**Persistent phylogeographic structure of an emerging virus on a homogeneous landscape**

Trieste Musial[1,2*], Scott Duke-Sylvester[3], Rolan Davis[4], Roman Biek[5*], Leslie A. Real[2,6]

[1] Population Biology, Ecology, and Evolution Program, Emory University, Atlanta, Georgia, United States of America

[2] Center for Disease Ecology, Emory University, Atlanta, Georgia, United States of America

[3] Department of Biology, University of Louisiana at Lafayette, Lafayette, Louisiana, United States of America

[4] Rabies Laboratory, Veterinary Diagnostic Laboratory, Kansas State University, Manhattan, Kansas, United States of America

[5] Boyd Orr Centre for Population and Ecosystem Health, Institute of Biodiversity, Animal Health and Comparative Medicine, College of Medical, Veterinary, and Life Sciences, University of Glasgow, Glasgow, United Kingdom

[6] Department of Biology, Emory University, Atlanta, Georgia, United States of America


* Corresponding authors

E-mail: tmusial@emory.edu; roman.biek@glasgow.ac.uk

1

**Abstract**

Landscape composition and structure influence animal movement, which in turn can affect transmission of their diseases. Spatio-temporal variation in host diffusion, caused by landscape heterogeneity, is thus expected to generate corresponding phylogeographic patterns in the pathogen. However, establishing causative links between genetic structure in pathogen populations and environmental variation does require appropriate null models. Here, we present an empirical example of the emergence and multi-decade persistence of phylogeographic structure on a homogeneous landscape in a rapidly diversifying pathogen in the absence of any apparent landscape heterogeneity. By applying phylogeographic inference to 173 sequences of a raccoon-specific strain of rabies virus, we reconstruct patterns of the virus' evolution and diffusion on the Florida peninsula, USA, from its first emergence in the 1940's to the present. Consistent with a lack of significant landscape heterogeneity relevant to raccoon movement in Florida, we found that the speed of rabies virus diffusion was spatially homogeneous across the peninsula. In contrast, we document the emergence of strong phylogeographic structure in the virus, in the form of five monophyletic lineages that diverged during the early years of colonization and now each occupy a distinct sub-region of Florida. Based on samples taken over multiple decades, we show that the spatial distribution of these lineages has changed little over the past four decades. This phylogeographic stability allowed us to retrospectively identify a small set of counties within Florida as the likely source of the virus strain that seeded a much larger rabies outbreak in the northeastern USA in the 1970s. Our results provide a rare empirical demonstration that spatial genetic structure can arise and be maintained in the absence of landscape heterogeneity, which has wider implications for the interpretation of phylogeographic data and the reconstruction of historical colonization patterns from molecular data.

2

**Introduction**

24
25
26          The geographic genetic organization of populations is a complex interaction between

27   patterns of organismal movement, breeding structure, habitat heterogeneity, and dispersal and

28   establishment.   On landscapes that are homogeneous with respect to individual movement

29   probabilities, genetic organization traditionally is considered to be determined by patterns of

30   local gene exchange and patterns of colonization by exogenous genotypes (Irwin 2002; Kuo &

31   Avise 2005; Waters 2011; Waters *et al.* 2013; Wright 1943, 1951).  Spatial genetic patterns on

32   these landscapes follow expectations of isolation-by-distance coupled to unpredictable stochastic

33   migration events and stochastic ecological extinctions (Irwin 2002; Marske *et al.* 2013; Neigel &

34   Avise 1993; Saunders *et al.* 1986).  In contrast, heterogeneous landscapes, with barriers to

35   movement or gene exchange and/or characterized by areas of inhospitable terrain, tend to lead to

36   more predictable spatial genetic relationships, delimited – and stabilized – by recognizable

37   landscape features (Avise *et al.* 1987; Kuo & Avise 2005; Manel *et al.* 2003).  Significant spatial

38   organization of genotypes independent of local gene exchange has long been considered *prima*

39   *facie* evidence of landscape heterogeneity associated with restrictions on movement and

40   establishment of genotypes distributed over space (Avise *et al.* 1987; Marske *et al.* 2013; Wright

41   1951).

42          Some research has indicated however that random local processes, as well as events

43   during colonization and expansion, can also generate enduring spatial genetic patterns.  When

44   this is the case, population spatial genetic structure may not be dictated solely by barrier

45   arrangement.  On simulated homogeneous landscapes, significant "phylogeographic breaks", i.e.

46   spatial partitioning of distinct phylogenetic lineages (Avise *et al.* 1987), can emerge

47   independently of underlying geographic barriers within established, continuous populations

3

Preprints

48    where dispersal rates are low and/or populations sizes are small (Baptestini *et al.* 2013; Irwin

49    2002; Kuo & Avise 2005).  Stochastic events will also shape the spatial arrangement of lineages

50    that diverge during continuous range expansion, and these processes can generate lasting

51    signatures that dominate long-term spatial genetic patterns (Excoffier & Ray 2008; Hallatschek

52    *et al.* 2007; Waters *et al.* 2013).  Phylogeographic breaks formed by stochastic processes contain

53    no information about the spatial ecology of the lineage in question but can easily be mistaken for

54    patterns generated by landscape-level processes (Crisp *et al.* 2011; Real *et al.* 2005a; Schwartz &

55    McKelvey 2008).  As of yet, the characterization of persistent genetic breaks on homogeneous

56    landscapes has been limited to simulation studies (e.g., Baptestini *et al.* 2013; Ibrahim *et al.*

57    1996; Irwin 2002; Kuo & Avise 2005) and microbial *in vitro* experiments (Hallatschek *et al.*

58    2007; Hallatschek & Nelson 2008).  Two key questions are thus: 1) under what conditions will

59    these relationships form on natural landscapes?  And 2) once formed, how stable are these spatial

60    patterns across natural landscapes that are highly homogeneous?  Here, we address these

61    questions using a rapidly evolving zoonotic pathogen, rabies virus, as a model.

62         Terrestrial rabies, caused by an RNA virus, is often used as a model system for inferring

63    the effects of landscape heterogeneity on phylogeographic patterns (e.g. Biek *et al.* 2007;

64    Brunker *et al.* 2012; Szanto *et al.* 2011).  The high mutation rate of rabies virus generates

65    substitutions on a scale similar to the rate of ecological changes, making this and other RNA

66    virus systems ideal for epidemiological reconstruction from genomic data (Biek *et al.* 2015;

67    Drummond *et al.* 2003; Pybus & Rambaut 2009).  Movement of rabies virus (like other directly-

68    transmitted viruses) generally occurs across landscapes and distances that are amenable to

69    independent host movement, though human-mediated long-distance translocation events are

70    known to have facilitated various rabies epizootics and range expansions (Nettles *et al.* 1979;

4

71    Talbi *et al.* 2010; Wilson *et al.* 1997).  Rabies lineages display strong regional

72    compartmentalization (e.g., Biek *et al.* 2007; Kuzmina *et al.* 2013; Real *et al.* 2005a; Szanto *et*

73    *al.* 2011), even in cases where a single strain is hosted by multiple species (Lembo *et al.* 2007;

74    Nadin-Davis *et al.* 1994).

75         Although rabies virus can infect any mammal, different variants are usually maintained

76    by a particular host species. One such host-specific rabies variant is raccoon rabies virus (RRV),

77    which is distributed across much of eastern North America. Here we investigate the conditions

78    under which phylogeographic breaks in RRV develop and persist on a homogeneous landscape,

79    the Florida peninsula.  RRV was first detected in Florida in 1947 (Bigler *et al.* 1973), and had

80    covered the eastern U.S. within 50 years (Childs *et al.* 2000).  The process of RRV expansion

81    was distinctly different between its initial establishment in Florida and its later spread throughout

82    the mid-Atlantic:  In the southeastern US, RRV was detected only sporadically in the decades

83    following its emergence in Florida and generated few detected, localized epizootic events (Bigler

84    *et al.* 1973; Kappus *et al.* 1970; Scatterday *et al.* 1960).  In fact, increased surveillance efforts

85    during the decade following its detection in 1947 determined that RRV had already reached an

86    enzootic state throughout peninsular Florida, having potentially gone undetected for much of its

87    early expansion (Kappus *et al.* 1970; Scatterday *et al.* 1960).  In contrast, the spread of RRV

88    through the mid-Atlantic states, likely originating from infected raccoons translocated from

89    southeastern states in the late 1970s (Nettles *et al.* 1979; Rupprecht & Smith 1994), proceeded

90    rapidly and conspicuously, generating one of the largest wildlife epizootics in history (Childs *et*

91    *al.* 2000).  Mountains, rivers, and major water bodies frequently have been recognized as barriers

92    to the movement of raccoons and, consequently, RRV (Biek *et al.* 2007; Smith *et al.* 2002;

93    Wheeler & Waller 2008).  Accordingly, phylogeographic analyses of RRV within its mid-

5

94    Atlantic range identified decreased viral velocity across mountain ranges (Biek *et al.* 2007). The

95    landscape of Florida, in contrast, exhibits few features that might serve as barriers to raccoon

96    movement. Florida's topography is relatively flat (with a max elevation of 105m above sea

97    level), and there is an abundance of favorable and continuous raccoon habitat (including large

98    swamps and long stretches of urbanized areas). Previous analyses based on microsatellite and

99    mtDNA data found evidence for a single, well-mixed raccoon population that covers Florida,

100    Georgia, Alabama, Tennessee, and South Carolina (Cullingham *et al.* 2008; Reeder-Carroll

101    2010), consistent with a lack of impediments to raccoon gene flow throughout this region.

