

Untangling the dynamics of persistence and colonization in microbial communities

Sylvia L. Ranjeva^{1,5}, Joseph R. Mihaljevic^{*1,4,5}, Maxwell B. Joseph², Anna R. Giuliano³,
and Greg Dwyer¹

¹Department of Ecology and Evolution, University of Chicago, Chicago, IL 60637

²Earth Lab, University of Colorado, Boulder, CO 80303

³Center for Immunization and Infection in Cancer Research (CIIRC), Moffitt Cancer
Center & Research Institute, Tampa, FL 33612

⁴Current Address: School of Informatics, Computing, and Cyber Systems, Northern
Arizona University, Flagstaff, AZ 86011

⁵These authors contributed equally to this work.

Running title: Community interactions through time

*Corresponding Author: Joseph R. Mihaljevic; School of Informatics, Computing, and Cyber Systems, 1295 S. Knoles Drive, Flagstaff, AZ 86011; joseph.mihaljevic@nau.edu; 928-523-5125

Abstract

1
2 A central goal of community ecology is to infer biotic interactions from observed distributions
3 of co-occurring species. Evidence for biotic interactions, however, can be obscured by shared
4 environmental requirements, posing a challenge for statistical inference. Here we introduce a
5 dynamic statistical model that quantifies the effects of spatial and temporal covariance in lon-
6 gitudinal co-occurrence data. We separate the fixed pairwise effects of species occurrences on
7 persistence and colonization rates, a potential signal of direct interactions, from latent pairwise
8 correlations in occurrence, a potential signal of shared environmental responses. We apply our
9 modeling approach to a pressing epidemiological question by examining how human papillo-
10 mavirus (HPV) types coexist. Our results suggest that while HPV types respond similarly to
11 common host traits, direct interactions are sparse and weak, so that HPV type diversity depends
12 largely on shared environmental drivers. Our modeling approach is widely applicable to micro-
13 bial communities and provides valuable insights that should lead to more directed hypothesis
14 testing and mechanistic modeling.

Introduction

15
16 A fundamental goal of community ecology is to understand how interactions between species
17 in a shared environment shape observed patterns of diversity over time. A key challenge in un-
18 derstanding community turnover is to disentangle effects of environmental drivers of species co-
19 occurrence from inter-species interactions, especially when the goal is to infer these mechanisms
20 from observational data [1, 2]. This challenge is also found in epidemiology, in which a major
21 goal is to understand the factors that allow pathogens to coexist [3]. As is the case with free-living
22 species, when determinants of environmental niches are shared among pathogen types, inferring
23 interactions is difficult [4]. Understanding the mechanisms of microbial community turnover
24 thus presents an ecological, statistical, and computational challenge, especially considering the

25 size of microbial and pathogen data sets [5, 6]. Ecological models of community turnover that
26 account for shared environmental drivers are thus important for understanding mechanisms that
27 underlie pathogen diversity.

28 For macroscopic organisms, null model analysis has historically been used to infer potential
29 species interactions from observational data sets, through the identification of statistically non-
30 random aggregations of species across multiple habitats [7, 8, 1, 9]. Similar approaches have been
31 used to develop computationally efficient algorithms that make it possible to infer large corre-
32 lation networks from microbial sequence data [5, 10]. Disentangling the simultaneous effects of
33 species interactions and environmental filters from survey data is nevertheless a challenge for
34 analyses of both macroscopic and microscopic communities [11, 2]. For example, highly mobile,
35 competing species should transiently aggregate in habitats with shared resources, even if com-
36 petitive exclusion is expected at equilibrium. Snap-shot surveys of co-occurrence can therefore
37 lead to biased interpretations of species interactions, but time-series data can help overcome this
38 problem.

39 In the microbial ecology literature, network inference models have only rarely been adapted to
40 incorporate time-series data from multiple localities. Available methods include local similarity
41 analysis [12, 11, 13] and generalized Lotka-Volterra modeling [14, 15]. While local similarity anal-
42 ysis can be used with incidence data, Lotka-Volterra modeling requires measures of abundance,
43 which are notoriously difficult to infer from sequence data, whereas relative abundances can
44 bias statistical analyses [16]. Local similarity analysis can infer microbial networks from observa-
45 tions of time-delays and temporal correlations between microbes and environmental covariates,
46 but it relies on multiple, independent tests with p-value corrections, instead of an integrated
47 analysis [12, 13]. Joint species distribution models provide a more comprehensive method for
48 identifying putatively interacting species from static ecological survey data, while accounting
49 for shared environmental drivers [17, 18, 19, 20, 21, 22]. These models use logistic regression
50 to estimate how environmental covariates affect species occupancy probabilities across a hetero-
51 geneous landscape. Species interactions are then inferred from residual correlations between

52 species occurrences. While joint-species models can generate hypotheses about static community
53 assemblages, most methods fail to capture important drivers of co-occurrence that emerge from
54 dynamic properties of the community dynamics [2]. For example, species co-occurrence may
55 be positively correlated across heterogeneous habitats, because of shared resources, but nega-
56 tively correlated across time, because of negative species interactions within sites (i.e. Simpson's
57 paradox, fig. 1).

58 Here we extend the joint-species modeling framework to infer more complex, biologically
59 realistic dynamics in a way that is computationally tractable for large microbial data sets. We
60 develop a statistical model of a dynamic, multi-species metacommunity in which species are
61 affected by each other's persistence and colonization probabilities, and by shared environmental
62 drivers. This approach can be readily applied to pathogenic microbe populations, in which
63 distinct pathogen types represent species coexisting within a heterogeneous landscape of host
64 organisms. In our method, we model correlations in species occupancy across habitats and across
65 time, resolving Simpson's paradox and accounting for latent environmental covariates. We also
66 estimate pairwise species effects on rates of colonization and persistence. Using synthetic data,
67 we demonstrate the ability of our model to accurately and precisely infer dynamics consistent
68 with Simpson's paradox, even with sparse occurrences. We then apply our model to data on
69 human papillomavirus (HPV), a pathogen of significant public health concern.

70 Human papillomavirus (HPV) is the most common sexually transmitted infection and a ma-
71 jor cause of cervical, genital, and oropharyngeal cancers, and it consists of over 200 types [23].
72 Uncertainty about the mechanisms underlying HPV type coexistence, and particularly about po-
73 tential HPV type interactions, reflects a crucial unknown. Four HPV types cause most disease
74 symptoms [24, 23, 25] and quadrivalent vaccination has demonstrated high efficacy in reducing
75 rates of cervical dysplasia and genital warts [26, 27]. A recent 9-valent HPV vaccine targets ad-
76 ditional oncogenic types [28]. Because the HPV vaccine is multivalent, it is possible that type
77 replacement will occur, in which non-vaccine types increase in frequency due to population-level
78 removal of vaccine-targeted types [29]. Type replacement following vaccination depends on inter-

79 actions between HPV types during natural infection, and particularly on inter-type competition
80 through cross-immunity [30]. Understanding the ecological mechanisms that underlie HPV type
81 diversity could therefore inform strategies for disease management and prevention. It has thus
82 far been difficult to distinguish HPV type interactions from the effects of shared host-specific risk
83 factors. Our dynamical community model allows us to investigate how type interactions and risk
84 factors together structure the HPV viral community.

85 In this study we address two questions, which differ in their scope. First, we use our full
86 model to ask which interactions between specific HPV types warrant future investigation? Sec-
87 ond, we ask a more ecological question: what are the dominant drivers of community compo-
88 sition across space and time? To address this second question, we build models of increasing
89 complexity, and we use model selection to determine whether HPV community patterns are
90 determined by putative interactions between HPV types, by host-level factors that determine
91 HPV distributions, or both. Our full model identified several interactions that warrant further
92 experimental investigation, including negative pairwise effects on persistence and colonization
93 probabilities. In addition, there is a strong signal of shared environmental drivers among HPV
94 types, highlighting the importance of host-specific risk factors in supporting coexistence. By
95 comparing models of varying complexity, however, we show that the dynamics of the HPV
96 community are most parsimoniously explained by shared environmental drivers, rather than
97 by strong pairwise interactions between HPV types. Pairwise species interactions thus do not
98 appear to drive community-wide patterns of co-occurrence in the HPV community. Our study
99 demonstrates the ability of our joint-species models to quickly and efficiently infer properties of
100 a large, real-world viral community, and the model could therefore be of broad usefulness in
101 understanding microbial communities.

Materials and Methods

HPV natural history

HPV types are classified based on the L1 viral capsid protein. A distinct HPV type is a variant whose L1 gene sequence is at least 10% dissimilar from any other HPV type [31]. The transmission and coexistence of individual HPV types depend on traits and risk factors of individual hosts [32, 33, 34, 35]. These include determinants of sexual behavior, including frequency of condom use, number of new and steady sexual partners, and sexual orientation; demography, including race and ethnicity; and non-sexual behavior, including smoking and alcohol consumption.

Interactions between HPV types could determine HPV diversity, though conclusive evidence of HPV type interactions is lacking [36, 37, 30]. As in any species, HPV type interactions may be synergistic, neutral, or competitive. Synergism occurs when one type facilitates infection by another, while competition occurs when one type prevents infection by another. Under competitive interactions, removal of one HPV type should lead to an increase in prevalence of the competing type in the host population, resulting in type replacement. Natural history surveys reporting elevated odds ratios for multiple to single infections with HPV have suggested that cross-immunity among HPV types is unlikely [38, 39, 40, 41]. Additionally, the genetic stability of HPV as a double-stranded DNA virus has been used to support arguments against the possibility of type replacement [42], on the grounds that rapid emergence of antigenic variants is unlikely [27]. Nevertheless, a recent increase in prevalence of non-vaccine types was found in young women following vaccination and in the United States [36], suggesting that type replacement may be occurring. Indeed, several models of HPV type interactions indicate that competition between HPV types is plausible under observed patterns of coinfections [43, 30] and have demonstrated the possibility of type-replacement after vaccination [43, 30, 44, 45].

