1

1   **Title:** The contribution of epidemiological predictors in unravelling the

2   phylogeographic history of HIV-1 subtype C in Brazil

3   **Running title:** The phylogeography of HIV-1C in Brazil

4   Tiago Gräf[1,2]#, Bram Vrancken[3], Dennis Maletich Junqueira[2,4,5], Rúbia Marília de

5   Medeiros[2,5], Marc A. Suchard[6,7], Philippe Lemey[3], Sabrina Esteves de Matos

6   Almeida[2,5], Aguinaldo Roberto Pinto[1]

7

8   1. Laboratório de Imunologia Aplicada, Departamento de Microbiologia,

9   Imunologia e Parasitologia, Universidade Federal de Santa Catarina,

10  Florianópolis, SC, Brazil; 2. Centro de Desenvolvimento Científico e Tecnológico,

11  Fundação Estadual de Produção e Pesquisa em Saúde, Porto Alegre, RS, Brazil; 3.

12  Department of Microbiology and Immunology, KU Leuven, Leuven, Belgium; 4.

13  Departamento de Ciências da Saúde, Uniritter Laureate International

14  Universities, Porto Alegre, RS, Brazil; 5. Programa de Pós-Graduação em Genética

15  e Biologia Molecular, Universidade Federal do Rio Grande do Sul , Porto Alegre,

16  RS, Brazil; 6. Departments of Biomathematics and Human Genetics, David Geffen

17  School of Medicine at UCLA, University of California, Los Angeles, CA, USA; 7.

18  Department of Biostatistics, UCLA Fielding School of Public Health, University of

19  California, Los Angeles, CA, USA.

20

21  #Address correspondence to Tiago Gräf, akograf@gmail.com

22    **Abstract**

23    The phylogeographic history of the Brazilian HIV-1 subtype C (HIV-1C) epidemic

24    is still unclear. Previous studies have mainly focused on the capital cities of

25    Brazilian federal states and the fact that HIV-1C infections increase at a higher

26    rate than subtype B in Brazil calls for a better understanding of the process of

27    spatial spread. A comprehensive sequence dataset sampled across 22 Brazilian

28    locations was assembled and analyzed. A Bayesian phylogeographic generalized

29    linear model approach was used to reconstruct the spatiotemporal history of

30    HIV-1C in Brazil considering several potential explanatory predictors of the viral

31    diffusion process. Analyses were performed on several subsampled datasets in

32    order to mitigate potential sample biases. We reveal a central role for Porto

33    Alegre city, the capital of the southernmost state, in the Brazilian HIV-1C

34    epidemic (HIV-1C_BR), and the northwards expansion of HIV-1C_BR could be

35    linked to source populations with higher HIV-1 burden and larger proportions of

36    HIV-1C infections. The results presented here bring new insights to the

37    continuing discussion about the HIV-1C epidemic in Brazil, and raise an

38    alternative hypothesis for its spatiotemporal history. The current work also

39    highlights how sampling bias can confound phylogeographic analyses and

40    demonstrates the importance of incorporating external information to protect

41    against this.

42    **Importance:** Subtype C is responsible for the largest HIV infection burden

43    worldwide, but our understanding of its transmission dynamics remains

44    incomplete. Brazil witnessed a relatively recent introduction of HIV-1C

45    compared to HIV-1B, but it swiftly spread throughout the South, where it now

46    circulates as the dominant variant. The northward spread is comparatively

47 slower and HIV-1B still prevails in this region. While epidemiological data and

48 viral genetic analyses have both independently shed light on the dynamics of

49 spread in isolation, their combination has not yet been explored. Here, we

50 complement publically available sequences and new genetic data from 13 cities

51 with epidemiological data to reconstruct the history of HIV-1C spread in Brazil.

52 The combined approach results in more robust reconstructions and can protect

53 against sampling bias. We found evidence for an alternative view on the HIV-1C

54 spatiotemporal history in Brazil, which, contrary to previous explanations,

55 integrates seamlessly with other observational data.

56

57 **Key Words:** Brazil; HIV-1 subtype C; phylogeography; generalized linear

58 models; epidemiological predictors.

59

60 **Introduction**

61    Upon emergence into the human population in central Africa around 1920

62 (1), HIV-1 group M has diversified into nine subtypes and numerous circulating

63 recombinant forms (CRFs) through a series of founder effects and recombination

64 events (2,3). Although HIV-1 subtype B (HIV-1B) dominates in many countries in

65 Europe and Americas (2), more than 50% of the infections worldwide are caused

66 by HIV-1 subtype C (HIV-1C), which is the most prevalent subtype in southern

67 Africa countries and India, and is increasing in prevalence in China and South

68 America (2,4).

69    The epidemic in Brazil is mainly driven by HIV-1B, followed by lower

70 frequencies of HIV-1C, F1 and BF1 recombinants (5). In the southern regions of

71 Brazil, however, the spread of HIV-1B is matched by HIV-1C, which co-circulates

72 in similar proportions and can even be responsible for up to 80% of the

73 infections (4). In addition, the two southernmost Brazilian states, Rio Grande do

74 Sul (RS) and Santa Catarina (SC), have the highest AIDS incidence in the country

75 (6).  These patterns have motivated several investigations into the history and

76 dynamics of the Brazilian HIV-1C epidemic (HIV-1C_BR), which estimated an

77 origin in the 1960s-1970s in the state of Paraná (PR) (7,8,9). Because the HIV-1C

78 incidence in more northern states has only recently begun to increase (4,10-14),

79 this suggests viral diffusion would be driven by unknown factors that promote a

80 fast dissemination to the south while constraining spread to the north.

81    HIV infections are characterized by a dynamic viral population of closely

82 related variants that can quickly adapt to changing selective pressures, which

83 manifests in a formidable speed at which genetic diversity accumulates within

84 hosts (15). This rapid accumulation of genetic diversity makes HIV-1 a prime

85    example of 'measurably evolving populations' (16). As a consequence, there has

86    been an important role for phylogenetic tools to shed light on the

87    epidemiological history of HIV. In fact, this has stimulated many developments in

88    the field of statistical phylodynamics, such as molecular clock models to

89    incorporate sampling time as calibration information (17, 18) and coalescent

90    models to infer the changes in viral population size over time (19-21). More

91    recently, such genealogy-based population genetic inferences have also been

92    complemented with state-of-the-art phylogeographic tools (22-24).