102         The raccoon rabies system in Florida offers a chance to explore long-term evolutionary

103    outcomes of invasion on a landscape devoid of spatial features that would be predicted to

104    maintain phylogeographic patterns. On this landscape, one might expect to see little or no

105    phylogeographic structure of viral strains and would instead predict lineages to spatially admix

106    over time, eliminating early spatial genetic differentiation. We are testing this hypothesis by

107    analyzing spatial evolutionary patterns of RRV expansion and lineage divergence in Florida.

108    Utilizing novel methods to reconstruct spatial patterns of viral diffusion rates within our study

109    area, we aimed to determine how RRV phylogeographic patterns arose and are maintained

110    following emergence of the virus and how this affects the interpretation of genetic data in the

111    context of RRV as well as biological invasion processes more generally.

112    **Materials and Methods**
113    **Sample collection and preparation**

114         We analyzed 173 brain tissue samples collected as part of ongoing rabies surveillance

115    efforts by state and local public health departments in Florida (*n*=164) and Alabama (*n*=9) from

116    1982 to 2012. Sampling fell within three general time periods – 37 samples were collected from

117    1982 to 1988, 40 from 1998 to 2004, and 96 from 2009 to 2012. Sample locations from 2003

       6

118     forward (*n*=103) were georeferenced to zip code; location data prior to 2003 were recorded as

119     county centroid (Fig 1).  Infection with the raccoon-specific rabies variant was initially

120     confirmed using an indirect assay with monoclonal antibodies for the nucleocapsid protein

121     (Smith *et al.* 1986). Total RNA was extracted from 50 – 100 mg of frozen brain tissue archived

122     after post mortem analysis of rabid raccoons.  RNA extraction was achieved using a hypotonic

123     lysis buffer (Smith *et al.* 1991) followed by precipitation with TRIzol Reagent (Invitrogen,

124     Carlsbad, CA, USA) and following the manufacturer's protocol.  Pellets were stored at -20°C in

125     DEPC-treated water until amplification.

126          RT-PCR was used to amplify the 5' portion of the rabies virus nucleoprotein (N) gene in

127     all samples using forward primer 21G (5'-ATGTAACACCTCTACA-3') (Orciari *et al.* 2001)

128     and reverse primer 304 (5'-ATGAGCAAGATCTTCGT-3').  Additionally, a portion of the

129     glycoprotein (G) gene was amplified in a subset of 20 samples, using primers and conditions

130     described in Biek *et al*. (2007).  Our amplification procedure utilized the SuperScript III One-

131     Step RT-PCR System with Platinum *Taq* High Fidelity (Invitrogen).  Amplicons were examined

132     by gel electrophoresis and samples with poor yield were reamplified prior to purification with the

133     Promega Wizard PCR kit (Promega, Madison, WI) and sequencing on an ABI 377 DNA

134     Sequencer (Applied Biosystems, Foster City, CA).  Final sequences had no indels and were

135     aligned manually and trimmed to final lengths of 591 nt at the N gene (*n*=173) and 1374 nt at the

136     G gene (n=20).  The N and G genes exhibit no significant incongruence (Biek *et al.* 2007) and

137     were concatenated for our analyses.

138     **Phylogenetic analysis**

139          Potential evolutionary models were explored using the phyml package in jModelTest 2

140     (Darriba *et al.* 2012; Guindon & Gascuel 2003), and we used BEAST's (Drummond *et al.* 2012)

7

141    path sampling procedure (Baele *et al.* 2012; Baele *et al.* 2013) to identify the best fit model;

142    computations were performed using XSEDE cloud-computing resources via CIPRES (Miller *et*

143    *al.* 2010; Towns *et al.* 2014).  A strict molecular clock with gamma-distributed rate variation

144    applied independently between codons at position 1+2 vs 3 ($HKY_{112} + CP_{112} + G_{112}$) (Hasegawa

145    *et al.* 1985; Yang 1994) was clearly the best model for our data.  We used BEAST to

146    simultaneously estimate phylogenetic relationships and ancestral dispersal patterns from our

147    samples.  Initial analyses using marginal likelihood estimation (Baele *et al.* 2012; Baele *et al.*

148    2013), identified a strict clock combined with a codon-specific SDR06 model (Shapiro *et al.*

149    2005) as appropriate for our data.  We modeled coalescence patterns with a Bayesian skyline

150    demographic prior, applying a piecewise linear smoother (Drummond *et al.* 2005) after

151    establishing its improved fit over the piecewise constant option.  The molecular clock was

152    parameterized with a normally distributed prior for the rate parameter with mean and variance set

153    to the evolutionary rate estimates reported in Biek et al. (2007).  Simultaneously, spatial

154    locations at internal nodes were reconstructed as continuous traits at each tree in the posterior:

155    Observed sample locations served as discrete priors for a Relaxed Random Walk model (RRW),

156    which allows the variance in diffusion rate from node to descendent to vary from branch to

157    branch within a topology (Lemey *et al.* 2009; Lemey *et al.* 2010).  We tested a number of models

158    for this diffusion process, and a gamma-distributed rate variation applied to an RRW was

159    selected based on MLE.  Because our samples were georeferenced to either a county or a zip

160    code, our sample locations included 55 spatial duplicates, to which we added random spatial

161    noise using a random "jitter" window of 0.5 before analysis with BEAST.

162            Finally, we ran our model for three independent runs, each of 300 million steps and

163    sampling from the posterior distribution of trees every 30,000 steps, which yielded 10,000

8

164    posterior trees per run and high ESS values (1000+) for each estimated parameter within the

165    runs.  Combining posterior outputs from the three runs caused ESS values to drop abruptly.). We

166    determined this was caused by a monophyletic clade within our tree converging on different

167    sister taxa between runs.  Since there were no practical differences between phylogenies holding

168    this clade in different positions (the timing of its MRCA did not change, nor did its placement

169    disrupt the branch lengths or topology of any other monophyly in our tree), we processed two

170    more runs and ultimately combined the three that agreed.  We chose burn-in lengths based on

171    posterior mixing patterns, discarding anywhere from 2,000 to 3,000 of the 10,000 trees per run

172    before sampling from the remainder for combining in LogCombiner (Drummond & Rambaut

173    2007).  The Maximum Clade Credibility (MCC) tree drawn from the combined posteriors was

174    summarized in TreeAnnotator.  For ease of comparisons we analyzed the N/G concatenated

175    sequences using the same parameters and models as for the N data.

**Reconstructing spatial patterns of diffusion**

177         Estimated ancestral locations and ages were used with sample locations and times to

178    produce a two-dimensional surface model of viral diffusion rates.  Calculated simply as the

179    geographic distance between a parent and daughter node and scaled by the time separating them,

180    these diffusion rates are conditioned on phylogenetic history. To accommodate uncertainty in the

181    estimated phylogeographic locations associated with nodes, as well as the phylogenetic

182    dependence of inferred locations from neighboring nodes, 1000 trees of the N-gene drawn

183    randomly from the posterior distributions logged by BEAST were used to build a data set of

184    ancestral node heights and spatial coordinates.  For each tree, we determined temporal and

185    geographic distances between parent and daughter nodes to obtain a branch-specific diffusion

186    rate for each edge.  We assigned these diffusion rate estimates to point locations by generating

9

187    50 uniformly distributed (X, Y) coordinate pairs along the great-circle line connecting parent and

188    daughter node.  One intermediate point was randomly selected from this distribution and

189    assigned a node height as a function of the height of the parent node, the distance of the

190    intermediate from the parent, and the branch-specific dispersal rate. Applying this process to the

191    sample of 1000 trees, yielded 344,000 georeferenced, time-stamped point estimates of viral

192    diffusion rates.