Data

125

126 We fit models of HPV type dynamics to data from the HPV Infection in Men (HIM) study
127 [32, 33, 46], a multinational cohort study of HPV infection in men with no prior diagnosis of
128 genital cancer or other sexually transmitted infections. The HIM study enrolled over 4000 men
129 between 2005 and 2009 from three cities: Tampa, Florida, USA; Cuernavaca, Mexico; and Sao
130 Paulo, Brazil. Detailed study methods are described elsewhere [32]. Briefly, the HIM study
131 tracked PCR-confirmed infections with 37 types of HPV in men over a mean of 5 years of follow-
132 up, recording behavioral and demographic information for all participants. The data for each
133 individual consist of binary time series describing infection status with respect to each type over
134 a median of 10 clinic visits, at median intervals of 6.0 months (variance = 0.7 months).

135 For the present analysis, we included the 3656 participants with no reported diagnosis of
136 HIV and PCR samples for each HPV type at all clinic visits (see Appendix). We limited our
137 analysis to ten of the HPV types available in the HIM dataset: the nine HPV types included in
138 the most recent HPV vaccine [28]) and HPV84, a type that has shown high prevalence in several
139 studies among men [23, 47]. Of the ten types analyzed, seven oncogenic or high-risk types -
140 HPV16, HPV18, HPV31, HPV33, HPV45, HPV52, and HPV58 - have a demonstrated connection
141 to cervical cancer, while three nononcogenic or low-risk types - HPV6, HPV11, and HPV84 - have
142 been implicated in benign anogenital lesions [48]. Overall, our study includes 30,525 data points:
143 one point per patient per virus type per visit.

144

Statistical Model

145 Our goal is to extend current joint-species modeling techniques to biological processes that may
146 be needed to understand community dynamics. Currently, only a limited number of joint-species
147 modeling techniques are available for longitudinal survey data. Sebastian-Gonzalez et al. [20] ex-
148 tended the joint-species modeling framework to allow for multiple community surveys through
149 time by modeling the fixed, pairwise effects of species co-occurrence between subsequent time

150 points. Dorazio [49] introduced a model that separately estimated rates of species colonization
151 and persistence from sequential community surveys. Although this latter model specifies the
152 processes of extinction and colonization that can explain occupancy dynamics over time, it does
153 not account for the residual dependence among species that can result from species interactions.
154 Here we describe a statistical model that is tailored to the repeated surveys of patients in the
155 HIM dataset, thereby combining the methods of Sebastian-Gonzalez et al. [20] and Dorazio [49]
156 in a computationally tractable way.

157 Our data consist of observations made in I patients, who can harbor up to J HPV types (in
158 our case limited to 10 types), sampled over a maximum of T sequential visits to the clinic. Ob-
159 servations of the HPV dataset are therefore aggregated as binary presence/absence data in the
160 $I \times J \times T$ incidence array \mathbf{Y} , such that $Y_{i,j,t}$ indicates the presence or absence of HPV type j in
161 patient i at visit t . Importantly, however, this model generalizes to metacommunities sampled
162 repeatedly through time. Specifically, the model structure is the same as considering a metacom-
163 munity made up of I discrete habitats or sites, which harbor up to J species from the regional
164 species pool, and that are surveyed over a maximum of T time points.

We fit a multivariate probit regression model to the binary presence/absence data in \mathbf{Y} , which has been used in other joint-species modeling approaches [21]. Probit regression relates a linear predictor to occupancy probabilities using a standard normal cumulative distribution function. In this model, the probability that a binary random variable is equal to one (i.e. $P(Y = 1)$) is equal to the probability that the latent variable z is greater than zero. The linear predictor μ completely determines the latent variable z and can be a function of one or more covariates and their effects. As part of the probit definition, the residual variance of z is equal to one. In general then, we are interested in understanding how linear predictors influence the probability that an HPV type occurs in a given patient. A generalized probit model with a single covariate x is

formulated for the i^{th} sample as:

$$\begin{aligned}
 Y &\in \{0, 1\}, \\
 P(Y_i = 1) &= P(z_i > 0), \\
 z_i &\sim N(\mu_i, 1), \\
 \mu_i &= \beta x_i.
 \end{aligned} \tag{1}$$

165 Our model extends the generalized probit model by assuming that occurrence probabilities
 166 are affected by both patient-level effects and potential interactions between HPV types. We
 167 therefore build upon the general case of the probit model (Eq. 1) to model observations of the
 168 dynamic HPV metacommunity. To account for temporal dynamics, we assume that the linear
 169 predictor $\mu_{i,j,t}$ for each observation depends on observation-specific probabilities of persistence
 170 and colonization:

$$\mu_{i,j,t} = \alpha_j + \mathbf{Y}_{i,1:J,t-1} \mathbf{B}_j^{(\phi)'} (Y_{i,j,t-1}) + \mathbf{Y}_{i,1:J,t-1} \mathbf{B}_j^{(\gamma)'} (1 - Y_{i,j,t-1}) + \epsilon_{patient_{i,j}} + \epsilon_{visit_{i,j,t}} \tag{2}$$

171 Here, α_j is an adjustment to account for among-type variation in commonness. The presence
 172 of a given HPV type can affect the probability of persistence or colonization of other types, with
 173 a one time-step lag. If HPV type j was present in patient i on the previous clinic visit ($t - 1$), then
 174 persistence effects are represented by the product $\mathbf{Y}_{i,1:J,t-1} \mathbf{B}_j^{(\phi)'}$, where $\mathbf{Y}_{i,1:J,t-1}$ is a row vector of
 175 length J containing the presence/absence states of strains $j = 1, \dots, J$ in patient i on the previous
 176 visit ($t - 1$), and $\mathbf{B}_j^{(\phi)'}$ is a column vector of length J containing pairwise interaction coefficients.
 177 These coefficients thus specify how HPV type composition at the previous visit affects persistence
 178 (ϕ) of type j . If type j was absent in patient i on visit $t - 1$, colonization effects are represented
 179 by the product $\mathbf{Y}_{i,1:J,t-1} \mathbf{B}_j^{(\gamma)'}$ ($1 - Y_{i,j,t-1}$), where $\mathbf{B}_j^{(\gamma)'}$ ($1 - Y_{i,j,t-1}$) is a column vector of length J ,
 180 again containing pairwise interaction coefficients. These coefficients thus specify how HPV type
 181 composition at the previous visit affects the colonization (γ) of type j . Both interaction matrices
 182 ($\mathbf{B}^{(\phi)}$ and $\mathbf{B}^{(\gamma)}$) are $J \times J$ dimensional, and $\mathbf{B}_j^{(\phi)}$ and $\mathbf{B}_j^{(\gamma)}$ represent the row vectors acquired by
 183 extracting row j .

Lastly, patient-level and visit-level adjustments are specified as $\epsilon_{patient_{i,j}}$ and $\epsilon_{visit_{i,j,t}}$, respectively. The multivariate patient-level random effect $\epsilon_{patient}$ allows pairwise correlations in HPV type occurrence across patients, thereby describing pairwise similarities in environmental requirements. In the case of the HIM data, $\epsilon_{patient}$ therefore controls for shared determinants of host risk, such as host behavioral covariates, that could confound estimates of HPV type interactions. The random visit-level effect ϵ_{visit} allows for pairwise correlations in HPV type occurrence across clinic visits that are not explained by the fixed temporal effects. $\epsilon_{patient}$ and ϵ_{visit} allow for residual pairwise correlations in co-occurrence that are not explained by the fixed, pairwise effects. Following the definition of the multivariate probit density, $\epsilon_{patient}$ and ϵ_{visit} are nested effects, such that the same $\epsilon_{patient}$ is added to all of that patient's visits, such that the variances of $\epsilon_{patient}$ and ϵ_{visit} must sum to one (i.e. $z \sim N(\mu, 1)$). These random effects are therefore structured as follows:

$$\begin{aligned}\epsilon_{patient} &\sim N(0, \Sigma_{patient}) \\ \epsilon_{visit} &\sim N(0, \Sigma_{visit}) \\ \sigma_{patient_j}^2 + \sigma_{visit_j}^2 &= 1\end{aligned}\tag{3}$$

184 where $\Sigma_{patient}$ and Σ_{visit} are $J \times J$ variance-covariance matrices, constrained so that the j^{th} variance
 185 parameters from the two matrices sum to one, for $j = 1, \dots, J$. Therefore, $\rho_{patient_{p,q}}$ represents the
 186 pairwise correlation between HPV types that is measured among patients, which is derived from
 187 the variance-covariance matrix $\Sigma_{patient}$. Then, $\rho_{visit_{p,q}}$ represents the pairwise correlation between
 188 HPV types that is measured between visits and within patients (i.e. longitudinally), which is
 189 derived from the variance-covariance matrix Σ_{visit} .

We also model fixed effects of the time between visits (TBV) on persistence and colonization, to allow for the variability in when patients visited the clinic. The median TBV was 6.0 months with variance = 0.7 months, which we centered and scaled for use in the model. We allowed for fixed effects of TBV on the HPV type-specific probability of persisting ($\beta_j^{(TBV,\phi)}$) and the probability of colonizing ($\beta_j^{(TBV,\gamma)}$). We hypothesized that the probability that an HPV type colonizes a patient increases with TBV, due to a longer period of risk, while the probability that

a HPV type persists in the patient decreases with TBV, due to a longer time in which clearance may occur. The structure of these fixed effects is:

$$\mathbf{Z}_{i,j,t-1}\beta_j^{(\text{TBV},\phi)}(Y_{i,j,t-1}) + \mathbf{Z}_{i,j,t-1}\beta_j^{(\text{TBV},\gamma)}(1 - Y_{i,j,t-1}) \quad (4)$$

190 In this formula, \mathbf{Z} is an $I \times T$ matrix that holds the centered and scaled values of TBV for
191 each patient. This formula is added to μ_{ijt} .

192 *Model inference*

193 We coded our Bayesian model in *Stan* [50], an efficient, generalizable, statistical programming
194 language, which employs adaptive Hamiltonian Monte Carlo (HMC) for model inference. We
195 used vague priors for all parameters, although as mentioned earlier, we constrained the patient-
196 and visit-level standard deviations to sum to one, to conform to the definition of the multivariate
197 probit. We also included priors on the HPV type-specific, baseline probabilities of occurrence,
198 α_j , that allowed us to assume that all types are rare across patients and clinic visits. Indeed, the
199 most common type, HPV84, was still only present in 8.3% of all observations.

