green) and uncoupled (yellow) models with one introduction and static R for three departments. Navy bars denote incidence data. Large bands denote the 75% CrI, darker band denotes the 50% CrI, and thick line denotes the median forecast. Each row is from the same time point. Time points were chosen to be equally spaced out through the epidemic, with the first set of forecasts from the week of the first case, the second set of forecasts generated at 24 weeks after the first case was reported in Colombia, and the third set of forecasts generated at 48 weeks after the first case was reported in Colombia.



Figure S20: Fitted R and forecasts in two departments for models with two ZIKV introductions, cell phone mobility informed human movement, and a dynamic or static R. Models with dynamic R are depicted in red and models with a static R are depicted in blue. In plots of R, each line is a draw from the posterior, with a bold median line; horizontal black line depicts R = 1. The first set of forecasts in the middle column are from the peak week in both Risaralda and Córdoba, 24-28 weeks after the first case was reported in Colombia, and the second set are from 44-48 weeks after the first case was reported in Colombia. Vertical grey bars depict the forecasts and data considered when assessing short-term forecast performance.



Figure S21: Incidence by department with temperature trends. Blue bars denote weekly Zika incidence and red line denotes temperature trends.



Figure S22: Incidence by department with mosquito occurrence probability trends. Blue bars denote weekly Zika incidence and green line denotes mosquito occurrence trends.



Figure S23: Observed versus forecasted incidence at 1-, 2-, 3-, and 4week ahead intervals for EM-weighted and equally weighted ensemble models for Antioquia, Norte de Santander, Cauca, and Amazonas. Plotted departments reflect differences in population, epidemic size, and geographic regions of Colombia and are denoted by point type. Point shape denotes department. Point color indicates time at which the forecast was made (visually denoted in inset plot and color bar). Point is the median value and lines are the 50% credible interval. 1:1 line is in grey. The root mean square errors for the EM-weighted and equally weighted forecasts shown here are 0.705 and 0.640, respectively.



Figure S24: Magnitude of the 50% uncertainty bounds (as shown in Fig. S23) for 1-, 2-, 3-, and 4- week ahead forecasts in five different data assimilation periods for the EM-weighted and equally weighted ensemble models for Amazonas, Antioquia, Cauca, and Norte de Santander. The four points per data assimilation period represent the 1-, 2-, 3-, and 4- week ahead forecasts for each of the four departments denoted by color. Smoothed loess lines are shown to demonstrate how the magnitude of uncertainty changes through time for each department.



Figure S25: Observed incidence with initial expectation maximization algorithm-weighted ensemble forecast with no weeks of data yet assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S26: Observed incidence with expectation maximization algorithmweighted ensemble forecast with 12 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S27: Observed incidence with expectation maximization algorithmweighted ensemble forecast with 24 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S28: Observed incidence with expectation maximization algorithmweighted ensemble forecast with 36 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S29: Observed incidence with expectation maximization algorithmweighted ensemble forecast with 48 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



Figure S30: Observed incidence with initial equally weighted ensemble forecast with no weeks of data yet assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



12 weeks of data assimilated into model

Figure S31: Observed incidence with equally weighted ensemble forecast with 12 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



24 weeks of data assimilated into model

Figure S32: Observed incidence with equally weighted ensemble forecast with 24 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



36 weeks of data assimilated into model

Figure S33: Observed incidence with equally weighted ensemble forecast with 36 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.



48 weeks of data assimilated into model

Figure S34: Observed incidence with equally weighted ensemble forecast with 48 weeks of data assimilated into forecasting models. Navy bars indicate Zika incidence data at weekly time interval [14]. Light green band denotes 75% credible interval, darker green band denotes 50% credible interval, and the dark green line denotes median ensemble forecast. Vertical line indicates the point at which the forecast was made, with data to the left of the line assimilated into the model. Forecasts to the right of the vertical line change as more data is assimilated into the model, while model fits to the left of the vertical line do not change.