93    Phylodynamic analyses of HIV-1 have culminated in a relatively rich statistical

94    account of its evolutionary and epidemiological history (1).

95         While these statistical models and inference tools have proven invaluable

96    for testing hypotheses using virus genetic data (25, 26), they are limited in their

97    ability to link epidemic processes to pathogen evolution because non-genetic

98    data is usually not directly incorporated into the models.  For phylogeography,

99    this has recently been addressed by extending a Bayesian discrete phylogenetic

100   diffusion approach in order to incorporate covariates in the process of spread

101   (27). This approach estimates phylogeographic history while identifying the

102   contribution of several potential explanatory variables (predictors) of spatial

103   spread and allows for cross-talk between the spatial genetic distribution and the

104   relevant predictors: the predictors are selected for the ability to explain the

105   location transition history, but by helping to inform the process parameters they

106   can also assist in shaping the ancestral reconstructions. This approach has

107   already proven useful to elucidate the drivers of both human and swine influenza

108   dispersal (27, 28).

109     In the present study we reconstruct the phylogeographic history of HIV-

110     1C in Brazil incorporating newly obtained sequence data. While previous studies

111     mostly included sequences from state capital cities, we here expanded the spatial

112     sampling by including HIV-1C sequences from 10 rural locations in the

113     southernmost states RS and SC. Our study demonstrates for the first time that

114     augmenting the molecular sequence data with relevant epidemiological

115     information can contribute to the robustness of phylogeographic

116     reconstructions.

117

118     **Methods**

119     **Patients, samples and new sequences:**

120     A total of 360 HIV-1 seropositive patients from 13 cities from the states of

121     SC and RS (Figure 1) were enrolled in this study, which was approved by the

122     ethics committees of the Federal University of Santa Catarina and the Foundation

123     of Research and Production in Health of the Rio Grande do Sul state. Between

124     October 2009 and February 2014 blood samples were collected and HIV-1

125     envelope (HXB2 6846-7360 bp) and polymerase (HXB2 2274-3545 bp)

126     fragments were amplified from whole cellular DNA by nested-PCR and

127     sequenced as described elsewhere (29).   Sequences were subtyped using the

128     REGA, RIP and SCUEAL online subtyping tools (30-32) and by performing

129     maximum likelihood phylogenetic inference incorporating HIV-1 subtype

130     reference sequences available from the Los Alamos HIV sequence database

131     ([www.hiv.lanl.gov](www.hiv.lanl.gov)).   Recombinant   sequences   were   identified   through

132     bootscanning analysis using Simplot 3.5.1 (33) (see Supplementary Text S1, for

133     methodological details). The sequences generated in the present study were

134   deposited in GenBank under accession numbers: KR065788-KR066336 and

135   KP224476-KP224501.

136

**Sequence dataset compilation:**

138       Briefly, we complemented our new sequence data with all publically

139   available Brazilian HIV-1C sequences (HIV-1C_BR) from *pol* and *env* genes ($n$ =

140   385). Non-Brazilian HIV-1C sequences were selected using BLAST+ (34). To this

141   purpose, we created a local BLAST database that contained all HIV1-C sequences

142   minus those from Brazil. We performed a similarity search on this database using

143   every HIV-1C_BR sequence as a query, and the 50 best hits for each search were

144   logged. Duplicate entries were removed from these hits and compiled as the

145   international database.  After extensive data cleaning this resulted in datasets with

146   1522 *pol* and 621 *env* sequences that were down-sampled to around 500

147   sequences each to reduce the computational burden in subsequent Bayesian

148   statistical analyses (see Supplementary Text S1, for full details of the followed

149   procedure). Six additional sub-samplings containing only Brazilian sequences

150   were made for *pol* and *env*, to allow assessing the robustness of the

151   phylogeographic reconstructions (see below). For this we aimed at reducing

152   sampling bias by creating three random down-samples in two groups: a) Rand10

153   - with a maximum of ten sequences by location and b) Rand20 - with a maximum

154   number of 20 sequences by location (see Supplementary Table S1, for HIV-1C_BR

155   sequences in each dataset).

156

**Phylogenetic divergence time estimation and population dynamics**

**inference**

Time-scaled phylogenetic trees were reconstructed using a Bayesian

inference method implemented in the BEAST/BEAGLE software (35,36). All

analyses were performed using the GTR+I+$\Gamma_4$ nucleotide substitution model and

an uncorrelated lognormal relaxed molecular clock under the Bayesian Skyride

coalescent model (18,37). Due to the low temporal signal of the datasets, the use

of an informative prior on the tMRCA of the Brazilian subype C clade was

required. For this purpose, we specified a normal distribution with mean (1976)

and standard deviation (5.1) based on previous estimates of the time of

introduction of subtype C in Brazil (8). When exact sampling dates were

unknown, the dates were integrated out over a known sampling time interval

(38). Monte Carlo Markov Chains (MCMC) were run sufficiently long to ensure

stationarity and adequate effective sample size (ESS >200) as diagnosed by

Tracer (http:// beastbioedacuk/Tracer). Maximum clade credibility (MCC) trees

were summarized using the TreeAnnotator tool and visualized in Figtree v1.4.0

(http://tree.bio.ed.ac.uk/software/figtree/). A representative sample of 1000

trees was collected and used as an empirical tree distribution for estimating the

virus migration processes (see below). To ensure that subsequent

phylogeographic analyses are based on histories specific to Brazil we pruned

sequences clustering outside the HIV-1C_BR cluster or non-Brazilian sequences

clustering inside HIV-1C_BR cluster from these trees using PAUP (http://

http://paup.csit.fsu.edu) (see Text S1, for methodological details).