193        Rate estimates and locations were used to model viral diffusion across Florida.  A

194    Generalized Additive Model (GAM – in R package mgcv (Wood 2011)) with penalized thin-

195    plate regression smoothing was used to fit a model characterizing the relationship between

196    diffusion rates and spatial coordinates.  Predicted diffusion rates were then interpolated on a

197    continuous grid of coordinates drawn from the rest of Florida.  The optimal model function, a

198    Gamma distribution with a logarithmic link, was selected by AIC.

199    **Simulations**

200        As our genetic samples were drawn discontinuously in space and time, we examined the

201    effect of sampling regime on our inference of the diffusion process.  Simulated sequences were

202    evolved on a grid that allowed equal probabilities of individuals diffusing to neighboring points

203    consistent with a homogeneous landscape (see Duke-Sylvester, et al. (2013) for details), and

204    sampled in spatiotemporal clusters representing different forms of sampling bias as well as the

205    empirical sampling scheme used in our study (see SI).

206    **Phylogeographic analysis**

207        Our analysis revealed the existence of five major RRV lineages in Florida that tended to

208    cluster spatially (see Results). We used Delaunay triangulation of sample locations to visualize

209    and track the distributions of these lineages over time.  Delaunay triangulation is a method of

10

210    calculating the two-dimensional geometric approximation of each point in a set, drawing

211    polygons around each point such that every point in the set is closer to its own polygon than it is

212    to any other point.  We generated Delaunay diagrams for sample subsets corresponding to our

213    three sampling time periods.  For each period, we aggregated Delaunay cells by lineage

214    membership and calculated the mean center of the resulting polygon, as well as the distance it

215    had moved since the prior time slice.

216    **Reconstructing the source of mid-Atlantic RRV**
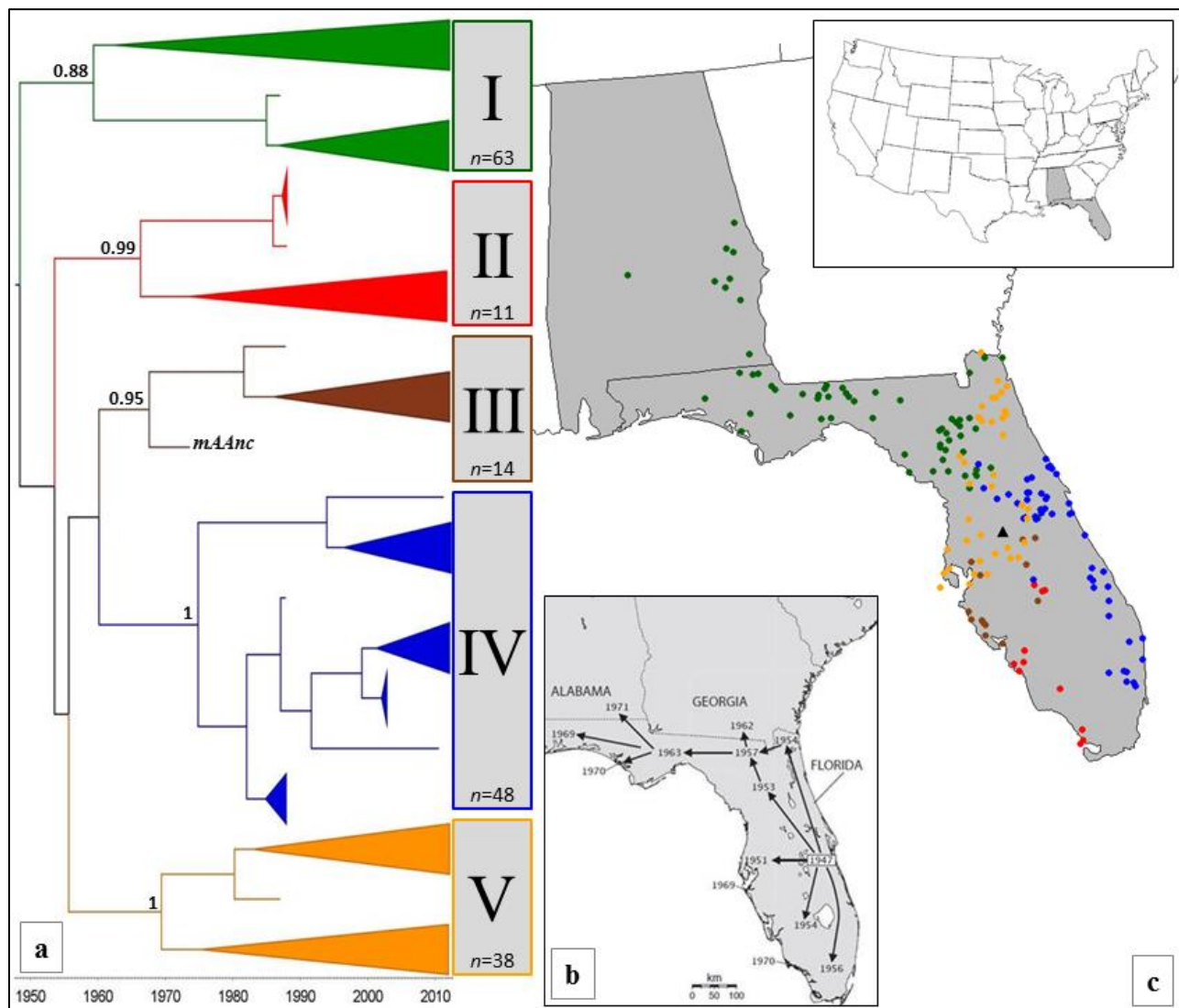
217         The large RRV epizootic that emerged in the Virginias in the 1970's and rapidly spread

218    throughout the mid-Atlantic states is thought to have been caused by the translocation of

219    raccoons from Florida (Childs *et al.* 2000; Nettles *et al.* 1979; Rupprecht & Smith 1994).

220    However, the origin and timing of this translocation event has never been identified.  To infer the

221    putative sequence of the virus that sparked the mid-Atlantic epizootic we applied marginal

222    reconstruction in baseml (Yang 2007) to the maximum clade credibility tree of concatenated N

223    and G sequences collected throughout the mid-Atlantic in an earlier study (Biek *et al.* 2007) to

224    obtain the ancestral sequence at the root node.  This sequence was incorporated into our

225    phylogenetic analyses of the N-gene and the concatenated N/G sequence data sets from Florida.

226    **Results**
227    **Phylogenetic analysis**

228         Our 173 samples exhibited approximately 98% overall sequence identity and included

229    137 unique variants of the N gene.  Phylogenetic analyses inferred five major clades, each

230    supported by posterior probabilities of ~0.9 or more (Fig 1a).  Relationships among these five

231    clades were more difficult to resolve, with posterior probabilities of more ancestral nodes

232    ranging from 0.3 to 0.7.  The estimated rate of evolution was $3.1 \times 10^{-4}$ substitutions/site/yr,

233    yielding a date for the most recent common ancestor of 1944 (95% HPD=1933 to 1961).  All five

11

234    clades diverged rapidly in the years immediately following emergence of RRV; the 95%

235    confidence intervals for the node height of all inferred lineage ancestors overlap (Table 1). We

236    estimated that four clades were established by the late 1970s; the 95% HPD interval for the

237    divergence time of clade IV extends to the early 1980s (Table 1, Fig 1a).  All five clades were

238    extant throughout the 30 years of our sample collection.  Analysis of concatenated N and G gene

239    sequences on a subset of our sample reaffirmed the patterns found with the N gene alone (Table

240    1 and SI 2).