200 *Testing the model with synthetic data*

201 Using synthetic data, we tested the ability of our model to: (1) infer dynamics consistent with
202 Simpson's Paradox, meaning opposite correlations in among-patient effects versus among-visit
203 effects, (2) infer dynamics given observations of rare species, reflective of the HIM data, and
204 (3) infer weak inter-species interactions, as are likely in nature. We generated a synthetic data
205 set roughly half the size of the HIM data set to demonstrate the ability of our model to correctly
206 estimate model parameters from a sparser data set. We therefore simulated data for a community
207 of ten hypothetical pathogen strains sampled in 1500 patients, in which each patient was sampled
208 10 times. We assumed low but variable baseline probabilities of occurrence for each strain, with
209 the baseline occurrence set to the baseline prevalence of the ten least prevalent HPV types in the

210 HIM dataset. We further assumed positive patient-level correlations and negative observation-
211 level correlations, such that correlations were equal across pathogen strain pairs ($\rho_{patient_{p,q}} = 0.5$,
212 $\rho_{visit_{p,q}} = -0.1$). Pairwise effects on persistence and colonization $\beta_{p,q}^{\phi}$ and $\beta_{p,q}^{\gamma}$ were drawn from
213 normal distributions. All of our code for generating the synthetic data, as well as the data set
214 itself, is available in our open-source repository https://bitbucket.org/jrmihalj/hpv_jsdm.

215 *Fitting the model to the HIM data*

216 Our first goal was to use our model to identify any interactions between HPV types that might
217 warrant future epidemiological investigations. We therefore fit our full model and quantified the
218 posterior distributions of the pairwise effects of HPV types on colonization and persistence rates.
219 Our second goal was to understand the relative contributions of environmental effects, such
220 as host-specific risk factors, and pairwise inter-type interactions to HPV community dynamics.
221 We therefore fit four nested models of varying complexity. Model 1 has fixed, pairwise effects
222 between HPV types, model 2 has residual correlations that account for environmental effects, and
223 model 3, our full model, has both. Model 4 includes only baseline occurrence probabilities α_j ,
224 and is therefore a type of null model. All of these models include the effects of the time between
225 visits (TBV). We then compared the models' out-of-sample predictive abilities using the leave-
226 one-out information criterion (LOO-IC), estimated using Pareto-smoothed importance sampling
227 in the R package "loo" [51]. Compared to the Watanabe-Akaike information criterion (WAIC),
228 which is asymptotically equal to LOO-IC, the LOO-IC has been found to be more robust when
229 using vague priors [52], as in our models. We considered two models to be substantially different
230 if their LOO-IC values differed by 3, which is the common convention [53]. In practice, for a data
231 set this large, small changes in overall goodness-of-fit could lead to very large changes in the
232 likelihood when integrated across the many data points, and thus large differences in LOO-IC.
233 We therefore emphasize that we use this model selection procedure as a heuristic to guide our
234 understanding of community dynamics, rather than as a robust hypothesis test.

Results

Model validation with synthetic data

When we tested our model with synthetic data, it accurately and precisely inferred dynamics consistent with Simpson's Paradox, even when the data were sparse (Fig. 2). The model correctly inferred the low baseline probabilities of species occurrence (Fig. 2 A) and all patient-level correlations (Fig. 2 B). It also accurately estimated the majority of negative correlations at the observation level, although some inferred pairwise correlations were indistinguishable from zero (Fig. 2 C). This latter effect was not surprising, because we assumed a weak negative correlation ($\rho_{visit} = -0.1$). Importantly, although the model's estimates of the magnitude of simulated correlations were sometimes incorrect, the model was unbiased with respect to the direction of the simulated correlations. The model also correctly estimated persistence ($\beta_{p,q}^{\phi}$) and colonization ($\beta_{p,q}^{\gamma}$) under both strong and weak interactions (Fig. 2 D,E). Finally, the model accurately recovered the effects of the time between visits on both persistence and colonization probabilities, which we assumed were the same for all pathogen strains (Fig. S2).

Metacommunity dynamics of HPV and model comparisons

In our full model, there were only a few interactions between HPV types that were worthy of future investigation, including several weakly negative effects on colonization probability (Fig. 3). Importantly, including these fixed effects and the random effects of patient-level and observation-level correlations led to a substantial improvement relative to a null model that accounts only for type-specific baseline occurrence probabilities, suggesting that the biology added to our model helps explain HPV community composition relative to the null model (Table 1). Based on LOO-IC selection, however, the most parsimonious model included only the random effects of patient-level and observation-level correlations, without pairwise interactions between the HPV types (Table 1). Pairwise inter-type interactions can thus be identified by our model, but the effect of

259 these interactions is not strong enough to substantially mediate the overall community compo-
260 sition in this subset of 10 HPV types. The best model, which did not include these pairwise
261 interactions, gives qualitatively similar insights for the random effects, meaning the patient-level
262 and observation-level correlations, as our full model (Fig. S4).

263 The best model captured important qualitative aspects of the HPV dynamics, as well. The
264 inferred baseline occurrence probability recovered the observed rank order of prevalence of the
265 ten HPV types (Fig. 3A). The model confirmed that increasing values of TBV had positive effects
266 on colonization probabilities ($\beta_j^{(TBV,\gamma)} > 0$) for all HPV types, but it had negative effects on
267 persistence probabilities ($\beta_j^{(TBV,\phi)} < 0$) for all but two HPV types (Figs. S3, S4).

268 Patient-level correlations were positive for all but one pair of HPV types (Fig. 3B). These
269 positive correlations suggest that there are shared environmental drivers across human hosts,
270 in the form of risk factors. In the case of HPV52 and HPV58 (Fig. 3C), there are both positive
271 patient-level and negative observation-level correlations. Positive observation-level correlations,
272 or correlations within individuals over time, likely signal affinity for co-transmission, because in
273 the models these effects are in addition to the pairwise effects on persistence and colonization.
274 Negative observation-level correlations thus signal reduced affinity for co-transmission. How-
275 ever, the negative observation-level correlations between HPV52 and HPV58 must be interpreted
276 with caution, as they could reflect the masking of HPV58 detection by HPV52, a problem that
277 has been documented in the linear array genotyping test used in the HIM study [42].

278 Discussion

279 Our results suggest that HPV type coexistence is strongly driven by shared environmental char-
280 acteristics. While the full model is able to estimate even sparse and weak (putative) interactions
281 between HPV types, our model selection procedure suggests that these interactions are not im-
282 portant for explaining overall patterns of community turnover in HPV. The influence of patient-
283 level correlations on HPV community dynamics suggests that HPV types segregate among hosts

284 with shared traits. It is therefore likely that human subpopulations exist that could promote HPV
285 type coexistence across space and time. This finding is consistent with epidemiological evidence
286 of type-specific differences in the risk factors that promote HPV transmission [54, 55], and with
287 another recent modeling study that characterized subtle differences in the profile of host-specific
288 risk factors that affect infection with each type [56].

289 Model selection shows that pairwise inter-type interactions that affect colonization and per-
290 sistence probabilities do not influence overall patterns of community turnover in this HPV data
291 set. However, the full model identified several putative interactions worthy of future epidemio-
292 logical investigations. In particular, it is possible that interactions could mediate the occurrence
293 patterns of specific pairs of HPV types, even though model selection suggests that pairwise in-
294 teraction effects have no meaningful effects on the HPV community dynamics as a whole. In
295 other words, the community-level patterns could swamp out the patterns of specific HPV pairs.
296 Further, by limiting our analysis to a subset of ten HPV types, it is possible that we by chance did
297 not include HPV types that have larger effects on the community. Also, our model only estimates
298 pairwise effects, and future studies could account for higher order interactions, which have been
299 shown to be important in diverse competitive networks [57].

300 The results of our analysis complement the results of a previous, mechanistic model of HPV
301 dynamics fitted to 6 HPV types of the HIM dataset [56]. The authors of this previous work
302 formulated an epidemiological model that allowed for homologous immunity, a form of within-
303 species competition, as well as the effects of 11 host-specific covariates. The best-fit version of this
304 model included no homologous immunity for any of the six HPV types (HPV84, HPV62, HPV89,
305 HPV16, HPV51, and HPV6), finding instead that previous infection with any type significantly
306 increases the risk of re-infection with the same type. In our statistical model, this effect is further
307 confirmed by the positive baseline persistence probabilities ($\beta_{p,p}^{\phi}$) across all ten HPV types ana-
308 lyzed. That study [56] also detected no pairwise interaction between two taxonomically similar
309 types, HPV16 and HPV31, which had been hypothesized to compete through cross-immunity
310 [58, 59]. Furthermore, the risk of initial infection with any HPV type was concentrated among

311 high-risk subpopulations, which were linked to host-specific covariates. Taken together, the re-
312 sults of this previous analysis [56] suggest that both intra-specific and inter-specific competition
313 are weak or absent in the HPV viral community, such that stabilizing competitive mechanisms
314 cannot explain HPV diversity. Instead, diversity may depend on sustained infection within high-
315 risk subpopulations specific to each HPV type. These findings are consistent with our finding
316 that inter-type interactions have little effect on HPV community dynamics (Table 1). Further-
317 more, by showing how host-specific traits define niches that are used by different HPV types, the
318 previous work [56] supports the importance of shared among-patient traits to explain patterns
319 of co-occurrence.

320 While the different quantitative approaches between the previous study [56] and our study
321 provide complementary results, there are important differences in the methods, applications,
322 and conclusions. Ranjeva et al. [56] tested mechanistic biological models about type-specific
323 HPV dynamics, whereas our approach allowed for the identification of statistical patterns in the
324 community dynamics of multiple types. Also, our method can be generalized to any metacom-
325 munity that is sampled through time, rather than being specific to a pathogen community that
326 interacts via cross-immunity, as modeled by Ranjeva et al. [56]. Indeed, our statistical frame-
327 work is agnostic to the specific mechanisms of interactions. Instead our model specifies latent
328 mechanisms that affect probabilities of persistence and colonization, which are estimated from
329 the occurrence data.