180

181     **Phylogeny-trait association:**

182          We tested for significant phylogenetic clustering by location in two

183     different ways. First, we calculated the Association Index (AI), Parsimony Score

184     (PS) and Monophyletic Clade (MC) measures using BaTS (39). For our second

185     approach we introduced the use of the path sampling (PS) and stepping-stone

186     (SS) sampling marginal likelihood estimators as implemented in BEAST (40,41) to

187     evaluate the extent to which the topology is required as a correlation structure to

188     explain the traits. For this we specified a discrete symmetric (reversible) model of

189     location transitioning and incorporate Bayesian stochastic search variable

190     selection (BSSVS) procedure (22) fitting the trait diffusion process on: 1) a fixed

191     MCC tree summarized from the Bayesian phylogenetic analysis of the complete

192     dataset; 2) a star-like tree with the same trait annotations as in the MCC tree.

193

194     **Phylogeographic inference with epidemiological predictors**

195          To assess the impact of potential explanatory variables (predictors) of the

196     viral diffusion process on phylogeographic reconstructions, we made use of the

197     recent generalized linear model (GLM) extension of Bayesian discrete

198     phylogeographic models (27). This allows reconstructing the spatial diffusion

199     history throughout the tree while simultaneously evaluating the contribution of

200     various potential predictors. Support for predictors is estimated using a BSSVS

201     procedure, and the contribution of each predictor is quantified as a GLM

202     coefficient that has an impact (effect size) in the transition rate among the

203     locations.

204          Using this approach, we tested the following predictors (see Text, SDC 1,

205     for methodological details):

206    1) Geographic distance: the great-circle distance among each pair of cities;

207    2) Passenger air traffic: the number of passengers traveling between each pair of

208       airports;

209    3) HIV population size: the total number of AIDS notifications in a period of 10

210       years reported in each city;

211    4) HIV prevalence: (HIV population size / city population size) X 100,000

212       habitants;

213    5) HIV-1C population size: HIV population size times the proportion of HIV-1C as

214       reported in the literature (4,10-14);

215    6) HIV-1C prevalence: (HIV-1C population size / city population size) X 100,000

216       habitants;

217    7) Sample size: the number of sequences by location.

218       Because not all sampling locations have an airport, we specified a

219    different geographic partitioning for evaluating predictor 2 (passenger air

220    traffic). This partitioning is not well suited for the epidemiological predictors,

221    which led us to test predictor 2 in separate analyses including only sample size

222    as an additional potential predictor.

223       GLM analyses were run in BEAST using previously recommended prior

224    specifications on the set of empirical trees obtained by the Bayesian phylogenetic

225    analysis (27). Bayes Factors (BF) were calculated to determine the support for

226    the inclusion of each predictor in the model and predictor contributions are

227    reported as effect sizes conditional on the effect being included in the model.

228       A phylogeographic analysis with BSSVS was performed with asymmetric

229    transition rates being informed by the predictors supported by the GLM analysis.

230    In other words, for each sub-sampled dataset, we used the rate estimates for

231   prior specification based on the corresponding GLM analysis. SPREAD software

232   was used to identify the well-supported transition rates based on BFs > 3 (42).

233   We complemented this analysis with Markov jump estimation of the number of

234   location transitions throughout the evolutionary history (43). RStudio

235   (http://www.rstudio.org/) was used to calculate the Bayes factors and effect

236   sizes, and to summarize the posterior densities of the highly supported

237   transitions from the BEAST log files.

238

239   **Results**

240   **Sequence dataset compilation**

241        We sequenced 140 *pol* and 202 *env* HIV-1C isolates in 13 locations in SC

242   and RS states, 10 of which have not been sampled before (Figure 1). By

243   combining the generated sequence data with publicly available Brazilian and

244   international HIV-1C sequences we were able to compile comprehensive *pol* and

245   *env* based datasets for reconstructing the spatiotemporal history of HIV-1C in

246   Brazil. In summary, the complete *pol* dataset contained 380 Brazilian and 120

247   international sequences while the *env* dataset totaled 293 Brazilian and 170

248   international sequences (see Supplemental Dataset 1, for complementary

249   information about sequences retrieved from public databanks). The Brazilian *pol*

250   sequences are distributed over 21 locations and the *env* sequences represent 17

251   locations totalizing 22 locations represented with *pol* or *env* sequences, most of

252   them in SC and RS (15/21 for *pol* and 14/17 for *env*). Considering the complete

253   Brazilian dataset, sequences represent the time period between 2002 and 2014

254   (see Supplementary Table S1).

255

256 **Phylogeny-trait association**

257     Because our preliminary analyses suggested a considerable degree of

258 phylogenetic mixing by location, we formally tested whether the datasets

259 containing only Brazilian sequences still supported spatial population structure.

260 The hypothesis of a panmictic population could be rejected for the *pol* and *env*

261 datasets by the association index (AI) and parsimony score (PS) statistics

262 ($p$<0.05), but the monophyletic clade (MC) scores revealed that for 12/21 (57%)

263 of *pol* locations and 12/17 (71%) of *env* locations random clustering could not be

264 rejected (see Supplementary Table S2, for MC scores). The results of the

265 approach based on model testing also provided strong support against the

266 absence of phylogenetic association by sampling location in the *pol* and *env*

267 datasets (Bayes factors of 74 and 39 respectively).

268

269 **Inconsistencies in root state estimates**

270     The results of the phylogeographic reconstruction showed, with strong

271 agreement between most datasets and models applied (50/56 analyses), that the

272 epidemic ignited in SC or RS. Its exact location of introduction could, however,

273 not be unambiguously determined using only virus genetic data. Whereas in the

274 complete *pol* and *env* datasets Florianópolis (FLP) was consistently estimated as

275 the most likely location at the root, other cities - most notably Porto Alegre

276 (POA) (7/48) and Criciúma (CRI) (7/48) - were implicated in 60% (29/48) of the

277 analyses based on the Rand10 and Rand20 subsampled datasets (Table 1).

278

279 **Predictors of viral spread**

280    Using a phylogeographic GLM approach, we evaluated which measures

281    predict the rates of location exchange in the complete and subsampled datasets

282    (Table 2). In the *pol* and *env* complete datasets only the origin and destination

283    sample size yielded strong Bayes factor (BF) support, reflecting the fact that only

284    sample sizes and their heterogeneity are needed to explain the number of

285    location transitions. We also found strong support for destination sample size in

286    the model as a predictor in all *pol* and *env* subsampled datasets (Rand10 and

287    Rand20). This indicates that despite the more homogeneous distribution of

288    sequences by sampling locations in these subsampled datasets, the remaining

289    heterogeneity still has an impact on the phylogeographic reconstructions.