241
242



        12

**Fig 1. Evolutionary relationships of raccoon rabies virus samples collected from southeastern US, with historic patterns of spread.** (**A**) Maximum clade credibility (MCC) tree from Bayesian coalescent analysis of N gene sequences, scaled to time. Node labels report posterior probabilities at each clade's ancestor. Groups within each clade are collapsed at posterior values <0.95. The sequence of the reconstructed ancestor of the mid-Atlantic RRV epizootic is labeled *mAAnc*. (**B**) RRV reported cases from the date and location of emergence, redrawn after Bigler et al. (1973). The first reported case is shown by a boxed date. (**C**) RRV positive samples sequenced at the N gene with their corresponding lineage assignments by color. The estimated tree root location is shown as a black triangle.
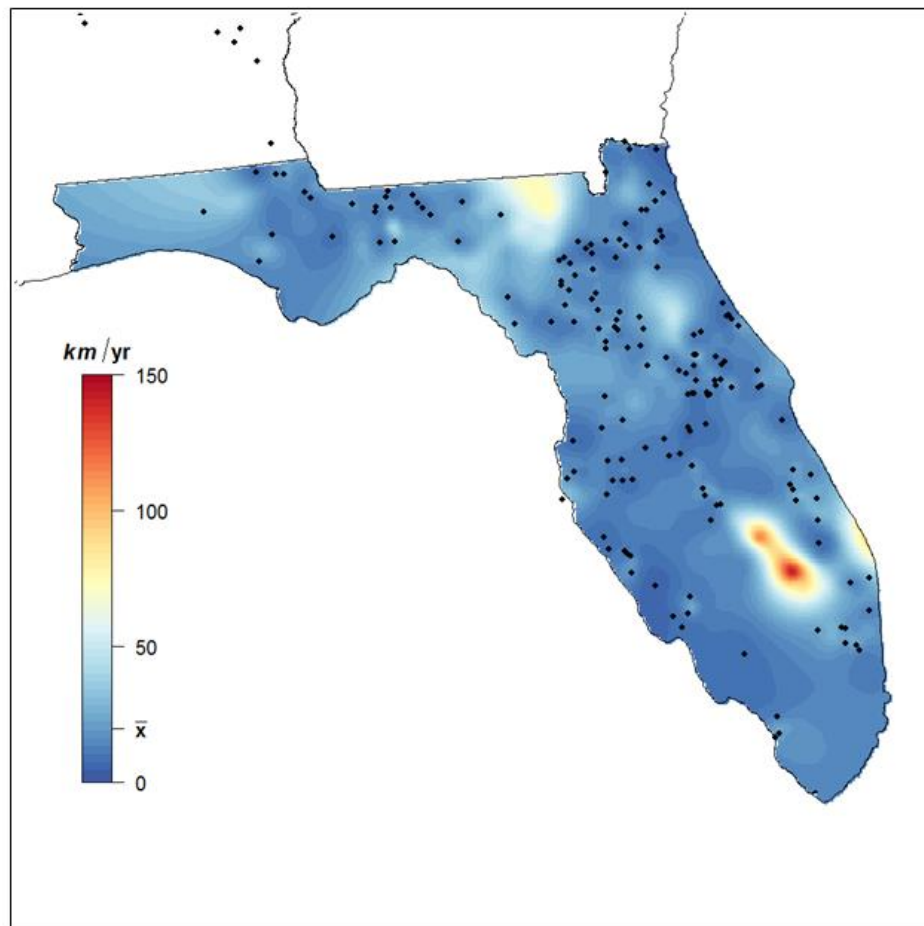
243
244

245
246
247

| Ancestor node | Estimated year of clade divergence [95% HPD] | |
| | N-gene | N+G-genes |
|---|---|---|
| Tree root | 1948 [1934, 1961] | 1957 [1940, 1969] |
| I | 1959 [1947, 1969] | 1962 [1947,1974] |
| II | 1966 [1954, 1976] | 1972 [1957, 1983] |
| III | 1968 [1959, 1973] | 1972 [1962, 1979] |
| IV | 1975 [1965, 1983] | 1985 [1977, 1991] |
| V | 1969 [1959, 1965] | 1971 [1957, 1981] |

248
249
250
251
252
253

254

**Table 1.** Median divergence times of each clade ancestor estimated using a time-scaled evolutionary rate applied to genetic data. All samples were sequenced at the nucleoprotein gene locus (N-gene), while a subset were analyzed at both the N and glycoprotein genes (N+G gene).

255

256 **Spatial patterns of diffusion**

257    The spatiotemporal history of the N-gene phylogeny was reconstructed under a relaxed

258 random walk model in continuous space (Lemey *et al.* 2010). Median viral diffusion rate

259 estimated from posterior trees was 6.47 km/year ($1^{st}$ quartile=3.11, $3^{rd}$ quartile=14.32), a value

260 consistent with rates reported for RRV movement elsewhere (Biek *et al.* 2007; Lemey *et al.*

261 2010). Extracting the temporal and spatial locations of ancestral nodes from the posterior

262 distribution of trees allowed us to interpolate a surface of parent-to-daughter diffusion rates

263 across Florida (Fig 2), while accounting for uncertainty in the phylogeny and the coalescent

264 model used to estimate those rates. The spatial distribution of estimated diffusion rates across

265 the Florida peninsula did not reveal any obvious areas of heterogeneity, even under the relaxed

13

266  random walk model, with the exception of minor clusters of elevated diffusion rates (Fig 2).

267  These clusters corresponded to areas where we lacked samples. We confirmed using simulations

268  that sequences evolved on a continuous landscape yielded similar spatial patterns of diffusion

269  when sampled discontinuously (SI Fig 1c).  We therefore propose that the "hotspots" of

270  increased diffusion rates seen in the empirical data are likely the result of sampling gaps.
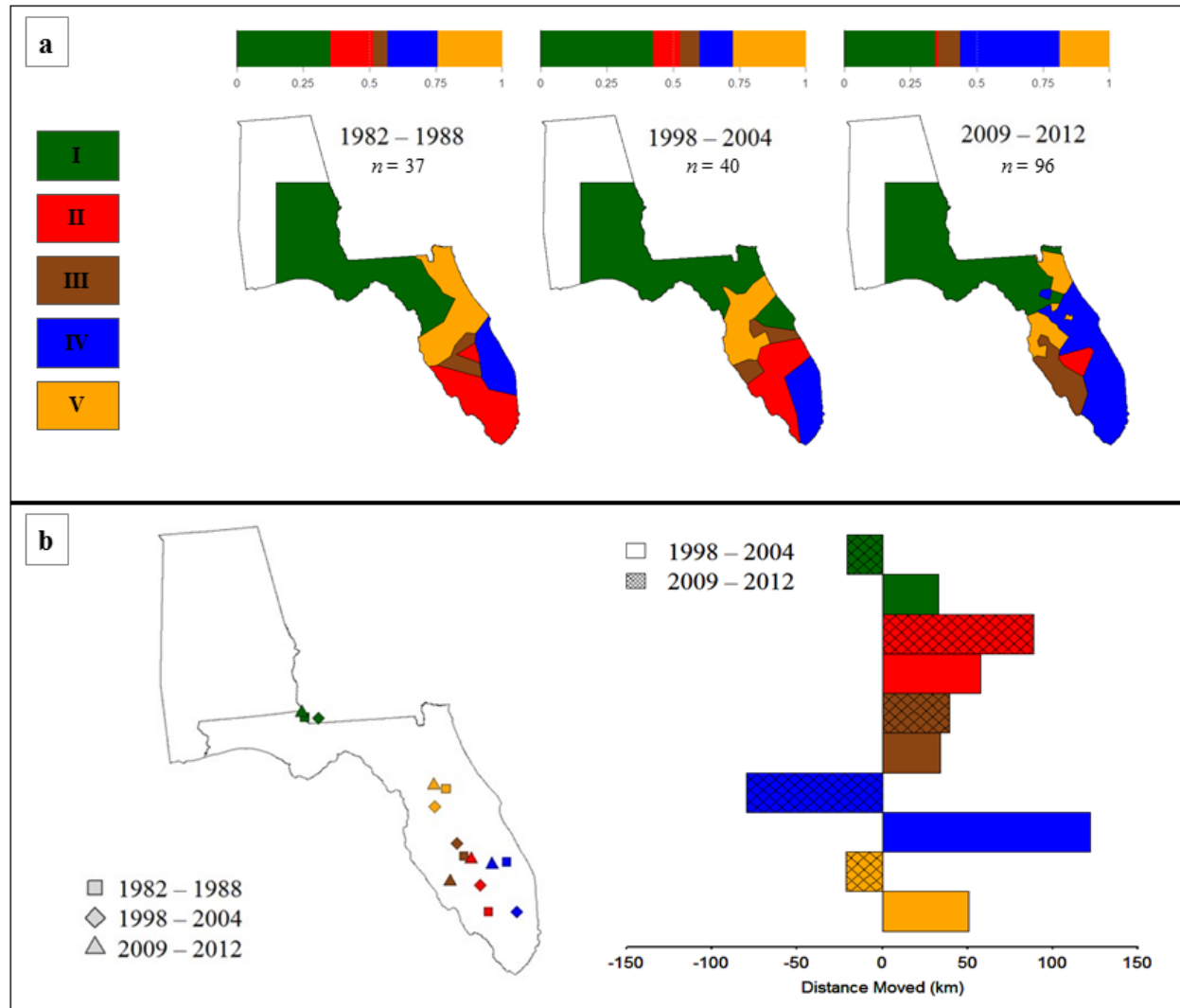


271
272

**Fig 2.  Predicted averaged annual rate of raccoon rabies virus diffusion throughout the Florida peninsula.**  Predicted diffusion was interpolated using spatial coordinates as predictor variables in a smoothed GAM of diffusion rates extracted from the posterior distribution of 1000 phylogenetic trees.  The mean diffusion rate is indicated at x-bar.  Increased diffusion rates are warmer in color.  Sample locations are marked with black circles.