330 We have shown that a relatively simple statistical model can be used to infer community
331 dynamics, even in a system with rare species occurrences. Sparsity of observational data in real-
332 world metacommunities generally limits the power of statistical models to correctly infer ecolog-
333 ical effects [49, 60, 61]. We showed that our model can be used to infer opposing environmental
334 and temporal dynamics from communities of rare species, and to detect weak interactions among
335 rare species, which are the most common types of interactions in nature [62]. Inferring residual
336 correlations with rare species requires a substantial amount of data, but, in the age of affordable,
337 high-throughput sequencing technologies, such data can often be obtained easily. Moreover, our

338 model accounts for the effects of unobserved environmental drivers, specifically host-specific
339 risk-factors in the case of the HPV data, without having to specify covariates explicitly. This may
340 be useful for analyzing large microbial communities, such as microbiome communities, in which
341 the environmental drivers are unknown.

342 In classical joint-species distribution models, residual correlations in species occurrence are
343 used to infer species interactions, but such residual correlations can arise instead from shared
344 covariate responses that are not explicitly included in the model structure [21, 2]. Our model,
345 however, does not rely on residual correlations to infer interspecies interactions *per se*. We use
346 species occupancy at the previous time step to estimate lagged, pairwise effects of species' oc-
347 currences on the probabilities of persistence and colonization of cohabitating species. Residual
348 correlations in our models instead account for latent environmental covariates, such as unmea-
349 sured host-specific traits. Although our statistical modeling approach can thus identify impor-
350 tant signatures of species interactions, mechanistic models and experimentation are nevertheless
351 required to rigorously test hypotheses about species interactions. Furthermore, we estimate in-
352 terspecies effects on persistence and colonization using a one-timestep lag, which requires that
353 the timescale of the species interactions be equal to the timescale of observations. This assump-
354 tion may not always hold. Our method is therefore best used to refine testable hypotheses
355 from observed dynamics of large community assemblages, such as microbiome assemblages, in
356 a computationally-feasible manner, rather than as a final step in inferring interactions.

357 A final caveat is that our models do not allow for dynamics that occur between observa-
358 tions. Given two consecutive observations of a species, our models instead assume that there is
359 either persistence over the entire interval, or that at most one extinction or colonization has oc-
360 curred. This assumption may result in bias in communities that are poorly sampled relative to the
361 timescale of the dynamics. Indeed, recent evidence shows that standard joint-species distribution
362 modeling approaches cannot accurately capture simulated predator-prey dynamics, especially if
363 habitats are relatively homogeneous, probably because of non-linear dynamics [2]. This problem
364 is likely to be important for non-linear host-pathogen dynamics as well, and should be a subject

365 of future simulation efforts. Our dataset however spans a wide diversity of patients, and includes
366 the effects of the time between visits, which should limit this type of bias.

367 **Competing Interests**

368 The authors declare no competing financial interests

Literature Cited

369

370 [1] Gotelli NJ, Graves GR. Null Models in Ecology. Washington, D.C.: Smithsonian Institution
371 Press; 1996.

372 [2] Zurell D, Pollock LJ, Thuiller W. Do joint species distribution models reliably detect inter-
373 specific interactions from co-occurrence data in homogenous environments? *Ecography*.
374 InPress;.

375 [3] Keeling MJ, Pejman R. Formulating the Deterministic SIR Model. In: *Modeling Infectious*
376 *Diseases in Humans and Animals*. 2nd ed. Princeton, NJ: Princeton University Press; 2008. p.
377 15–52. Available from: https://books.google.com/books?hl=en&lr=&id=G8enmS23c6YC&oi=fnd&pg=PP2&dq=force+of+infection+keeling+and+rohani&ots=rFKSDv9toM&sig=hHVhDvYUpjrd4lkCgLe3bkH8m_4#v=onepage&q=forceofinfectionkeelingandrohani&f=false.
378
379
380

381 [4] Godsoe W, Franklin J, Blanchet FG. Effects of biotic interactions on modeled species' distri-
382 bution can be masked by environmental gradients. *Ecology and evolution*. 2017;7(2):654–664.

383 [5] Faust K, Raes J. Microbial interactions: From networks to models. *Nature Reviews Microbi-*
384 *ology*. 2012;10(8):538–550. Available from: <http://dx.doi.org/10.1038/nrmicro2832>.

385 [6] Seabloom EW, Borer ET, Gross K, Kendig AE, Lacroix C, Mitchell CE, et al. The commu-
386 nity ecology of pathogens: coinfection, coexistence and community composition. *Ecology*
387 *Letters*. 2015;18:401–415. Available from: <http://doi.wiley.com/10.1111/ele.12418>.

388 [7] Diamond JM. Assembly of species communities. *Ecology and evolution of communities*.
389 1975;p. 342–444.

390 [8] Connor EF, Simberloff D. The Assembly of Species Communities: Chance or Competition?
391 *Ecology*. 1979;60(6):1132. Available from: <http://www.jstor.org/stable/1936961?origin=crossref>.
392

- 393 [9] Gotelli NJ, McCabe DJ. Species co-occurrence: A meta-analysis of JM Diamond's assembly
394 rules model. *Ecology*. 2002;83(8):2091–2096.
- 395 [10] Fisher CK, Mehta P. Identifying keystone species in the human gut microbiome from
396 metagenomic timeseries using sparse linear regression. *PLoS ONE*. 2014;9(7):1–10.
- 397 [11] Weiss S, Van Treuren W, Lozupone C, Faust K, Friedman J, Deng Y, et al. Correlation
398 detection strategies in microbial data sets vary widely in sensitivity and precision. *ISME*
399 *Journal*. 2016;10(7):1669–1681. Available from: [http://dx.doi.org/10.1038/ismej.2015.](http://dx.doi.org/10.1038/ismej.2015.235)
400 235.
- 401 [12] Ruan Q, Dutta D, Schwalbach MS, Steele JA, Fuhrman JA, Sun F. Local similarity analy-
402 sis reveals unique associations among marine bacterioplankton species and environmental
403 factors. *Bioinformatics*. 2006;22(20):2532–2538.
- 404 [13] Xia LC, Ai D, Cram JA, Liang X, Fuhrman JA, Sun F. Statistical significance approxima-
405 tion in local trend analysis of high-throughput time-series data using the theory of Markov
406 chains. *BMC Bioinformatics*. 2015;16(1):1–14. Available from: [http://dx.doi.org/10.1186/](http://dx.doi.org/10.1186/s12859-015-0732-8)
407 [s12859-015-0732-8](http://dx.doi.org/10.1186/s12859-015-0732-8).
- 408 [14] Stein RR, Bucci V, Toussaint NC, Buffie CG, Räscht G, Pamer EG, et al. Ecological Modeling
409 from Time-Series Inference: Insight into Dynamics and Stability of Intestinal Microbiota.
410 *PLoS Computational Biology*. 2013;9(12):31–36.
- 411 [15] Carrara F, Giometto A, Seymour M, Rinaldo A, Altermatt F. Inferring species interactions in
412 ecological communities: A comparison of methods at different levels of complexity. *Methods*
413 *in Ecology and Evolution*. 2015;6(8):895–906.
- 414 [16] Cardona C, Weisenhorn P, Henry C, Gilbert JA. Network-based metabolic analysis and mi-
415 crobial community modeling. *Current Opinion in Microbiology*. 2016;31:124–131. Available
416 from: <http://dx.doi.org/10.1016/j.mib.2016.03.008>.

- 417 [17] Ovaskainen O, Hottola J, Siitonen J. Modeling species co-occurrence by multivariate logistic
418 regression generates new hypotheses on fungal interactions. *Ecology*. 2010;91(9):2514–2521.
419 Available from: <http://www.esajournals.org/doi/abs/10.1890/10-0173.1>.
- 420 [18] Ovaskainen O, Tikhonov G, Norberg A, Blanchet FG, Duan L, Dunson D, et al. How to
421 make more out of community data? A conceptual framework and its implementation as
422 models and software. *Ecology Letters*. 2017;20:561–576.
- 423 [19] Ovaskainen O, Abrego N, Halme P, Dunson D. Using latent variable models to identify large
424 networks of species-to-species associations at different spatial scales. *Methods in ecology
425 and evolution*. 2015;.
- 426 [20] Sebastián-González E, Sánchez-Zapata JA, Botella F, Ovaskainen O. Testing the heterospe-
427 cific attraction hypothesis with time-series data on species co-occurrence. *Proceedings of
428 the Royal Society B: Biological Sciences*. 2010;277:2983–2990.
- 429 [21] Pollock LJ, Tingley R, Morris WK, Golding N, O'Hara RB, Parris KM, et al. Understanding
430 co-occurrence by modelling species simultaneously with a Joint Species Distribution Model
431 (JSDM). *Methods in Ecology and Evolution*. 2014 may;5(5):397–406. Available from: <http://doi.wiley.com/10.1111/2041-210X.12180>.
- 433 [22] Joseph MB, Preston DL, Johnson PTJ. Integrating occupancy models and structural equation
434 models to understand species occurrence. *Ecology*. 2016;97(3):765–775.
- 435 [23] Ma Y, Madupu R, Karaoz U, Nossa CW, Yang L, Yooseph S, et al. Human papillomavirus
436 community in healthy persons, defined by metagenomics analysis of human microbiome
437 project shotgun sequencing data sets. *Journal of Virology*. 2014 5;88(9):4786–97. Available
438 from: <http://www.ncbi.nlm.nih.gov/pubmed/24522917>.
- 439 [24] Frazer IH. Interaction of human papillomaviruses with the host immune system: A well
440 evolved relationship. *Virology*. 2009;384(2):410–414. Available from: [http://dx.doi.org/
441 10.1016/j.virol.2008.10.004](http://dx.doi.org/10.1016/j.virol.2008.10.004).