290    Two predictors, "origin HIV prevalence" and "origin subtype C population

291    size", were included in all *pol* and *env* Rand10 and Rand20 datasets with Bayes

292    factor estimates ranging from moderate (BF=6) to strong (BF=39) support and

293    with positive mean conditional effect sizes (Figure 2). Hence, locations with

294    higher HIV prevalence and larger HIV-1C populations tend to act as sources for

295    onwards spread.

296    In addition to epidemiological predictors we also tested geographical

297    distance or air transportation data (in a separate analysis, data not shown) as a

298    predictors of HIV-1C diffusion, but these did not result in noticeable support by

299    any of the analyzed datasets.

300    Interestingly, incorporating relevant epidemiological information into the

301    phylogeographic reconstructions resulted in consistent root state estimates:

302    using the GLM model we find POA as the modal root state in all (12/12) *pol* and

303    *env* Rand10 and Rand20 datasets. Only in the complete datasets, where the

304     sampling bias is more severe, Florianópolis (FLP) was still estimated as the

305     modal root state.

306          To assess the robustness of the phylogeographic reconstructions with

307     respect to the root height prior (see Methods), we also performed the ancestral

308     reconstruction using genealogies estimated under priors that specified a mean

309     tMRCA that was 10 years older and younger respectively. We find that

310     differences in tree depths did not impact the outcome: POA is consistently the

311     modal root state and the same predictors find substantial Bayes factor support in

312     all *pol* and *env* Rand10 and Rand20 datasets of the extended and shortened

313     histories.

314

315     **Porto Alegre as a central hub of the HIV-1C epidemic**

316          We subsequently estimated the most likely migration patterns using an

317     asymmetrical phylogeographic analysis with BSSVS and priors on the location

318     exchange rate priors that are based on the GLM rate estimates. The robustness of

319     the ancestral reconstructions was somewhat lower because in this analysis the

320     predictors can only influence the analysis through the prior specification: POA

321     was found to be the root state location in 10/12 *pol* and *env* Rand10 and Rand20

322     datasets (data not shown). Nonetheless, POA was strongly linked to all other

323     locations (Bayes factors $\geq$ 3) while only a few additional well-supported

324     transitions were found. Because this suggests a central role for POA in the

325     Brazilian HIV-1C dissemination, we address its role in more detail.

326          The arrival of HIV-1C in POA was estimated in 1973 (1966 – 1980, 95%

327     HPD) for *pol* and 1971 (1963-1978, 95%HPD) for *env*, and the spread to other

328     cities started around 1980. The timing of these events reveals a consistent

329 pattern. Nearby locations within RS (Rio Grande and Uruguaiana cities) were

330 initially affected, followed by export to south and southeast state capitals in the

331 early 1980s (e.g. to Florianópolis, Curitiba, Rio de Janeiro and São Paulo). More

332 distant locations were affected at a later stage, first in Central-West region

333 (Campo Grande) in middle 1980s and later in North region (Palmas and Manaus)

334 in the late 1980s and early 1990s.  Only two exceptions to this pattern are found

335 (one in the *pol* and one in the *env* datasets): the capital city Goiânia, where HIV-

336 1C appears to have been introduced from POA in 1981 (*pol*), and Rio de Janeiro,

337 where the introduction of HIV-1C has been more recent according to the *env*

338 datasets (see Supplementary Table S3, for time of first transition from Porto

339 Alegre).

340       More insights into the temporal pattern of spread were obtained by

341 mapping the density of location transitions from POA to the other state capital

342 cities through time. This reveals a period of higher density of viral influx 25 to 30

343 years ago to the South region capital cities Florianópolis and Curitiba. Among the

344 sampled capitals in the Southeast region, a similar pattern emerged for Rio de

345 Janeiro but there is a more evenly distributed transition density through time to

346 São Paulo. Such a shift of transition density towards more recent times is slightly

347 noticeable for the capital cities of the Central-West and North regions (see

348 Supplementary Figure S1, for transitions by time from Porto Alegre).

349

350 **Discussion**

351       We reconstructed the phylogeographic history of HIV-1C in Brazil using a

352 comprehensive set of *pol* and *env* subtype C sequences from 22 different cities, of

353 which 10 were sampled for the first time. Using a new model-testing based

354    approach and by calculating several phylogeny-trait association measures using

355    BaTS (39) we could reject random mixing in both datasets. However, as seen in

356    MC scores, not all locations contributed equally to phylogenetic signal resulting

357    in a considerable degree of uncertainty in the phylogeographic inferences.

358    Nevertheless, after balancing the number of samples per location to mitigate the

359    confounding effects of sampling biases, we were able to identify support for two

360    epidemiological predictors of the viral spatial diffusion process.

361        Specifically, we found higher migration intensity from cities with larger

362    numbers of HIV-1C infected patients and higher HIV prevalence. Interestingly,

363    this is in agreement with a pattern of HIV-1C spread towards the north of Brazil,

364    where the prevalence of HIV is smaller and only few cases of HIV-1C infection are

365    found (4,6). An intriguing result illustrating the complexity of modeling human

366    mobility is that neither "geographical distance" nor "passenger air traffic"

367    predicted viral spread. The sample size of source and/or recipient locations, on

368    the other hand, were always included in the model (in isolation or together,

369    Table 2). Samples sizes are expected to predict the number of transitions to some

370    extent, and it was not our intention to formally demonstrate this. Rather, we

371    wanted to avoid that other predictors would be supported simply because of

372    correlation with sample sizes. In other, words we do not expect the support for

373    HIV prevalence and subtype C population size in the origin locations is artifact of

374    the potential correlation with sample size as this is already accommodated

375    explicitly in the GLM analysis.