273

14

**Phylogeographic structure**

The inferred clades exhibited strong geographic structure (Fig 1c), as our sample region

can be partitioned into areas dominated by a certain clade.  For example, clade IV is generally

restricted to eastern Florida while sequences grouped in clade I are found in Alabama and the

Florida panhandle (Fig 1c). All five clades can be found in the area around the estimated location

of the tree's root (Fig 1c), positioned at a similar latitude but farther(?) to the west of the earliest

reported case of RRV (Fig 1b). The geographic distribution of the five clades over time was

remarkably stable, with mean centers during the second and third sampling period that remained

within 100km to their estimated locations in the 1982 – 1988 sampling interval (Fig 3a, b).   The

geographic centers of clades I and V moved the least (<30 km total) (Fig 3b). Unlike clade I,

however, clade V appears to have experienced some fragmentation by the expansion of clade IV

(Fig 3a).  The mean center of clade IV shifted a large distance from its estimated location, only

to shift back the following sampling interval (Fig 3a, b).  Only the mean centers of clades II and

III exhibited repeated movement away from their respective locations in the 1982 – 1988

sampling interval (Fig 3a, b).

15

**Fig 3. Stable geographic distribution of raccoon rabies virus (RRV) lineages over time.**
(**A**) Clade boundaries estimated from Delaunay triangulation at each general sampling interval. Colored squares indicate lineage number. The colored bars along the top indicate the proportion of samples assigned to a given lineage in each sampling interval. (**B**) Geometric mean center of each clade's triangulated area at each sampling interval, with the distances each mean center moved from the first sampling interval (1982 – 1988). Positive distances indicate a spatial shift away from the center's location in the first interval, negative distances are movements back toward that location.

**Mid-Atlantic origin**

Included in our phylogenetic analyses was a viral sequence of the reconstructed ancestor

of the mid-Atlantic RRV epizootic (Biek *et al.* 2007); the mid-Atlantic outbreak emerged in

1977 from a point far north of the RRV variant's original range in the southeastern states, likely

16

302     due to translocation of an infected raccoon from Florida to West Virginia (Nettles *et al.* 1979;

303     Rupprecht & Smith 1994).  However, the exact spatial origin of the translocated raccoon has

304     never been clarified.  Our phylogenetic analysis placed the reconstructed ancestral sequence

305     firmly within clade III, for both the N-gene-based phylogeny and the phylogeny constructed

306     using both N and G genes (Figs 1a and SI 2). Both analyses indicated that this virus diverged

307     early from the clade III ancestor. Clade III was one of the last clades to diverge and is restricted

308     to a relatively small area in central western Florida. Our results thus suggest that the mid-Atlantic

309     epizootic was caused by an infected raccoon originating from this area.

**Discussion**

311          The results presented here provide evidence for the development and maintenance of

312     spatial aggregates of rabies viral lineages on a landscape with few features that might impede

313     host racoon movement.  A key characteristic of rabies virus – indeed, of all rapidly-evolving

314     RNA pathogens – is that there is a direct link between viral population dynamics and their

315     molecular evolution, as mutations are fixed at the same tempo as population dynamics occur

316     (Grenfell *et al.* 2004; Holmes 2004, 2009).  This characteristic is central to the expectation that

317     ecological processes are recorded in viral RNA genomes in near "real-time" (Drummond *et al.*

318     2003; Holmes 2008; Holmes 2009; Real *et al.* 2005a), yielding phylogenies that are rich in

319     epidemiological information and predicted by viral ecology (Grenfell *et al.* 2004; Pybus *et al.*

320     2012; Real *et al.* 2005b).  Without landscape features that significantly influence host movement,

321     however, the spatial distribution and subsequent phylogeography of viral populations will be a

322     result of host dynamics that occur at the local scale or that were written into the viral phylogeny

323     during longer-term population processes such as colonization and epizootic expansion.  We

324     reconstructed evolutionary and ecological aspects of RRV back to the time and location of its

325     original emergence in raccoons.  The timescale of our inferred genealogy covers the known time

17

326   period of RRV existence – the first reported case of rabies in raccoons occurred in 1947 (Bigler

327   *et al.* 1973; Scatterday *et al.* 1960) (Fig 1b), which corresponds well with the estimated

328   emergence date of our phylogenetic tree root.  Reconstructed dispersal patterns suggest

329   homogeneous viral diffusion across the state, a result that is consistent with the fact that the

330   terrain of Florida presents few known barriers to raccoon movement, as well as with previous

331   research supporting the presence of a single statewide raccoon population (Cullingham *et al.*

332   2008; Reeder-Carroll 2010).  Early cases of raccoon rabies in Florida exhibited no spatial

333   correlation – initial cases were sporadic and spatially disjoint – but RRV is believed to have

334   covered the state by 1965 (Bigler *et al.* 1973).  Localized outbreaks began to emerge in the late

335   1960s (Bigler *et al.* 1973; Kappus *et al.* 1970).  We identified five distinct viral clades had

336   diverged within ~20 years (by the early 1970s) following the mid-20th century emergence of

337   RRV, with a phylogenetic structure that appears to be dominated by spatial radiation but in the

338   absence of landscape-level effects.  Despite the lack of landscape-level constraints on host

339   movement, the distribution of the RRV lineages following their emergence has remained

340   spatially compartmentalized and remarkably stable throughout the history of raccoon rabies in

341   Florida.

342           Terrestrial rabies viral variants are frequently analyzed within the context of landscape

343   heterogeneity and its effects on viral diffusion processes.  In a few cases, phylogeographic

344   structuring of terrestrial variants has been observed without the apparent influence of geographic

345   barriers:  Kuzmina et al. (2013) noted the presence of several spatially restricted, distinct

346   phylogenetic lineages of skunk-specific rabies in Texas, remarking that there were no "obvious

347   natural barriers restricting virus spread."  An early investigation of fox rabies in Ontario, Canada

348   suggested that viral lineages were segregated among distinct habitats (Nadin-Davis *et al.* 1999),

18

349    however reanalysis of those data against a null model of isolation-by-distance found no evidence

350    of genetic structuring due to ecological variables (Real *et al.* 2005a).  There are many more

351    examples of landscape heterogeneity driving the epidemiology and phylogeographic structure of

352    terrestrial rabies (e.g., Biek *et al.* 2007; Cullingham *et al.* 2009; Smith *et al.* 2002; Talbi *et al.*

353    2010), which have identified a range of environmental predictors of rabies spread.  However, our

354    work shows that the existence of genealogical breaks in the virus phylogeny does not allow to

355    conclude the existence of barriers to viral gene flow and transmission.

356            During population expansion, random events can generate spatial compartments of

357    distinct evolutionary lineages (Cavalli-Sforza *et al.* 1993; Edmonds *et al.* 2004; Hallatschek *et al.*

358    2007; Waters *et al.* 2013).  As new genetic lineages diverge and radiate outward, individuals

359    located at their spatial centers are separated from sister clades by increasing distances, a

360    phenomenon that should help solidify phylogeographic partitions as they form (Irwin 2002;

361    Waters *et al.* 2013).  Beyond colonization processes, it is also possible for phylogeographic

362    structure to develop in established populations:  Irwin (2002) noted that phylogeographic breaks

363    emerged at random locations within simulated established and continuous populations modeled

364    with low dispersal distances, even in the absence of physical barriers to gene flow.  All the viral

365    clades detected here diverged while rabies expanded through a novel host system, suggesting that

366    colonization and expansion processes drove the phylogeographic patterns.  Spatial

367    compartmentalization of the five clades observed here remained stable throughout the history of

368    RRV in our study area, with each clade restricted to distinct geographic areas that changed very

369    little in their apparent arrangement.  Our results suggest that the mean centers of each clade

370    move very little – some shifting of clades' geographic centers was recorded throughout our

371    sampling period, but the general locations of these centers were often preserved over multiple

19

372     decades.  Even clades experiencing fragmentation and/or contraction exhibited a constancy in

373     mean center location, suggesting sampling variation as a possible cause rather than true shifts in

374     distribution. This stability, as well as the absence of lineage turnover, suggest that spatial genetic

375     patterns in rabies virus are preserved through time, long after the initial invasion process. It also

376     highlights the overriding importance of local host movement processes, resulting in limited

377     spatial admixture, in the maintenance of RRV.  Interestingly, Hallatscheck *et al*. (2007)

378     described similar spatial genetic partitions among expanding microbial colonies.  The authors

379     noted the inherent randomness to this process:  The spatial arrangement of microbial diversity

380     was ultimately determined by a few cells at the expanding edge contributing to the clear majority

381     of future generations – and as these colonies were composed of genetically identical individuals

382     expanding across a neutral landscape, it was not selection but random genetic drift that

383     determined which microbial lines went on to dominate the gene pool (Hallatschek *et al.* 2007).