138 References

- [1] Sebastian Funk et al. "Comparative Analysis of Dengue and Zika Outbreaks Reveals Differences by Setting and Virus". In: *PLOS Neglected Tropical Diseases* 10.12 (Dec. 2016), pp. 1–16. DOI: 10.1371/journal.pntd.0005173.
 URL: https://doi.org/10.1371/journal.pntd.0005173.
- [2] Samir Bhatt et al. "The global distribution and burden of dengue". In: Nature 496.7446 (Apr. 2013), pp. 504-507. ISSN: 1476-4687. DOI: 10.1038/ nature12060. URL: https://doi.org/10.1038/nature12060.
- [3] Monaise M.O. Silva et al. "Accuracy of dengue reporting by national surveillance system, Brazil". In: *Emerging Infectious Diseases* 22.2 (Feb. 2016). DOI: https://dx.doi.org/10.3201/eid2202.150495. URL: https://wwwnc. cdc.gov/eid/article/22/2/15-0495_article.
- [4] Deborah P. Shutt et al. "Estimating the reproductive number, total outbreak
 size, and reporting rates for Zika epidemics in South and Central America".
 In: *Epidemics* 21 (Dec. 2017), pp. 63-79. ISSN: 1755-4365. URL: http://www.
 sciencedirect.com/science/article/pii/S1755436517300257.
- [5] Rachel J Oidtman, Guido España, and T Alex Perkins. "Co-circulation and misdiagnosis led to underestimation of the 2015-2017 Zika epidemic in the Americas". In: medRxiv (2020). DOI: https://doi.org/10.1101/19010256.
 URL: https://www.medrxiv.org/content/10.1101/19010256v1.
- [6] Sean M. Moore et al. "Leveraging multiple data types to estimate the size of the Zika epidemic in the Americas". In: *PLOS Neglected Tropical Diseases*14.9 (Sept. 2020), pp. 1–25. DOI: 10.1371/journal.pntd.0008640. URL: https://doi.org/10.1371/journal.pntd.0008640.
- M. U. G. Kraemer et al. "Utilizing general human movement models to predict the spread of emerging infectious diseases in resource poor settings". In: Scientific Reports 9.1 (2019), p. 5151. ISSN: 2045-2322. DOI: 10.1038/s41598-019-41192-3.
- [8] T. Alex Perkins et al. "Model-based projections of Zika virus infections in childbearing women in the Americas". In: *Nature Microbiology* 1.9 (2016), p. 16126. DOI: 10.1038/nmicrobiol.2016.126. URL: https://doi.org/10. 1038/nmicrobiol.2016.126.
- Erin A. Mordecai et al. "Detecting the impact of temperature on transmission of Zika, dengue, and chikungunya using mechanistic models". In: *PLOS Neglected Tropical Diseases* 11.4 (Apr. 2017), pp. 1–18. DOI: 10.1371/journal.
 pntd.0005568. URL: https://doi.org/10.1371/journal.pntd.0005568.
- Moritz UG Kraemer et al. "The global distribution of the arbovirus vectors *Aedes aegypti* and *Ae. albopictus*". In: *eLife* 4 (June 2015), e08347. ISSN: 2050-084X. DOI: 10.7554/eLife.08347. URL: https://doi.org/10.7554/eLife.
 08347.
- [11] Oliver J. Brady et al. "Global temperature constraints on Aedes aegypti and Ae. albopictus persistence and competence for dengue virus transmission". In: *Parasites & Vectors* 7.1 (July 2014), p. 338. ISSN: 1756-3305. DOI: 10.1186/ 1756-3305-7-338. URL: https://doi.org/10.1186/1756-3305-7-338.
- [12] Miranda Chan and Michael A. Johansson. "The Incubation Periods of Dengue Viruses". In: *PLOS ONE* 7.11 (Nov. 2012), pp. 1–7. DOI: 10.1371/journal.
 pone.0050972. URL: https://doi.org/10.1371/journal.pone.0050972.

[13] Allison Black et al. "Genomic epidemiology supports multiple introductions and cryptic transmission of Zika virus in Colombia". In: *BMC Infectious Diseases* 19.1 (2019), p. 963. DOI: 10.1186/s12879-019-4566-2. URL: https: //doi.org/10.1186/s12879-019-4566-2.

[14] Amir S. Siraj et al. "Spatiotemporal incidence of Zika and associated environmental drivers for the 2015-2016 epidemic in Colombia". In: Scientific Data 5.1 (2018), p. 180073. DOI: 10.1038/sdata.2018.73. URL: https://doi.org/10.1038/sdata.2018.73.