376        To explore how sampling heterogeneity also impacts ancestral

377    reconstructions, we analyzed six random down-sampled datasets in parallel with

378    the complete *pol* and *env* datasets. This highlighted a substantial variability in the

379     root state estimates (Table 1) and confirms that the sampling scheme can indeed

380     have a profound effect on the inferred location state probabilities at the internal

381     nodes of the tree. The impact of sampling biases was most likely aggravated by

382     the relatively high degree of mixing observed in the *pol* and *env* datasets (see

383     Supplementary Table S2, for MC scores). The geographical partitioning is also an

384     important factor in discrete phylogeographic analyses because this determines

385     the level of spatial detail that can be recovered. Whereas previous studies

386     investigating the spread of HIV-1C in Brazil discretized locations according to

387     federal states or geopolitical regions (7-9), we opted for a higher-resolution

388     scheme and defined cities as the locations of interest. This allowed us to include

389     more precise predictors in the GLM model.

390     We were able to largely resolve sampling-bias related inconsistencies by

391     informing the phylogeographical reconstructions with relevant epidemiological

392     information. Our results consistently identify POA, and not the state of Paraná (7-

393     9), as the point of introduction. Several lines of evidence support this hypothesis.

394     The population in the metropolitan area of POA has about 4 million inhabitants,

395     the largest in the South region, and the AIDS incidence rate in POA and its

396     metropolitan area is the highest rate in Brazil (6). This suggests that the virus

397     found ideal circumstances for transmission and explains why the HIV-1C

398     prevalence in Paraná's capital Curitiba is much lower (~22%) compared to POA

399     (~40% and up to ~60% if the proportion of CRF31_BC - a local circulating form

400     with a small subtype B insertion in a subtype C backbone- is considered) (4, 44).

401     Differences in risk-group associations between the subtype B and C

402     epidemics in Brazil also seem to support our findings. Whereas in POA the

403     association between men-having-sex-with-men (MSM) and HIV-1B disappeared

404    in more recent sampling because of an expansion of HIV-1C in heterosexual

405    (HET) and MSM groups (45), compartmentalized epidemics are still observed in

406    other cities from the South region, including in Paraná, which could be explained

407    by a later introduction of HIV-1C (4,46-48).

408        Finally, a central role for POA is also reflected in the high support for

409    transitions from POA to all other locations and the reconstructed temporal

410    pattern of dissemination. After its introduction in the early 1970s, HIV-1C started

411    spreading to other cities in the beginning of the 1980s, first to nearby locations

412    and then to locations progressively further away. It is interesting to note that we

413    could recover a noticeably higher fraction of recent jumps from POA to São Paulo

414    when compared to transitions from POA to other South or South-East region

415    capital cities, which points towards a strong longstanding epidemiological link

416    between both cities.

417        Although our analysis provides support for POA as the central

418    dissemination point of HIV-1C in Brazil, some caution is required when analyzing

419    the number of transitions in star-like trees such as those typically found for HIV-

420    1. The absence of clear phylogenetic structure deeper in the trees also offers

421    little opportunity to capture clear spatial structuring and transitions beyond

422    those out of the location state at the root. In the current work, our sampling

423    strategy focused on broad geographic coverage rather than on a dense sampling,

424    and a small sample from a large and diverse population that has grown

425    exponentially through time, generally results in star-like topologies. Thus,

426    despite the support for a central role of POA, we can recover little detail on viral

427    spread beyond transitions out of this location.

428        In conclusion, we present a comprehensive reconstruction of the spatial

429    and temporal dynamics of HIV-1C in Brazil based on *pol* and, for the first time,

430    *env* sequence data, and included data from 10 newly sampled cities. By

431    augmenting the viral genetic information with epidemiological data, we revealed

432    a central role for POA city in the spread of HIV-1C in Brazil. In addition, we also

433    identified locations with high HIV prevalence and large subtype C population

434    sizes as key in the epidemic expansion towards the north of Brazil.

## Acknowledgements:

447      The authors have no conflicts of interest.

## References

1. **Faria NR, Rambaut A, Suchard M a., Baele G, Bedford T, Ward MJ, Tatem AJ, Sousa JD, Arinaminpathy N, Pepin J, Posada D, Peeters M, Pybus OG, Lemey P**. 2014. The early spread and epidemic ignition of HIV-1 in human populations. Science **346**:56–61.

2. **Tebit DM, Arts EJ**. 2011. Tracking a century of global expansion and evolution of HIV to drive understanding and to combat disease. Lancet Infect Dis **11**:45–56.

3. **Ariën KK, Vanham G, Arts EJ**. 2007. Is HIV-1 evolving to a less virulent form in humans? Nat Rev Microbiol **5**:141–151.

4. **Gräf T, Pinto AR**. 2013. The increasing prevalence of HIV-1 subtype C in Southern Brazil and its dispersion through the continent. Virology **435**:170–178.

5. **Inocencio LA, Pereira AA, Sucupira M, Fernandez J, Jorge CP, Souza DF, Fink HT, Diaz RS, Becker IM, Suffert TA, Arruda MB, Macedo O, Simão MB, Tanuri A**. 2009. Brazilian Network for HIV Drug Resistance Surveillance: a survey of individuals recently diagnosed with HIV. J Int AIDS Soc **12**:20.

6. **Brazilian Ministry of Health**. 2014. AIDS Epidemiological Bulletin July 2013-June 2014. Brasília, DF.

7. **Veras NMC, Gray RR, de Macedo Brigido LF, Rodrigues R, Salemi M**. 2011. High-resolution phylogenetics and phylogeography of human immunodeficiency virus type 1 subtype C epidemic in South America. J Gen Virol **92**:1698–1709.

8. **Delatorre E, Couto-Fernandez JC, Guimarães ML, Vaz Cardoso LP, de Alcantara KC, Martins de Araújo Stefani M, Romero H, Freire CCM, Iamarino A, de A Zanotto PM, Morgado MG, Bello G**. 2013. Tracing the origin and northward dissemination dynamics of HIV-1 subtype C in Brazil. PLoS One **8**:e74072.

473    9. **Bello G, Zanotto PMA, Iamarino A, Gräf T, Pinto AR, Couto-Fernandez JC,**

474    **Morgado MG**. 2012. Phylogeographic analysis of HIV-1 subtype C dissemination

475    in Southern Brazil. PLoS One **7**:e35649.