384     Further, as the colonized range increases in size, core areas of distinct genetic lines are more

385     likely to survive through time, buffered from turnover by average dispersal distance being less

386     than the area of a genetic "sector" or spatial compartment (Hallatschek *et al.* 2007; Irwin 2002).

387     Genetic lines or clades that are isolated from the expanding edge of the range, however, are

388     consequently excluded from leaving much trace on the final phylogeny (Hallatschek *et al.* 2007).

389     Using a measurably evolving pathogen, we are able to provide an empirical example of what has

390     previously been demonstrated *in vitro* (Hallatschek *et al.* 2007) or using simulations (Irwin

391     2002).

392         A major goal of population genetics is to determine which of natural selection or neutral

393     processes is driving a population's evolution.  We find this system to be a clear example of

394     neutral processes dominating population genetic structure. The stability of the phylogeographic

20

395    patterns detected here enabled us to recover the potential source location of the mid-Atlantic

396    rabies epizootic that began decades after RRV had reached an enzootic state in Florida.

397    Crucially, we were able to identify this location even though our samples for this study were

398    themselves collected long after RRV establishment in both Florida and in the mid-Atlantic states,

399    as well as to reconstruct the source sequence of the mid-Atlantic epizootic.  The clade from

400    which the mid-Atlantic source appears to originate is not located near the expanding northern

401    edge of RRV in the southeast.  Instead, Clade III is surrounded by either water or by RRV sister

402    clades.  This event was enabled by human-mediated long-distance-dispersal (Nettles *et al.* 1979),

403    without which we would predict that raccoon rabies would have continued its slow expansion up

404    the east coast, and that Clade I – with exclusive access to the expanding edge of RRV in Florida

405    – would have gone on to dominate future lineages.

406            RNA viruses in general have been identified as likely systems for evolutionary patterns

407    driven by dispersal processes rather than selection, as they tend to cause acute infections that

408    preclude co-divergence with their host species (Holmes 2004).  Nevertheless, the rapid rate of

409    rabies evolution has led to the expectation that its spatial genetic patterns will reflect its host

410    population structure.  Here, however, we find a significant and stable lack of spatial genetic

411    congruence between RRV and its host, which displays no similar spatial compartmentalization

412    (Cullingham *et al.* 2008; Reeder-Carroll 2010).  Host contact rates and viral transmission

413    dynamics are different, therefore, from those processes that drive host gene flow and migration

414    patterns.  Additionally, the ecological dynamics recoverable from this viral phylogeny my not

415    have occurred on the same timescale as its sampling scheme.  It is much more likely that the

416    sweeping patterns set up during the disease's initial expansion are being recovered, while the

417    stochastic and short-term dynamics of local movements are too transient to affect detectable

21

418   phylogeographic changes (Carroll *et al.* 2007; Holmes 2004). Raccoon rabies control efforts that

419   rely on landscape variables to assist in restraining viral movement (Elmore *et al.* 2017; Russell *et*

420   *al.* 2005) will potentially not be applicable to systems such as this, and alternative methods of

421   control will be necessary.

**Acknowledgements**

**Data accessibility**
The DNA sequences used for this paper have been submitted to Genbank and are available under accession nos. xxx-xxx.

AUTHOR CONTRIBUTIONS
TM, LAR and RB were involved in study design and concept. RD provided molecular data. TM carried out the analysis. SDS provided analytical tools and technical advice. TM wrote the manuscript with significant contributions from RB. All authors viewed and revised the final manuscript.

**References**

Avise JC, Arnold J, Ball RM*, et al.* (1987) Intraspecific Phylogeography: The Mitochondrial DNA Bridge Between Population Genetics and Systematics. *Annual Review of Ecology and Systematics* **18**, 489-522.

Baele G, Lemey P, Bedford T*, et al.* (2012) Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Molecular Biology and Evolution* **29**, 2157-2167.

Baele G, Li WL, Drummond AJ, Suchard MA, Lemey P (2013) Accurate model selection of relaxed molecular clocks in bayesian phylogenetics. *Molecular Biology and Evolution* **30**, 239-243.

Baptestini EM, de Aguiar MA, Bar-Yam Y (2013) Conditions for neutral speciation via isolation by distance. *Journal of Theoretical Biology* **335**, 51-56.

Biek R, Henderson JC, Waller LA, Rupprecht CE, Real LA (2007) A high-resolution genetic signature of demographic and spatial expansion in epizootic rabies virus. *Proc Natl Acad Sci U S A* **104**, 7993-7998.

Biek R, Pybus OG, Lloyd-Smith JO, Didelot X (2015) Measurably evolving pathogens in the genomic era. *Trends in Ecology & Evolution* **30**, 306-313.

Bigler WJ, McLean RG, Trevino HA (1973) Epizootiologic aspects of raccoon rabies in Florida. *American Journal of Epidemiology* **98**, 326-335.

Brunker K, Hampson K, Horton DL, Biek R (2012) Integrating the landscape epidemiology and genetics of RNA viruses: rabies in domestic dogs as a model. *Parasitology* **139**, 1899-1913.

22

Carroll SP, Hendry AP, Reznick DN, Fox CW (2007) Evolution on ecological time-scales. *Functional Ecology* **21**, 387-393.

Cavalli-Sforza LL, Menozzi P, Piazza A (1993) Demic expansions and human evolution. *Science* **259**, 639-646.

Childs JE, Curns AT, Dey ME, *et al.* (2000) Predicting the local dynamics of epizootic rabies among raccoons in the United States. *Proc Natl Acad Sci U S A* **97**, 13666-13671.

Crisp MD, Trewick SA, Cook LG (2011) Hypothesis testing in biogeography. *Trends in Ecology & Evolution* **26**, 66-72.

Cullingham CI, Kyle CJ, Pond BA, Rees EE, White BN (2009) Differential permeability of rivers to raccoon gene flow corresponds to rabies incidence in Ontario, Canada. *Molecular Ecology* **18**, 43-53.

Cullingham CI, Kyle CJ, Pond BA, White BN (2008) Genetic structure of raccoons in eastern North America based on mtDNA: implications for subspecies designation and rabies disease dynamics. *Canadian Journal of Zoology* **86**, 947-958.

Darriba D, Taboada GL, Doallo R, Posada D (2012) jModelTest 2: more models, new heuristics and parallel computing. *Nature Methods* **9**, 772.

Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG (2003) Measurably evolving populations. *Trends in Ecology & Evolution* **18**, 481-488.

Drummond AJ, Rambaut A (2007) BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evolutionary Biology* **7**, 214.

Drummond AJ, Rambaut A, Shapiro B, Pybus OG (2005) Bayesian coalescent inference of past population dynamics from molecular sequences. *Molecular Biology and Evolution* **22**, 1185-1192.

Drummond AJ, Suchard MA, Xie D, Rambaut A (2012) Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Molecular Biology and Evolution* **29**, 1969-1973.

Duke-Sylvester SM, Biek R, Real LA (2013) Molecular evolutionary signatures reveal the role of host ecological dynamics in viral disease emergence and spread. *Philos Trans R Soc Lond B Biol Sci* **368**, 20120194.

Edmonds CA, Lillie AS, Cavalli-Sforza LL (2004) Mutations arising in the wave front of an expanding population. *Proc Natl Acad Sci U S A* **101**, 975-979.

Elmore SA, Chipman RB, Slate D, *et al.* (2017) Management and modeling approaches for controlling raccoon rabies: The road to elimination. *PLoS Negl Trop Dis* **11**, e0005249.