- 442 [25] Todd RW, Roberts S, Mann CH, Luesley DM, Gallimore PH, Steele JC. Human papillo-
443 mavirus (HPV) type 16-specific CD8+ T cell responses in women with high grade vulvar
444 intraepithelial neoplasia. *International Journal of Cancer*. 2004;108(6):857–862.
- 445 [26] Dunne EF, Unger ER, McQuillan G, Swan DC, Patel SS, Markowitz LE. Prevalence of HPV
446 Infection. 2014;297(8):813–819.
- 447 [27] Markowitz LE, Hariri S, Lin C, Dunne EF, Steinau M, McQuillan G, et al. Reduction
448 in human papillomavirus (HPV) prevalence among young women following HPV vac-
449 cine introduction in the United States, National Health and Nutrition Examination Sur-
450 veys, 2003-2010. *The Journal of infectious diseases*. 2013 8;208(3):385–93. Available from:
451 <http://jid.oxfordjournals.org/cgi/content/long/jit192v1>.
- 452 [28] Joura EA, Giuliano AR, Iversen OE, Bouchard C, Mao C, Mehlsen J, et al. A 9-Valent HPV
453 Vaccine against Infection and Intraepithelial Neoplasia in Women. *New England Journal*
454 *of Medicine*. 2015 2;372(8):711–723. Available from: [http://www.nejm.org/doi/10.1056/](http://www.nejm.org/doi/10.1056/NEJMoa1405044)
455 [NEJMoa1405044](http://www.nejm.org/doi/10.1056/NEJMoa1405044).
- 456 [29] Bernard E, Pons-Salort M, Favre M, Heard I, Delarocque-Astagneau E, Guillemot D,
457 et al. Comparing human papillomavirus prevalences in women with normal cytology
458 or invasive cervical cancer to rank genotypes according to their oncogenic potential: a
459 meta-analysis of observational studies. *BMC infectious diseases*. 2013;13(1):373. Available
460 from: <http://www.ncbi.nlm.nih.gov/pubmed/23941096>[http://www.pubmedcentral.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3751808)
461 [nih.gov/articlerender.fcgi?artid=PMC3751808](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3751808)[http://www.pubmedcentral.nih.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3751808&tool=pmcentrez&rendertype=abstract)
462 [gov/articlerender.fcgi?artid=3751808&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3751808&tool=pmcentrez&rendertype=abstract).
- 463 [30] Durham DP, Poolman EM, Ibuka Y, Townsend JP, Galvani AP. Reevaluation of epidemio-
464 logical data demonstrates that it is consistent with cross-immunity among human papillo-
465 mavirus types. *Journal of Infectious Diseases*. 2012;206(8):1291–1298.

- 466 [31] Bernard HU, Burk RD, Chen Z, van Doorslaer K, Hausen Hz, de Villiers EM. Classification
467 of papillomaviruses (PVs) based on 189 PV types and proposal of taxonomic amendments.
468 Virology. 2010 5;401(1):70–79. Available from: [https://www.sciencedirect.com/science/
469 article/pii/S0042682210001005](https://www.sciencedirect.com/science/article/pii/S0042682210001005).
- 470 [32] Giuliano AR, Lazcano-Ponce E, Villa LL, Flores R, Salmeron J, Lee JH, et al.
471 The human papillomavirus infection in men study: human papillomavirus preva-
472 lence and type distribution among men residing in Brazil, Mexico, and the United
473 States. *Cancer Epidemiology, Biomarkers & Prevention*. 2008 8;17(8):2036–43. Available
474 from: [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3471778&tool=
475 pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3471778&tool=pmcentrez&rendertype=abstract).
- 476 [33] Giuliano AR, Lazcano E, Villa LL, Flores R, Salmeron J, Lee JH, et al. Circumcision
477 and sexual behavior: factors independently associated with human papillomavirus detec-
478 tion among men in the HIM study. *International Journal of Cancer*. 2009 3;124(6):1251–
479 7. Available from: [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=
480 3466048&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3466048&tool=pmcentrez&rendertype=abstract).
- 481 [34] Nyitray A, Carvalho R, Baggio M, Beibei L, Abrahamsen M, Papenfuss M, et al. Age-Specific
482 Prevalence of and Risk Factors for Anal Human Papillomavirus (HPV) among Men Who
483 Have Sex with Women and Men Who Have Sex with Men : The HPV in Men (HIM) Study.
484 *Journal of Infectious Diseases*. 2011;203(1):49–57.
- 485 [35] Nyitray AG, Carvalho Da Silva RJ, Baggio ML, Smith D, Abrahamsen M, Papenfuss M,
486 et al. Six-month incidence, persistence, and factors associated with persistence of anal
487 human papillomavirus in men: The HPV in men study. *Journal of Infectious Diseases*.
488 2011;204(11):1711–1722.
- 489 [36] Kahn Ja, Brown DR, Ding L, Widdice LE, Shew ML, Glynn S, et al. Vaccine-type human pa-
490 pillomavirus and evidence of herd protection after vaccine introduction. *Pediatrics*. 2012

491 8;130(2):249–56. Available from: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3408690&tool=pmcentrez&rendertype=abstract>.

493 [37] Wheeler CM, Castellsagué X, Garland SM, Szarewski A, Paavonen J, Naud P, et al. Cross-
494 protective efficacy of HPV-16/18 AS04-adjuvanted vaccine against cervical infection and
495 precancer caused by non-vaccine oncogenic HPV types: 4-year end-of-study analysis of the
496 randomised, double-blind PATRICIA trial. *Lancet Oncology*. 2012;13(1):100–110.

497 [38] Rousseau MC, Pereira JS, Prado JC, Villa LL, Rohan TE, Franco EL. Cervical coinfection with
498 human papillomavirus (HPV) types as a predictor of acquisition and persistence of HPV
499 infection. *Journal of infectious diseases*. 2001 12;184(12):1508–17. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/11740725>.

501 [39] Liaw KL. A prospective study of human papillomavirus (HPV) type 16 DNA detection by
502 polymerase chain reaction and its association with acquisition and persistence of other HPV
503 types. *Journal of Infectious Diseases*. 2001;183(1):8–15. Available from: <http://dx.doi.org/10.1086/317638>.

505 [40] Chaturvedi AK. Prevalence and Clustering Patterns of Human Papillomavirus Genotypes
506 in Multiple Infections. *Cancer Epidemiology, Biomarkers & Prevention*. 2005;14(10):2439–
507 2445. Available from: <http://cebp.aacrjournals.org/cgi/doi/10.1158/1055-9965.EPI-05-0465>.

509 [41] Chaturvedi AK, Engels EA, Pfeiffer RM, Hernandez BY, Xiao W, Kim E, et al. Human
510 Papillomavirus and Rising Oropharyngeal Cancer Incidence in the United States. *Journal*
511 *of Clinical Oncology*. 2011 11;29(32):4294–4301. Available from: <http://ascopubs.org/doi/10.1200/JCO.2011.36.4596>.

513 [42] Tota JE, Ramanakumar AV, Jiang M, Dillner J, Walter SD, Kaufman JS, et al. Epidemio-
514 logic approaches to evaluating the potential for human papillomavirus type replacement

- 515 postvaccination. *American journal of epidemiology*. 2013 8;178(4):625–34. Available from:
516 <http://www.ncbi.nlm.nih.gov/pubmed/23660798>.
- 517 [43] Elbasha EH, Galvani AP. Vaccination against multiple HPV types. *Mathematical Biosciences*.
518 2005;197(1):88–117.
- 519 [44] Murall CL, McCann KS, Bauch CT. Revising ecological assumptions about Human papillo-
520 mavirus interactions and type replacement. *Journal of Theoretical Biology*. 2014;350:98–109.
521 Available from: <http://dx.doi.org/10.1016/j.jtbi.2013.12.028>.
- 522 [45] Poolman EM, Elbasha EH, Galvani AP. Vaccination and the evolutionary ecology of human
523 papillomavirus. *Vaccine*. 2008;26(Supplement 3):25–30.
- 524 [46] Giuliano AR, Lee JH, Fulp W, Villa LL, Lazcano E, Papenfuss MR, et al. Incidence and clear-
525 ance of genital human papillomavirus infection in men (HIM): a cohort study. *Lancet*. 2011
526 3;377(9769):932–40. Available from: [http://www.pubmedcentral.nih.gov/articlerender.
527 fcgi?artid=3231998&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3231998&tool=pmcentrez&rendertype=abstract).
- 528 [47] Han JJ, Beltran TH, Song JW, Klaric J, Choi YS. Prevalence of Genital Human Papillomavirus
529 Infection and Human Papillomavirus Vaccination Rates Among US Adult Men. *JAMA On-
530 cology*. 2017;3(6):810–816. Available from: [http://oncology.jamanetwork.com/article.
531 asp?doi=10.1001/jamaoncol.2016.6192](http://oncology.jamanetwork.com/article.aspx?doi=10.1001/jamaoncol.2016.6192).
- 532 [48] Giuliano AR, Nyitray AG, Kreimer AR, Pierce Campbell CM, Goodman MT, Sudenga SL,
533 et al. EUROGIN 2014 roadmap: Differences in human papillomavirus infection natural
534 history, transmission and human papillomavirus-related cancer incidence by gender and
535 anatomic site of infection. *International Journal of Cancer*. 2015;136(12):2752–2760.
- 536 [49] Dorazio RM, Kéry M, Royle JA, Plattner M. Models for inference in dynamic metacommunity
537 systems. *Ecology*. 2010 8;91(8):2466–75. Available from: [http://www.ncbi.nlm.nih.
538 gov/pubmed/20836468](http://www.ncbi.nlm.nih.gov/pubmed/20836468).