476    10.  **Brígido LFM, Ferreira JLP, Almeida VC, Rocha SQ, Ragazzo TG, Estevam**

477    **DL, Rodrigues R**. 2011. Southern Brazil HIV Type 1 C expansion into the state of

478    São Paulo, Brazil. AIDS Res Hum Retrovir **27**:339–344.

479    11.  **Cardoso LPV, Pereira GAS, Viegas AA, Schmaltz LEPR, Stefani MM de A**.

480    2010. HIV-1 primary and secondary antiretroviral drug resistance and genetic

481    diversity among pregnant women from central Brazil. J Med Virol **82**:351–357.

482    12.  **Carvalho BC, Cardoso LPV, Damasceno S, Stefani MM de A**. 2011.

483    Moderate prevalence of transmitted drug resistance and interiorization of HIV

484    type 1 subtype C in the inland north state of Tocantins, Brazil. AIDS Res Hum

485    Retrovir **27**:1081–1087.

486    13. **Ferreira AS, Cardoso LPV, Stefani MM de A**. 2011. Moderate prevalence of

487    transmitted drug resistance and high HIV-1 genetic diversity in patients from

488    Mato Grosso state, Central Western Brazil. J Med Virol **83**:1301–1307.

489    14. **da Silveira AA, Cardoso LPV, Francisco RBL, Stefani MM de A**. 2012. HIV

490    type 1 molecular epidemiology in pol and gp41 genes among naive patients from

491    Mato Grosso do Sul State, Central Western Brazil. AIDS Res Hum Retrovir

492    **28**:304–307.

493    15. **Rambaut A, Posada D, Crandall KA, Holmes EC**. 2004. The causes and

494    consequences of HIV evolution. Nat Rev Genet **5**:52–61.

495    16. **Drummond AJ, Pybus OG, Rambaut A, Forsberg R, Rodrigo AG**. 2003.

496    Measurably evolving populations. Trends Ecol Evol **18**:481–488.

497  17. **Rambaut A**. 2000. Estimating the rate of molecular evolution: incorporating

498  non-contemporaneous sequences into maximum likelihood phylogenies.

499  Bioinformatics **16**:395–399.

500  18. **Drummond AJ, Ho SYW, Phillips MJ, Rambaut A**. 2006. Relaxed

501  phylogenetics and dating with confidence. PLoS Biol **4**:699–710.

502  19. **Pybus OG, Rambaut A, Harvey PH**. 2000. An integrated framework for the

503  inference of viral population history from reconstructed genealogies. Genetics

504  **155**:1429–1437.

505  20. **Strimmer K, Pybus OG**. 2001. Exploring the demographic history of DNA

506  sequences using the generalized skyline plot. Mol Biol Evol **18**:2298–2305.

507  21. **Drummond AJ, Nicholls GK, Rodrigo AG, Solomon W**. 2002. Estimating

508  mutation parameters, population history and genealogy simultaneously from

509  temporally spaced sequence data. Genetics **161**:1307–1320.

510  22. **Lemey P, Rambaut A, Drummond AJ, Suchard MA**. 2009. Bayesian

511  phylogeography finds its roots. PLoS Comput Biol **5**:e1000520.

512  23. **Lemey P, Rambaut A, Welch JJ, Suchard MA**. 2010. Phylogeography takes a

513  relaxed random walk in continuous space and time. Mol Biol Evol **27**:1877–1885.

514  24. **Vaughan TG, Kuhnert D, Popinga A, Welch D, Drummond AJ**. 2014.

515  Efficient Bayesian inference under the structured coalescent. Bioinformatics

516  **30**:2272–2279.

517  25. **Pybus OG, Drummond AJ, Nakano T, Robertson BH, Rambaut A**. 2003.

518  The epidemiology and iatrogenic transmission of hepatitis C virus in Egypt: A

519  Bayesian coalescent approach. Mol Biol Evol **20**:381–387.

520    26. **Rambaut A, Pybus OG, Nelson MI, Viboud C, Taubenberger JK, Holmes**

521    **EC**. 2008. The genomic and epidemiological dynamics of human influenza A

522    virus. Nature **453**:615–619.

523    27. **Lemey P, Rambaut A, Bedford T, Faria N, Bielejec F, Baele G, Russell CA,**

524    **Smith DJ, Pybus OG, Brockmann D, Suchard MA**. 2014. Unifying viral genetics

525    and human transportation data to predict the global transmission dynamics of

526    human influenza H3N2. PLoS Pathog **10**:e1003932.

527    28. **Nelson MI, Viboud C, Vincent AL, Culhane MR, Detmer SE, Wentworth**

528    **DE, Rambaut A, Suchard MA, Holmes EC, Lemey P**. 2015. Global migration of

529    influenza A viruses in swine. Nat Commun **6**:6696.

530    29. **Librelotto CS, Gräf T, Simon D, Almeida SEM, Lunge VR**. 2015. HIV-1

531    epidemiology and circulating subtypes in the countryside of South Brazil. Rev

532    Soc Bras Med Trop **48**:249–257.

533    30. **de Oliveira T, Deforche K, Cassol S, Salminen M, Paraskevis D, Seebregts**

534    **C, Snoeck J, van Rensburg EJ, Wensing AMJ, van de Vijver DA, Boucher CA,**

535    **Camacho R, Vandamme A-M**. 2005. An automated genotyping system for

536    analysis of HIV-1 and other microbial sequences. Bioinformatics **21**:3797–3800.

537    31.  **Pond SLK, Posada D, Stawiski E, Chappey C, Poon AFY, Hughes G,**

538    **Fearnhill E, Gravenor MB, Brown AJL, Frost SDW**. 2009. An evolutionary

539    model-based algorithm for accurate phylogenetic breakpoint mapping and

540    subtype prediction in HIV-1. PLoS Comput Biol **5**:e1000581.

541    32. **Siepel AC, Halpern AL, Macken C, Korber BTM**. 1995. A computer program

542    designed to screen rapidly for HIV type 1 intersubtype recombinant sequences.

543    AIDS Res Hum Retrovir **11**:1413–1416.

544 33. **Lole KS, Bollinger RC, Paranjape RS, Gadkari D, Kulkarni SS, Novak NG,**

545 **Ingersoll R, Sheppard HW, Ray SC**. 1999. Full-length human immunodeficiency

546 virus type 1 genomes from subtype C-infected seroconverters in India, with

547 evidence of intersubtype recombination. J Virol **73**:152–160.