Excoffier L, Ray N (2008) Surfing during population expansions promotes genetic revolutions and structuration. *Trends in Ecology & Evolution* **23**, 347-351.

Grenfell BT, Pybus OG, Gog JR, *et al.* (2004) Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* **303**, 327-332.

Guindon S, Gascuel O (2003) A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Systematic Biology* **52**, 696-704.

Hallatschek O, Hersen P, Ramanathan S, Nelson DR (2007) Genetic drift at expanding frontiers promotes gene segregation. *Proc Natl Acad Sci U S A* **104**, 19926-19930.

Hallatschek O, Nelson DR (2008) Gene surfing in expanding populations. *Theoretical Population Biology* **73**, 158-170.
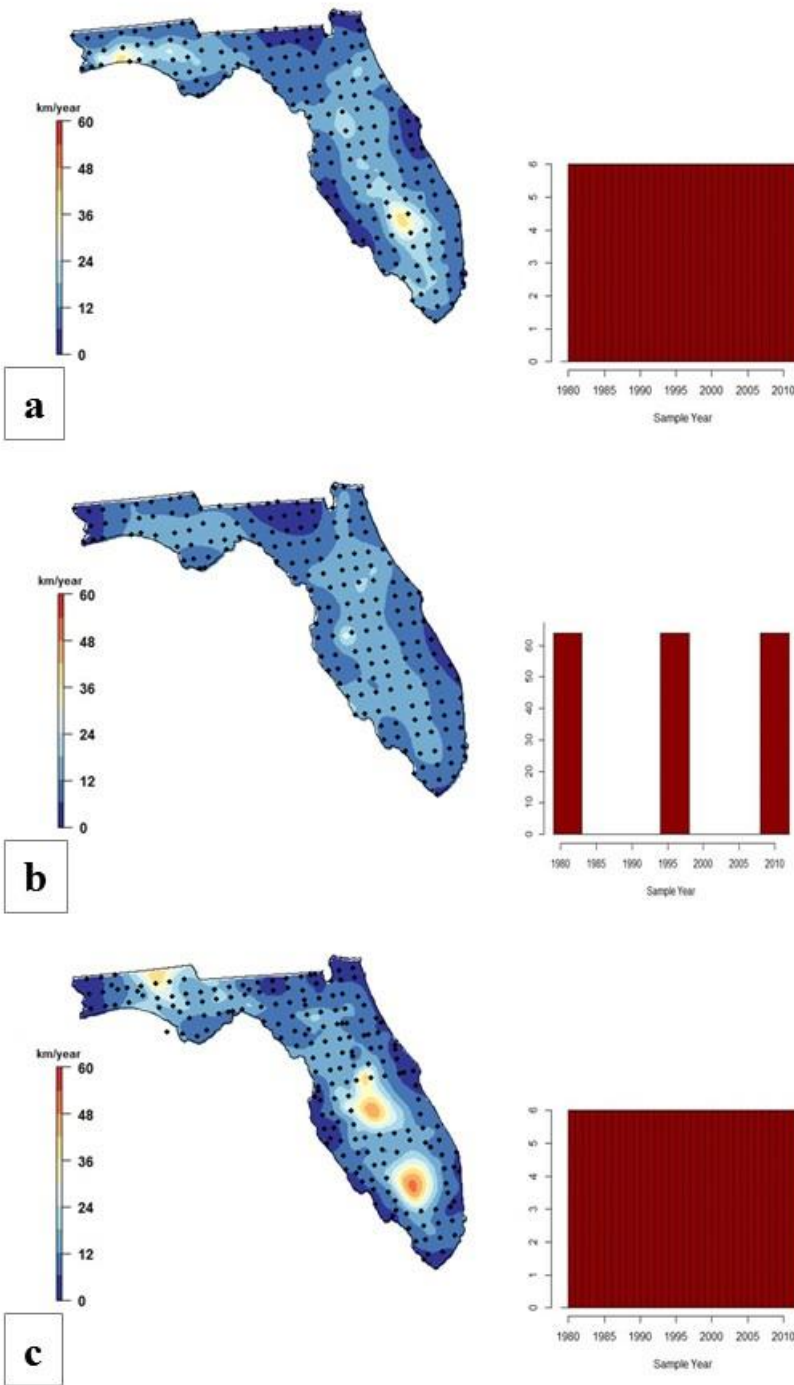
23

506  Hasegawa M, Kishino H, Yano T (1985) Dating of the human-ape splitting by a
507      molecular clock of mitochondrial DNA. *Journal of Molecular Evolution* **22**, 160-
508      174.
509  Holmes EC (2004) The phylogeography of human viruses. *Molecular Ecology* **13**, 745-
510      756.
511  Holmes EC (2008) Evolutionary history and phylogeography of human viruses. In:
512      *Annual Review of Microbiology*, pp. 307-328.
513  Holmes EC (2009) *The evolution and emergence of RNA viruses* Oxford University
514      Press.
515  Ibrahim KM, Nichols RA, Hewitt GM (1996) Spatial patterns of genetic variation
516      generated by different forms of dispersal during range expansion. *Heredity* **77**,
517      282-291.
518  Irwin DE (2002) Phylogeographic breaks without geographic barriers to gene flow.
519      *Evolution* **56**, 2383-2394.
520  Kappus KD, Bigler WJ, McLean RG, Trevino HA (1970) The raccoon an emerging rabies
521      host. *Journal of Wildlife Diseases* **6**, 507-509.
522  Kuo CH, Avise JC (2005) Phylogeographic breaks in low-dispersal species: the
523      emergence of concordance across gene trees. *Genetica* **124**, 179-186.
524  Kuzmina NA, Lemey P, Kuzmin IV*, et al.* (2013) The phylogeography and
525      spatiotemporal spread of south-central skunk rabies virus. *PLOS ONE* **8**, e82348.
526  Lembo T, Haydon DT, Velasco-Villa A*, et al.* (2007) Molecular epidemiology identifies
527      only a single rabies virus variant circulating in complex carnivore communities of
528      the Serengeti. *Proc Biol Sci* **274**, 2123-2130.
529  Lemey P, Rambaut A, Drummond AJ, Suchard MA (2009) Bayesian phylogeography
530      finds its roots. *PLoS Comput Biol* **5**, e1000520.
531  Lemey P, Rambaut A, Welch JJ, Suchard MA (2010) Phylogeography takes a relaxed
532      random walk in continuous space and time. *Molecular Biology and Evolution* **27**,
533      1877-1885.
534  Manel S, Schwartz MK, Luikart G, Taberlet P (2003) Landscape genetics: combining
535      landscape ecology and population genetics. *Trends in Ecology & Evolution* **18**,
536      189-197.
537  Marske KA, Rahbek C, Nogués-Bravo D (2013) Phylogeography: spanning the ecology-
538      evolution continuum. *Ecography* **36**, 1169-1181.
539  Miller MA, Pfeiffer W, Schwartz T (2010) Creating the CIPRES Science Gateway for
540      inference of large phylogenetic trees, 1-8.
541  Nadin-Davis SA, Casey GA, Wandeler AI (1994) A molecular epidemiological study of
542      rabies virus in central Ontario and western Quebec. *Journal of General Virology*
543      **75 ( Pt 10)**, 2575-2583.
544  Nadin-Davis SA, Sampath MI, Casey GA, Tinline RR, Wandeler AI (1999)
545      Phylogeographic patterns exhibited by Ontario rabies virus variants.
546      *Epidemiology and Infection* **123**, 325-336.
547  Neigel JE, Avise JC (1993) Application of a random walk model to geographic
548      distributions of animal mitochondrial DNA variation. *Genetics* **135**, 1209-1220.
549  Nettles VF, Shaddock JH, Sikes RK, Reyes CR (1979) Rabies in translocated raccoons.
550      *American Journal of Public Health* **69**, 601-602.
551  Orciari LA, Niezgoda M, Hanlon CA*, et al.* (2001) Rapid clearance of SAG-2 rabies virus
552      from dogs after oral vaccination. *Vaccine* **19**, 4511-4518.

24

Pybus OG, Rambaut A (2009) Evolutionary analysis of the dynamics of viral infectious disease. *Nature Reviews: Genetics* **10**, 540-550.

Pybus OG, Suchard MA, Lemey P*, et al.* (2012) Unifying the spatial epidemiology and molecular evolution of emerging epidemics. *Proc Natl Acad Sci U S A* **109**, 15066-15071.