- 539 [50] Carpenter B, Gelman A, Hoffman MD, Lee D, Goodrich B, Betancourt M, et al. *Stan*
540 : A Probabilistic Programming Language. *Journal of Statistical Software*. 2017;76(1). Avail-
541 able from: <http://www.jstatsoft.org/v76/i01/>.
- 542 [51] Vehtari A, Gelman A, Gabry J. loo: Efficient leave-one-out cross-validation and WAIC for
543 Bayesian models; 2016. R package version 1.1.0. Available from: [https://CRAN.R-project.](https://CRAN.R-project.org/package=loo)
544 [org/package=loo](https://CRAN.R-project.org/package=loo).
- 545 [52] Vehtari A, Gelman A, Gabry J. Practical Bayesian model evaluation using leave-one-out
546 cross-validation and WAIC. *Statistics and Computing*. 2016;27(5):1–20.
- 547 [53] Gelman A, Carlin J, Stern H, Dunson D. *Bayesian data analysis*. 2nd ed. Chap-
548 man and Hall; 2013. Available from: [http://books.google.com/books?hl=en&lr={&}id=ZXL6AQAAQBAJ&oi=fnd&pg=PP1&dq=Bayesian+data+analysis&ots=](http://books.google.com/books?hl=en&lr={&}id=ZXL6AQAAQBAJ&oi=fnd&pg=PP1&dq=Bayesian+data+analysis&ots=uNYiu1bGW5&sig=eZoeglTT1{ }60ozzh4Hrf5YQDvWE)
549 [uNYiu1bGW5&sig=eZoeglTT1{ }60ozzh4Hrf5YQDvWE](http://books.google.com/books?hl=en&lr={&}id=ZXL6AQAAQBAJ&oi=fnd&pg=PP1&dq=Bayesian+data+analysis&ots=uNYiu1bGW5&sig=eZoeglTT1{ }60ozzh4Hrf5YQDvWE).
- 551 [54] Albero G, Castellsagué X, Lin HY, Fulp W, Villa LL, Lazcano-Ponce E, et al. Male circum-
552 cision and the incidence and clearance of genital human papillomavirus (HPV) infection in
553 men: the HPV Infection in men (HIM) cohort study. *BMC Infectious Diseases*. 2014;14(1):75.
554 Available from: [http://www.scopus.com/inward/record.url?eid=2-s2.0-84893644201&](http://www.scopus.com/inward/record.url?eid=2-s2.0-84893644201&partnerID=tZ0tx3y1)
555 [partnerID=tZ0tx3y1](http://www.scopus.com/inward/record.url?eid=2-s2.0-84893644201&partnerID=tZ0tx3y1).
- 556 [55] Nyitray AG, Chang M, Villa LL, Carvalho RJ, Baggio ML. *us c t Ac ce p te d Ac ce p te d*
557 *us*. 2015;p. 1–26.
- 558 [56] Ranjeva SL, Baskerville EB, Dukic V, Villa LL, Lazcano-Ponce E, Giuliano AR, et al. Re-
559 curring infection with ecologically distinct HPV types can explain high prevalence and di-
560 versity. *Proceedings of the National Academy of Sciences of the United States of Amer-*
561 *ica*. 2017 12;114(51):13573–13578. Available from: [http://www.ncbi.nlm.nih.gov/pubmed/](http://www.ncbi.nlm.nih.gov/pubmed/29208707)
562 [29208707](http://www.ncbi.nlm.nih.gov/pubmed/29208707)<http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5754802>.

- 563 [57] Levine JM, Bascompte J, Adler PB, Allesina S. Beyond pairwise mechanisms of species
564 coexistence in complex communities. *Nature*. 2017 5;546(7656):56–64. Available from: <http://www.nature.com/doi/10.1038/nature22898>.
565
- 566 [58] Draper E, Bissett S, Howell-Jones R, Edwards D. Neutralization of non-vaccine human
567 papillomavirus pseudoviruses from the A7 and A9 species groups by bivalent HPV vaccine
568 sera. *Vaccine*. 2011;29(47):8585–8590. Available from: <http://www.sciencedirect.com/science/article/pii/S0264410X11014277>.
569
- 570 [59] Kemp TJ, Hildesheim A, Safaeian M, Dauner JG, Pan Y, Porras C, et al. HPV16/18 L1
571 VLP vaccine induces cross-neutralizing antibodies that may mediate cross-protection. *Vac-*
572 *cine*. 2011;29(11):2011–2014. Available from: <http://www.sciencedirect.com/science/article/pii/S0264410X1100017X>.
573
- 574 [60] Mihaljevic JR, Joseph MB, Johnson PT. Using multispecies occupancy models to improve the
575 characterization and understanding of metacommunity structure. *Ecology*. 2015;96(7):1783–
576 1792.
- 577 [61] Warton DI, Blanchet FG, O’Hara RB, Ovaskainen O, Taskinen S, Walker SC, et al.. So Many
578 Variables: Joint Modeling in Community Ecology. Elsevier Ltd; 2015. Available from: <http://dx.doi.org/10.1016/j.j.tree.2015.09.007>.
579
- 580 [62] Chesson P. Mechanisms of Maintenance of Species Diversity. *Annual Review of Ecology*
581 *and Systematics*. 2000 11;31(1):343–366. Available from: <http://www.annualreviews.org/doi/10.1146/annurev.ecolsys.31.1.343>.
582
- 583 [63] Gelman A. Inference and monitoring convergence. In: Gilks W, Richardson S, Spiegelhalter
584 D, editors. *Markov Chain Monte Carlo in Practice*. CRC Press; 1996. p. 131–143.

Figure Legends

585

586 **Figure 1:** Simpson's paradox demonstrated for two species that are sampled across ten habitat
587 sites, with each site surveyed fifteen times. **A** Species covary positively across sites (over space),
588 indicating response to similar habitat requirements. **B** Species covary negatively within sites over
589 time, indicating inter-specific competition. Probabilities of occurrence are on the probit scale.

590 **Figure 2:** Inference of model parameters from synthetic data simulated for 10 pathogen strains
591 across 1500 patients, where each patient was tested 10 times. **A** Recovery of baseline occurrence
592 probability for the 10 strains. Red vertical line gives the true value. **B** Recovery of positive,
593 pairwise correlations in among-patient random effects. Dashed line represents zero effect, while
594 dotted line represents the true value (0.5). **C** Recovery of weakly negative, pairwise correlations in
595 within-patient, observation-level random effects. Dashed line represents zero effect, while dotted
596 line represents the true value (-0.11). **D** Recovery of fixed inter-strain effects on probability of
597 strain persistence. **E** Recovery of fixed inter-strain effects on probability of strain colonization.

598 **Figure 3:** Inference of model parameters from the HIM data. **A** Estimate of the baseline occur-
599 rence probability for each HPV type. **B** Inferred correlations in among-patient random effects.
600 **C** Inferred correlations in within-patient, observation-level random effects. **D** Recovery of fixed
601 inter-type effects on the probability of type persistence. **E** Recovery of fixed inter-type effects on
602 probability of type colonization.

Tables

Table 1: Comparison of candidate models using leave-one-out cross-validation. The table shows whether fixed and/or random effects were included, the log-likelihood of the model fit (i.e. $\mathcal{L}(\theta|D)$), and the LOO-IC. The standard error (SE) in the LOO-IC is shown to emphasize that the LOO-IC is an estimated statistic with error, but also that none of our LOO-IC values overlap within $\pm 2SE$.

HPV Interactions	Among-patient and Among-visit Correlations	Log- Likelihood	LOO-IC	SE LOO-IC
	✓	-220310.1	708323.6	373.8
✓	✓	-279574.9	825439.5	329.8
✓		-433036.9	1109717.0	465.3
		-432977.8	1130039.0	681.3

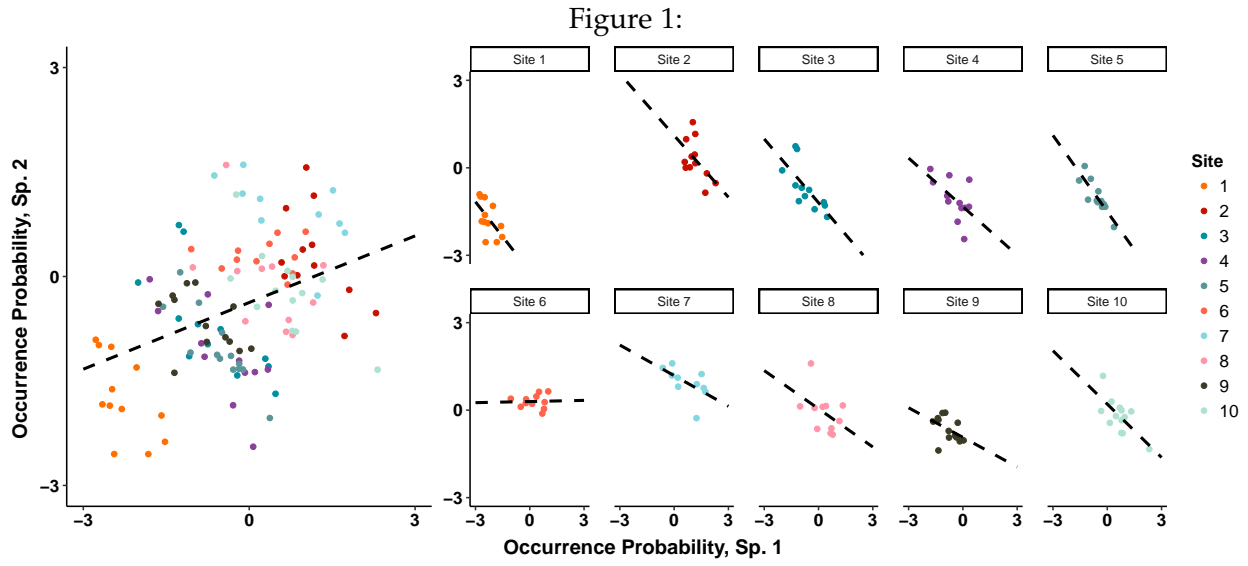


Figure 2:

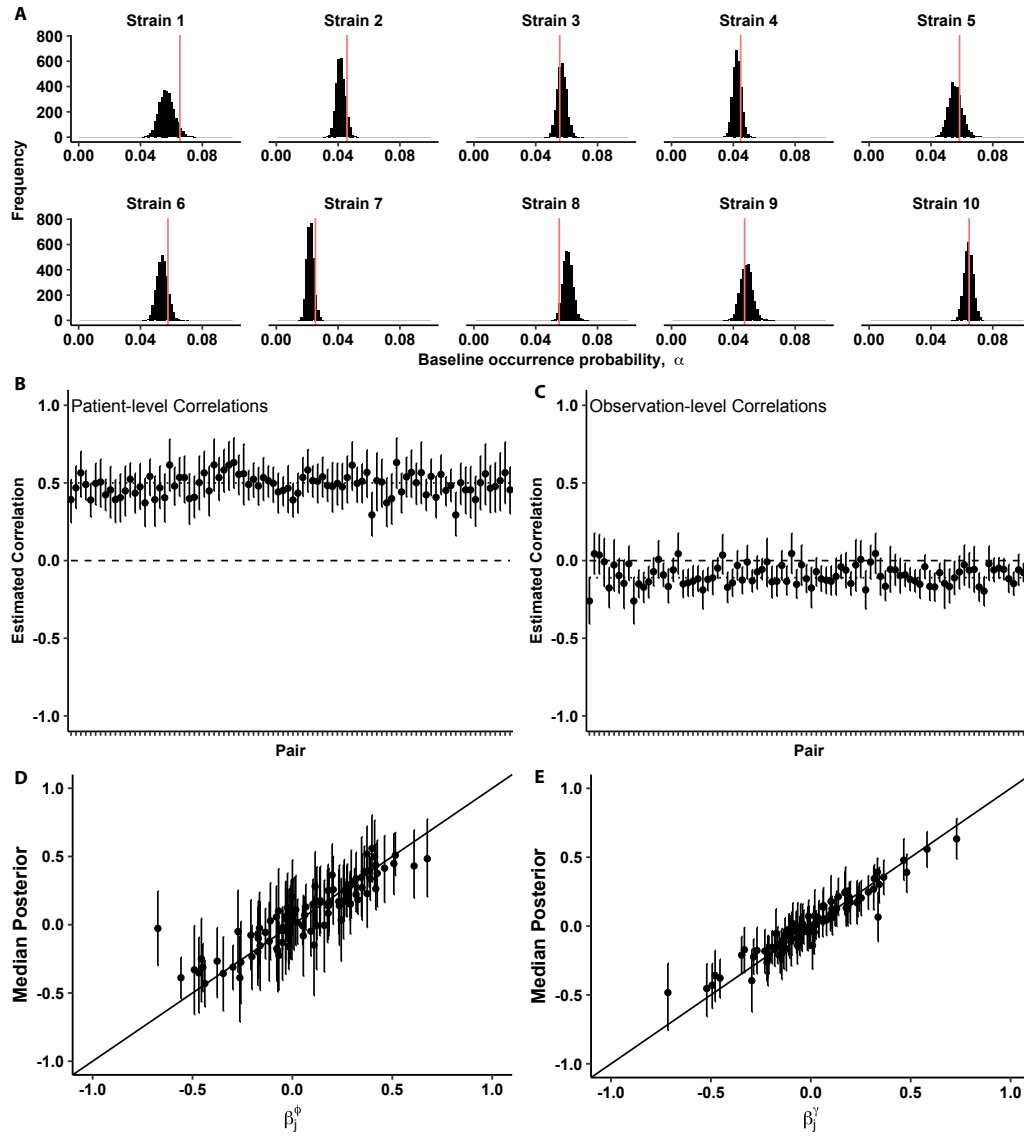
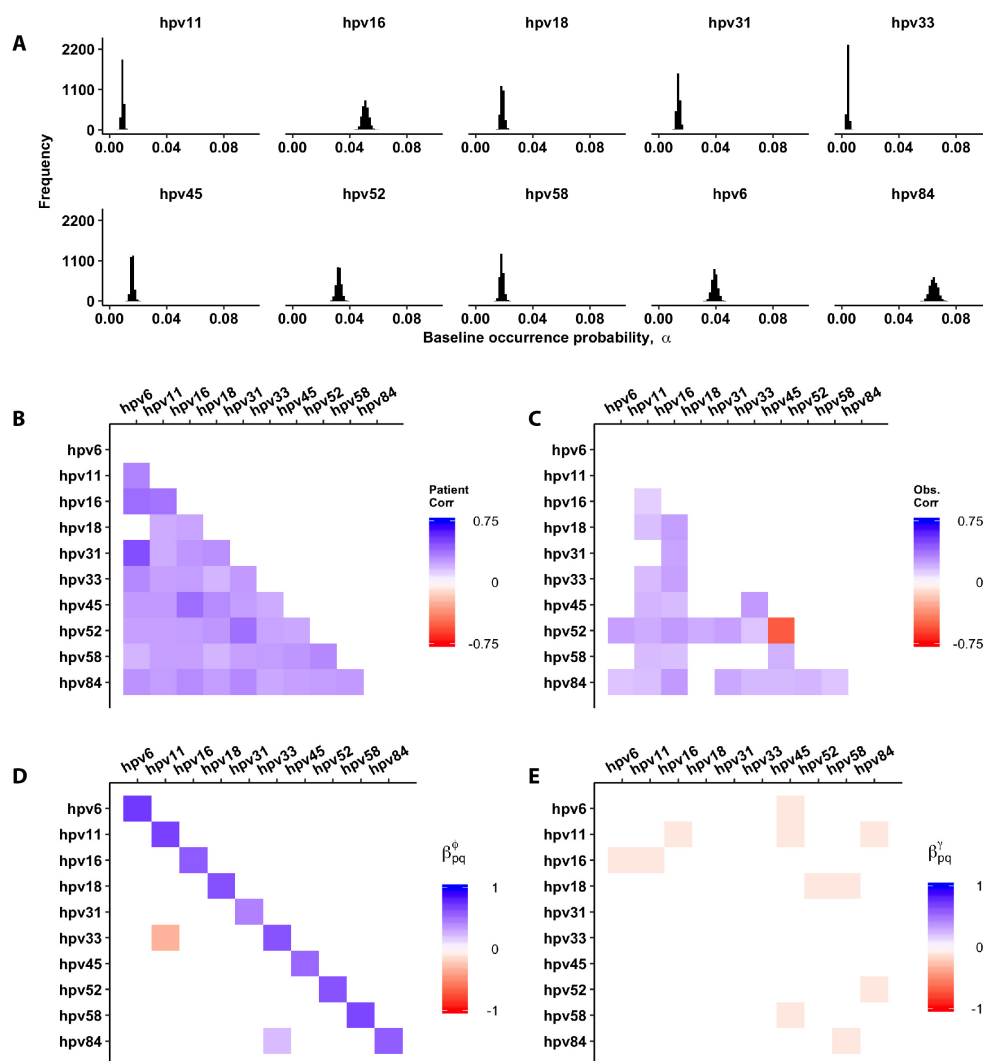


Figure 3:



Supporting Information (SI)

Subset of HIM data included in the analysis

We excluded individuals that failed to meet the full eligibility criteria described by the HIM study [32]. The criteria included: ages 18 to 70 years; residents of one of three sites — Sao Paulo, Brazil; Morelos, Mexico; or southern Florida, United States; no prior diagnosis of penile or anal cancers; no prior history of genital or anal warts; no symptoms of a sexually transmitted infection at baseline or recent treatment for a sexually transmitted infection; no history of participation in an HPV vaccine study; and no history of HIV or AIDS.

We identified 3,656 eligible participants from the 4,123 men enrolled in the HIM study as of October 2014. For each of the 10 HPV types that we analyzed, we include in our data the binary infection status of each man at each clinic visit. We also include the length of time between consecutive clinic visits.

Type-specific HPV prevalence over follow-up

We calculated the prevalence of the 10 HPV types included in the analysis at each visit (Fig. S1). Note that, because individuals varied in their visit dates, the prevalence at each visit is a time-averaged estimate. The data show that the expected distribution of HPV types in the metacommunity is consistent across visits.

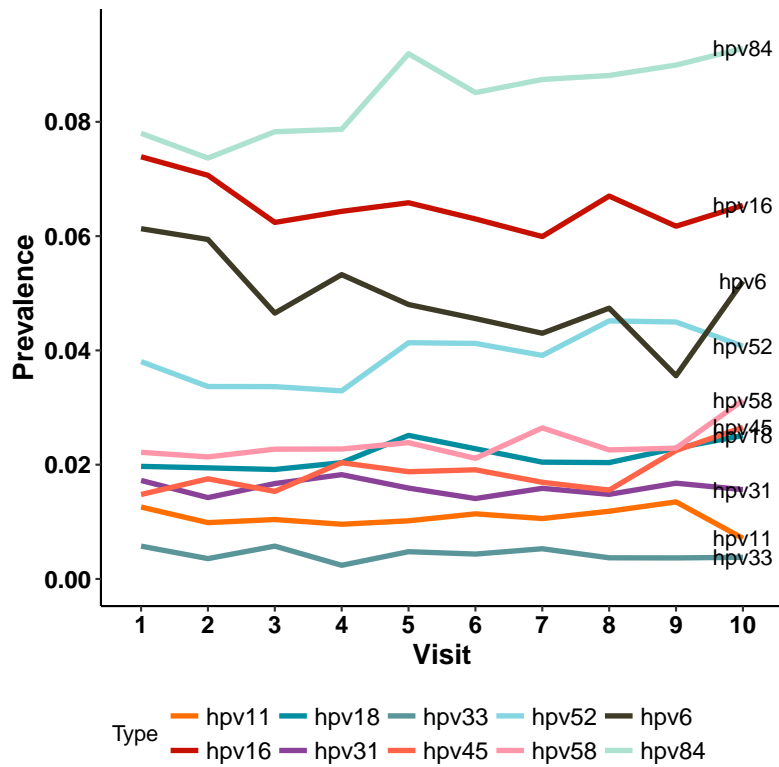


Figure S1: Observed visit-level prevalence of each of the 10 HPV types included in this analysis.