548 34. **Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K,**

549 **Madden TL**. 2009. BLAST+: architecture and applications. BMC Bioinformatics

550 **10**:421.

551 35. **Drummond AJ, Suchard MA, Xie D, Rambaut A**. 2012. Bayesian

552 phylogenetics with BEAUti and the BEAST 1.7. Mol Biol Evol **29**:1969–1973.

553 36. **Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO,**

554 **Huelsenbeck JP, Ronquist F, Swofford DL, Cummings MP, Rambaut A,**

555 **Suchard MA**. 2012. BEAGLE: An application programming interface and high-

556 performance computing library for statistical phylogenetics. Syst Biol **61**:170–3.

557 37. **Minin VN, Bloomquist EW, Suchard MA.** 2008. Smooth skyride through a

558 rough skyline: Bayesian coalescent-based inference of population dynamics. Mol

559 Biol Evol **25**:1459–1471.

560 38. **Shapiro B, Ho SYW, Drummond AJ, Suchard MA, Pybus OG, Rambaut A**.

561 2011. A bayesian phylogenetic method to estimate unknown sequence ages. Mol

562 Biol Evol **28**:879–887.

563 39. **Parker J, Rambaut A, Pybus OG**. 2008. Correlating viral phenotypes with

564 phylogeny: accounting for phylogenetic uncertainty. Infect Genet Evol **8**:239–46.

565 40. **Baele G, Lemey P, Bedford T, Rambaut A, Suchard MA, Alekseyenko A V**.

566 2012. Improving the accuracy of demographic and molecular clock model

567 comparison while accommodating phylogenetic uncertainty. Mol Biol Evol

568 **29**:2157–67.

569    41. **Baele G, Li WLS, Drummond AJ, Suchard MA, Lemey P**. 2013. Accurate

570    model selection of relaxed molecular clocks in bayesian phylogenetics. Mol Biol

571    Evol **30**:239–43.

572    42.    **Bielejec F, Rambaut A, Suchard MA, Lemey P**. 2011. SPREAD: spatial

573    phylogenetic reconstruction of evolutionary dynamics. Bioinformatics **27**:2910–

574    2.

575    43. **O'Brien JD, Minin VN, Suchard MA.** 2009. Learning to count: Robust

576    estimates for labeled distances between molecular sequences. Mol Biol Evol

577    **26**:801–814.

578    44.    **Passaes CPB, Bello G, Lorete RS, Matos Almeida SE, Junqueira DM,**

579    **Veloso VG, Morgado MG, Guimarães ML**. 2009. Genetic characterization of HIV-

580    1 BC recombinants and evolutionary history of the CRF31_BC in Southern Brazil.

581    Infect Genet Evol **9**:474–482.

582    45. **Almeida SEM, de Medeiros RM, Junqueira DM, Gräf T, Passaes CPB, Bello**

583    **G, Morgado MG, L Guimarães M**. 2012. Temporal dynamics of HIV-1 circulating

584    subtypes in distinct exposure categories in southern Brazil. Virol J **9**:306.

585    46.    **Raboni SM, Almeida SM De, Rotta I, Elisa C, Ribeiro L, Rosario D, Vidal**

586    **LR, Nogueira MB, Riedel M, Winhescki G, Ferreira KA, Ellis R**. 2010.

587    Molecular epidemiology of HIV-1 clades in Southern Brazil. Mem I Oswaldo Cruz

588    **105**:1044–1049.

589    47. **Gräf T, Passaes CPB, Ferreira LGE, Grisard EC, Morgado MG, Bello G,**

590    **Pinto AR**. 2011. HIV-1 genetic diversity and drug resistance among treatment

591    naïve patients from Southern Brazil: An association of HIV-1 subtypes with

592    exposure categories. J Clin Virol **51**:186–191.

593    48.  **Silveira J, Santos AF, Martínez AMB, Góes LR, Mendoza-Sassi R, Muniz**

594    **CP, Tupinambás U, Soares MA, Greco DB**. 2012. Heterosexual transmission of

595    human immunodeficiency virus type 1 subtype C in southern Brazil. J Clin Virol

596    **54**:36–41.

**Figure legends:**

**Figure 1. Administrative map of Brazil indicating the locations from where HIV-1C sequences were obtained.** Pie charts show the HIV-1C (black) percentage of infections relative to other HIV-1 strains (grey) in all cities with *pol* or *env* sequences included in this study. State name abbreviations are shown in bold. The inset shows an enlarged map with the sampling locations (black and red dots) in Santa Catarina and Rio Grande do Sul from which new sequence data were generated. Red dots: cities sampled for the first time. Black dots: sampling locations from where sequence data from other studies were also available. Brazilian regions are colored according to the legend. **Acronyms for states:** AM: Amazonas; GO: Goiás; MS: Mato Grosso do Sul; PR: Paraná; RJ: Rio de Janeiro; SC: Santa Catarina; SP: São Paulo; TO: Tocantins; RS: Rio Grande do Sul. **Acronyms for cities:** BLU: Blumenau; CHA: Chapecó; CPG: Campo Grande; CRA: Cruz Alta; CRI: Criciúma; CTB: Curitiba; CXS: Caxias do Sul; FLP: Florianópolis; GOI: Goiania; ITA: Itajaí; JOI: Joinville; LAJ: Lajeado; LGE: Lages; MAN: Manaus; PAL: Palmas; POA: Porto Alegre; RIG: Rio Grande; RJN: Rio de Janeiro; SPL: São Paulo; STI: Santiago; STL: Santana do Livramento; URU: Uruguaiana.

**Figure 2. Significant predictors of the Brazilian HIV-1C epidemic spread among sub-sampling datasets of *pol* and *env* genes.** Inclusion probabilities are represented as Bayes factors (BF) and a BF=3 threshold was used as positive indication of the predictor inclusion. The effect of each predictor, conditional to its inclusion, is represented by the posterior mean (black dot) and 95% HPD of the GLM coefficients in log scale.