Real LA, Henderson JC, Biek R*, et al.* (2005a) Unifying the spatial population dynamics and molecular evolution of epidemic rabies virus. *Proc Natl Acad Sci U S A* **102**, 12107-12111.

Real LA, Russell C, Waller L, Smith D, Childs J (2005b) Spatial dynamics and molecular ecology of North American rabies. *Journal of Heredity* **96**, 253-260.

Reeder-Carroll SA (2010) *Population genetics of raccoons in the Eastern United States with implications for rabies transmission and spread* Ph.D., Emory University.

Rupprecht CE, Smith JS (1994) Raccoon rabies: the re-emergence of an epizootic in a densely populated area. *Seminars in Virology* **5**, 155-164.

Russell CA, Smith DL, Childs JE, Real LA (2005) Predictive spatial dynamics and strategic planning for raccoon rabies emergence in Ohio. *PLoS Biology* **3**, e88.

Saunders NC, Kessler LG, Avise JC (1986) Genetic Variation and Geographic Differentiation in Mitochondrial DNA of the Horseshoe Crab, LIMULUS POLYPHEMUS. *Genetics* **112**, 613-627.

Scatterday JE, Schneider NJ, Jennings WL, Lewis AL (1960) Sporadic animal rabies in Florida. *Public Health Reports* **75**, 945-953.

Schwartz MK, McKelvey KS (2008) Why sampling scheme matters: the effect of sampling scheme on landscape genetic results. *Conservation Genetics* **10**, 441-452.

Shapiro B, Rambaut A, Drummond AJ (2005) Choosing appropriate substitution models for the phylogenetic analysis of protein-coding sequences. *Molecular Biology and Evolution* **23**, 7-9.

Smith DL, Lucey B, Waller LA, Childs JE, Real LA (2002) Predicting the spatial dynamics of rabies epidemics on heterogeneous landscapes. *Proc Natl Acad Sci U S A* **99**, 3668-3672.

Smith JS, Fishbein DB, Rupprecht CE, Clark K (1991) Unexplained rabies in three immigrants in the United States. A virologic investigation. *N Engl J Med* **324**, 205-211.

Smith JS, Reid-Sanden FL, Roumillat LF*, et al.* (1986) Demonstration of antigenic variation among rabies virus isolates by using monoclonal antibodies to nucleocapsid proteins. *Journal of Clinical Microbiology* **24**, 573-580.

Szanto AG, Nadin-Davis SA, Rosatte RC, White BN (2011) Genetic tracking of the raccoon variant of rabies virus in eastern North America. *Epidemics* **3**, 76-87.

Talbi C, Lemey P, Suchard MA*, et al.* (2010) Phylodynamics and human-mediated dispersal of a zoonotic virus. *PLoS Pathog* **6**, e1001166.

Towns J, Cockerill T, Dahan M*, et al.* (2014) XSEDE: accelerating scientific discovery. *Computing in Science & Engineering* **16**, 62-74.

Waters JM (2011) Competitive exclusion: phylogeography's 'elephant in the room'? *Molecular Ecology* **20**, 4388-4394.

Waters JM, Fraser CI, Hewitt GM (2013) Founder takes all: density-dependent processes structure biodiversity. *Trends in Ecology & Evolution* **28**, 78-85.
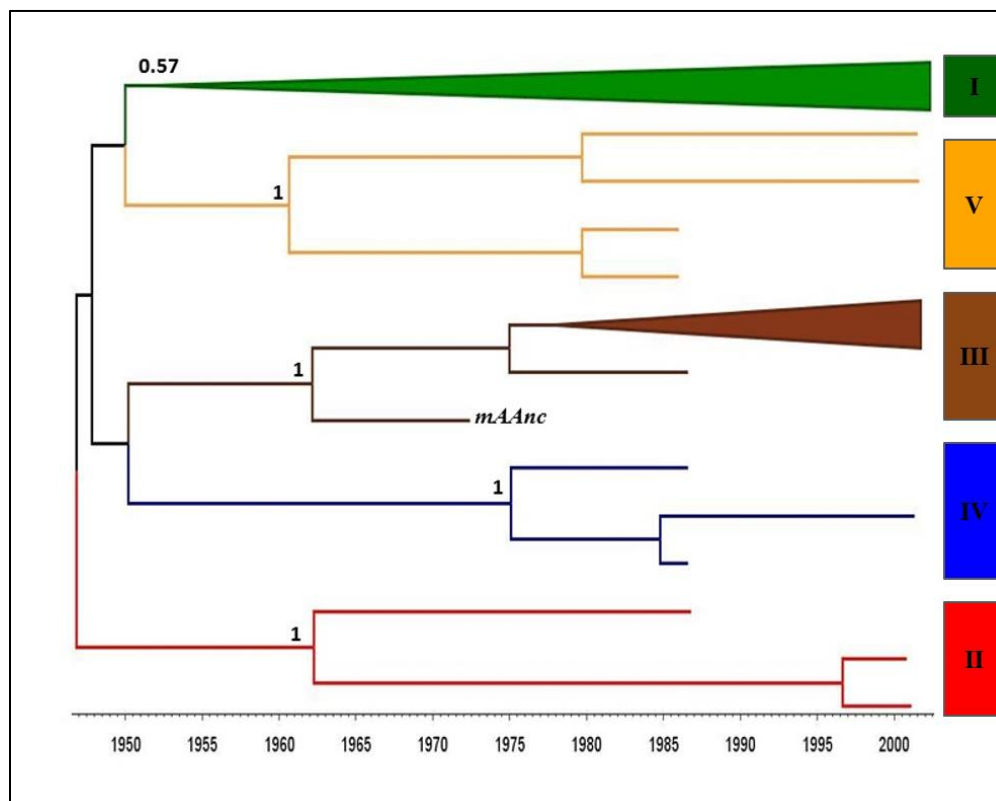
25

599  Wheeler DC, Waller LA (2008) Mountains, valleys, and rivers: The transmission of
600        raccoon rabies over a heterogeneous landscape. *J Agric Biol Environ Stat* **13**,
601        388-406.
602  Wilson ML, Bretsky PM, Cooper GH, Jr.*, et al.* (1997) Emergence of raccoon rabies in
603        Connecticut, 1991-1994: spatial and temporal characteristics of animal infection
604        and human contact. *American Journal of Tropical Medicine and Hygiene* **57**,
605        457-463.
606  Wood SN (2011) Fast stable restricted maximum likelihood and marginal likelihood
607        estimation of semiparametric generalized linear models. *Journal of the Royal
608        Statistical Society: Series B (Statistical Methodology)* **73**, 3-36.
609  Wright S (1943) Isolation by Distance. *Genetics* **28**, 114-138.
610  Wright S (1951) The genetical structure of populations. *Ann Eugen* **15**, 323-354.
611  Yang Z (1994) Maximum likelihood phylogenetic estimation from DNA sequences with
612        variable rates over sites: approximate methods. *Journal of Molecular Evolution*
613        **39**, 306-314.
614  Yang Z (2007) PAML 4: phylogenetic analysis by maximum likelihood. *Molecular
615        Biology and Evolution* **24**, 1586-1591.

616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645

26

**Supplemental Information**



a



b



c

*SI 1*

27

692  **SI Fig 1.  Diffusion rates estimated from simulated sequence evolution on homogeneous**
693  **landscapes and sampled under different schemes**.  Host movement was modeled on a regular
694  grid clipped to the shape of Florida, Georgia, and Alabama, with an index infection introduced at
695  the site of the first reported case of RRV (Fig 1b) following five years of host movement.  Host
696  movement served as the backdrop for disease transmission and sequence evolution for a period
697  of 75 years (see Duke-Sylvester et. al (2013) for details), after which sequences were sampled
698  from host locations throughout Florida during the final 35 years of the model.  Sampled
699  sequences were analyzed in BEAST to estimate ancestral locations and diffusion rates for spatial
700  interpolation of disease diffusion rates using a GAM.  Sequences were collected (**A**) uniformly
701  through space and time, (**B**) uniformly through space but discontinuously through time, and (**C**)
702  continuously through time but discontinuously in space.

703
704
705
706



707
*SI 2*

710
711  **SI Fig 2.  Evolutionary relationships of raccoon rabies virus among 20 samples using**
712  **concatenate N and G gene sequences (1965bp)**.  Maximum clade credibility (MCC) tree scaled
713  to time, with posterior probabilities of each clade's ancestor at nodes (relationships are collapsed
714  where posterior values are less than 0.95).  The sequence of the reconstructed ancestor of the
715  mid-Atlantic racoon rabies epizootic is indicated at individual *mAAnc*.
716

28