18

Stan model details

19 All of our code to run the Stan model is provided in our open-source repository ****LINK****, but we
20 will briefly describe the fitting routine here. For each nested model, we ran three MCMC chains in
21 parallel on the Gardner high performance computing (HPC) cluster at the University of Chicago
22 (Center for Research Informatics). Each chain ran for 5000 iterations with a 2000 iteration warm-
23 up period, and we thinned the samples by three, giving us a total of 1000 posterior samples from
24 each chain. Parameter samples were stored as tables in a SQLite database for later processing.
25 Due to the large number of columns of the log-likelihood table, we split this table into sub-
26 components before storage. We monitored convergence with the Gelman-Rubin (\hat{R}) statistic, and
27 we conducted several standard visual diagnostics to check MCMC chain performance [63, 53].
28 All models converged after 5000 iterations, and no problems were observed in the MCMC chains.

29

Time between visits

30 Here we display the effects of time between visit (TBV) on persistence and colonization probabilit-
31 ities for the synthetic data (Fig. S2) and for the HIM dataset, using the full model that includes
32 both correlations and fixed, pairwise interactions (Fig. S3).

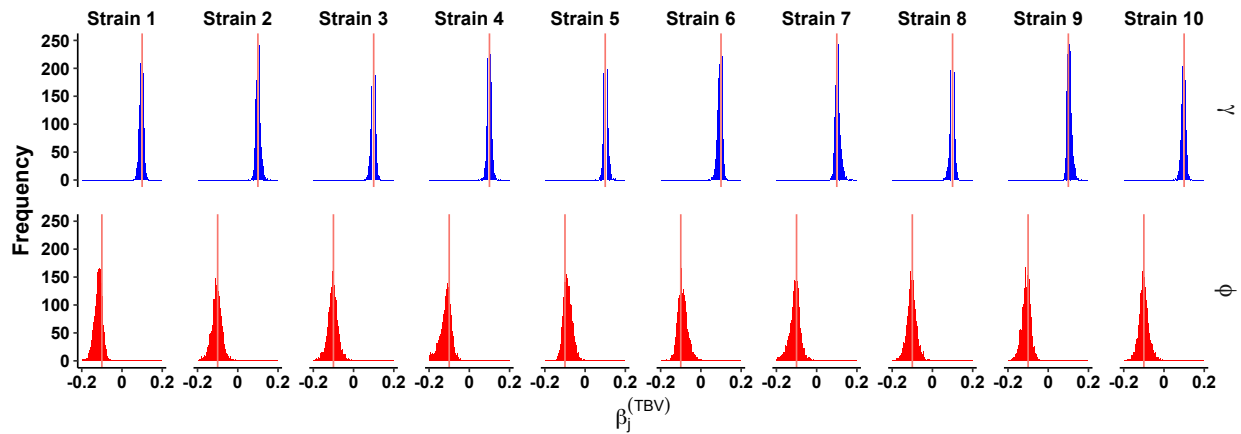


Figure S2: Effects of time between visit (TBV) on colonization (top row) and persistence (bottom row) probabilities for each of 10 simulated pathogen strains, from the synthetic data. These results are generated from the full model, which has both fixed effects of pairwise interactions, as well as patient-level and observation-level correlations among residuals. Blue histograms are effects greater than zero, while red histograms are effects less than zero, based on a 95% credible intervals (CI) that does not overlap zero. The true, simulated values are shown as red vertical lines

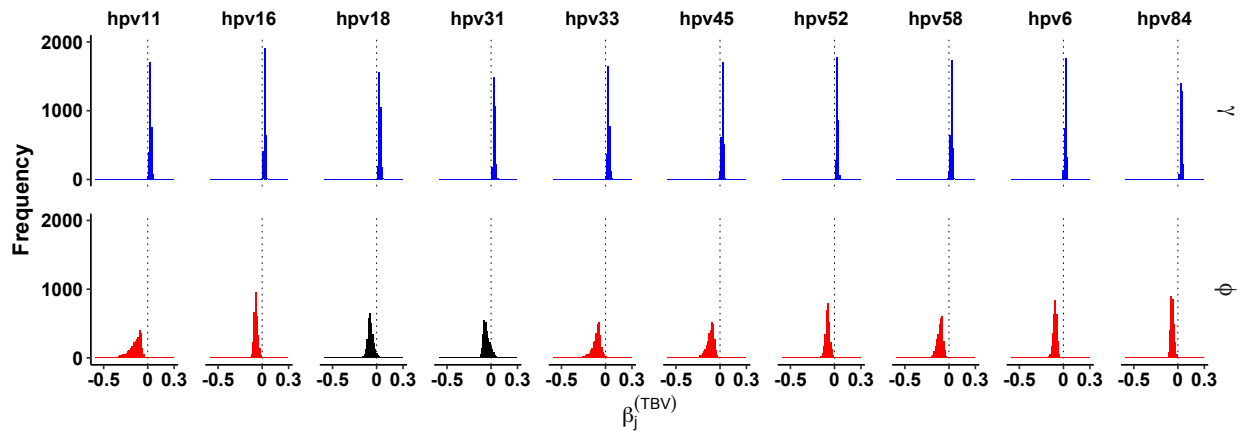


Figure S3: Effects of time between visit (TBV) on colonization (top row) and persistence (bottom row) probabilities for each HPV type. These results are generated from the full model, which has both fixed effects of pairwise interactions, as well as patient-level and observation-level correlations among residuals. Blue histograms are effects greater than zero, while red histograms are effects less than zero (marked as the vertical dotted line), based on a 95% credible intervals (CI) that does not overlap zero. Histograms with black bars have effects with 95% credible intervals (CI) that overlap zero.

33

Results from “best” model, with no pairwise interaction effects

34

35

36

37

38

The figure below displays the results from the most preferred model, which includes the random effects (i.e. patient-level and observation-level correlations among HPV types), but does not include pairwise effects on persistence and colonization probabilities (Fig. S4). Notably, this model is nearly identical to the full model in terms of baseline probabilities of occurrence (Fig. S4 A), the random effects (Fig. S4 B,C), and the effects of time between visit (TBV) (Fig. S4 D).

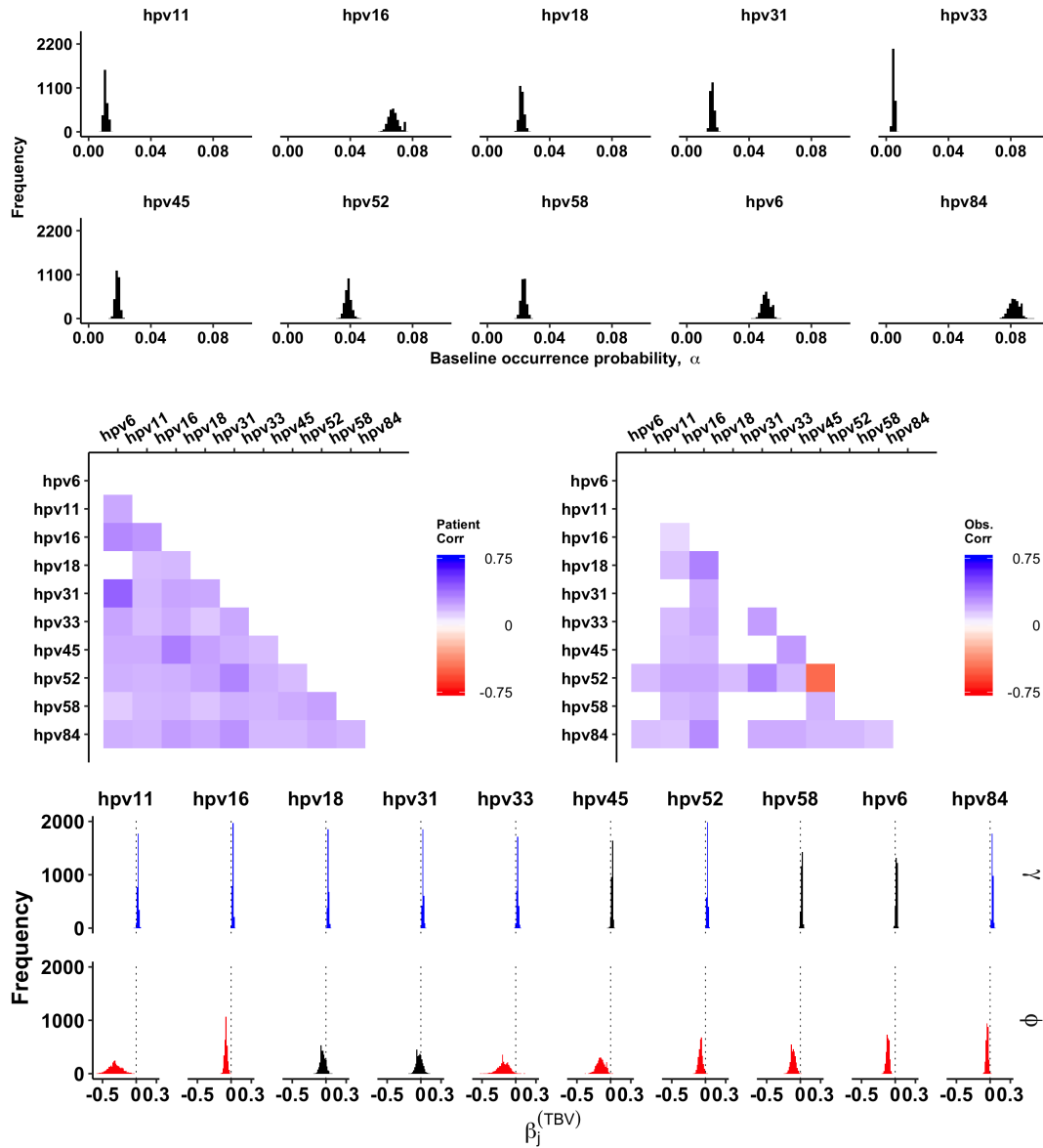


Figure S4: Inference of model parameters from the HIM data, using the “best” model, which only has patient-level and observation-level correlations among types, but does not have fixed effects on persistence and colonization. **A** Estimate of the baseline occurrence probability for each HPV type. **B** Inferred correlations in among-patient random effects. **C** Inferred correlations in within-patient, observation-level random effects. **D** Estimates of the effects of time between visit (TBV) on persistence and colonization probabilities. Colors and vertical lines in **D** are the same as in Fig. S3.