**Table1. Modal root state and posterior probability estimates resulting from different discrete Bayesian phylogeographic analyses applied to different datasets.**

| Method | Subsampling | | | | | | |
|---|---|---|---|---|---|---|---|
| | Complete | RAND10A | RAND10B | RAND10C | RAND20A | RAND20B | RAND20C |
| *pol* | | | | | | | |
| *Symmetric-BSSVS* | FLP (1.00) | SPL (0.99) | CTB (0.99) | RJN (0.99) | POA (0.99) | RIG (0.99) | RIG (0.99) |
| *Symmetric* | FLP (0.99) | CRI (0.99) | FLP (0.97) | CRI (0.97) | FLP (0.99) | FLP (0.98) | ITA (0.99) |
| *Asymmetric-BSSVS* | FLP (1.00) | CRI (1.00) | FLP (0.99) | CTB (0.99) | ITA (1.00) | FLP (0.99) | POA (0.99) |
| *Asymmetric* | FLP (0.99) | RJN (0.99) | FLP (0.99) | CTB (0.99) | ITA (1.00) | FLP (0.99) | FLP (0.99) |
| *env* | | | | | | | |
| *Symmetric-BSSVS* | FLP (0.97) | FLP (0.99) | POA (0.99) | POA (0.99) | FLP (0.99) | FLP (0.99) | CXS (0.99) |
| *Symmetric* | FLP (0.97) | FLP (0.99) | POA (0.99) | FLP (0.96) | FLP (0.99) | FLP (0.99) | CRI (0.99) |
| *Asymmetric-BSSVS* | FLP (0.99) | CRI (0.99) | POA (0.99) | LAJ (0.99) | CRI (1.00) | POA (0.99) | FLP (1.00) |
| *Asymmetric* | FLP (0.99) | BLU (0.99) | FLP (0.99) | CRI (0.69) | FLP (0.99) | CRA (0.99) | FLP (1.00) |

**Acronyms for cities:** BLU: Blumenau; CRI: Criciúma; CTB: Curitiba; CXS: Caxias do Sul; FLP: Florianópolis; ITA: Itajaí; LAJ: Lajeado; POA: Porto Alegre; RIG: Rio Grande; RJN: Rio de Janeiro.

**Table 2. Bayes factor support for an explanatory role in the HIV-1C_BR**

**diffusion process for all tested predictors in all datasets**

| Predictor | Dataset | | | | | | |
|---|---|---|---|---|---|---|---|
| | Complete | RAND10A | RAND10B | RAND10C | RAND20A | RAND20B | RAND20C |
| *pol* | | | | | | | |
| Geographical Distance | 0.1 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 |
| Origin Sample Size | **Inf** | 0.5 | 0.4 | 0.4 | 0.4 | 0.4 | 0.4 |
| Destination Sample Size | **Inf** | **6583.3** | **5758.4** | **1674.5** | **Inf** | **Inf** | **Inf** |
| Origin HIV pop. size | 0.6 | 1.4 | 1.1 | 1.6 | 1.2 | 1.1 | 1.2 |
| Dest. HIV pop. size | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Origin HIV prevalence** | 0.2 | **6.1** | **10.8** | **16.1** | **17.3** | **13.1** | **7.9** |
| Dest. HIV prevalence | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Origin HIV-1C pop. size** | 0.3 | **14.4** | **22.4** | **9.6** | **18.1** | **23.4** | **38.8** |
| Dest. HIV-1C pop. size | 0.0 | 0.0 | 0.1 | 0.1 | 0.0 | 0.0 | 0.0 |
| Origin HIV-1C prevalence | 0.3 | **7.2** | 1.3 | **4.1** | 1.1 | 1.4 | 1.0 |
| Dest. HIV-1C prevalence | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| *env* | | | | | | | |
| Geographical Distance | 0.0 | 0.1 | 0.1 | 0.1 | 0.0 | 0.0 | 0.0 |
| Origin Sample Size | **Inf** | 0.3 | 0.4 | 0.9 | 0.3 | 0.3 | 0.4 |
| Destination Sample Size | **Inf** | **20.6** | **32.5** | **29.4** | **Inf** | **Inf** | **Inf** |
| Origin HIV pop. size | 0.3 | 0.8 | 1.1 | 2.6 | 0.8 | 0.9 | 0.9 |
| Dest. HIV pop. size | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Origin HIV prevalence** | 0.3 | **19.6** | **18.9** | **15.3** | **22.3** | **26.0** | **23.4** |
| Dest. HIV prevalence | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **Origin HIV-1C pop. size** | 0.3 | **13.8** | **14.0** | **15.5** | **13.0** | **11.3** | **12.9** |
| Dest. HIV-1C pop. size | 0.0 | 0.1 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Origin HIV-1C prevalence | 0.4 | 0.4 | 0.6 | 1.0 | 0.7 | 0.6 | 0.5 |
| Dest. HIV-1C prevalence | 0.0 | 0.6 | 0.6 | 0.6 | 0.0 | 0.0 | 0.0 |

Epidemiological predictors included in all Rand10 and Rand20 datasets are in bold, as well as BF

≥3; Dest. : destination.

Brazilian Regions
- North
- Central-West
- Northeast
- Southeast
- South

AM

MAN

TO

PAL

GO

GOI

MS

CPG

SP

SPL

PR

CTB

SC

FLP

RS

RIG

RJ

RJN

N

CHA

BLU

JOI

ITA

LGE

CRA

CRI

STI

CXS

LAJ

URU

POA

STL

## pol



| Predictors | Complete | Rand10 | Rand20 |
|---|---|---|---|
| | | BF | BF | BF |
| Origin HIV-1C pop. size | | | 22 | 39 |
| Origin HIV prev. | | | 10 | 8 |
| Destination sample size | Inf | >150 | Inf |
| Origin sample size | Inf | | |

**Effect (In coefficient)**

## env



| Predictors | Complete | Rand10 | Rand20 |
|---|---|---|---|
| | | BF | BF | BF |
| Origin HIV-1C pop. size | | | 14 | 13 |
| Origin HIV prev. | | | 19 | 13 |
| Destination sample size | Inf | 32 | Inf |
| Origin sample size | Inf | | |

**Effect (In coefficient)**