

1 **Main Manuscript for**  
2 **Characterizing Spatial Epidemiology in a Heterogeneous**  
3 **Transmission Landscape Using a Novel Spatial**  
4 **Transmission Count Statistic**

5 Leke Lyu<sup>a</sup>, Gabriella Veytsel<sup>a</sup>, Guppy Stott<sup>a</sup>, Spencer Fox<sup>b</sup>, Cody Dailey<sup>a</sup>, Lambodhar Damodaran<sup>c</sup>,  
6 Kayo Fujimoto<sup>d</sup>, Pamela Brown<sup>e</sup>, Roger Sealy<sup>e</sup>, Armand Brown<sup>e</sup>, Magdy Alabady<sup>f</sup>, Justin Bahl<sup>a\*</sup>

7  
8 a. Institute of Bioinformatics, Department of Infectious Diseases, Department of Epidemiology and  
9 Biostatistics, Center for Ecology of Infectious Diseases, University of Georgia, Athens, GA, USA

10 b. Institute of Bioinformatics, Department of Epidemiology and Biostatistics, University of Georgia,  
11 Athens, GA, USA

12 c. Department of Pathobiology, School of Veterinary Medicine, University of Pennsylvania,  
13 Philadelphia, PA, USA

14 d. Department of Health Promotion and Behavioral Sciences, The University of Texas Health  
15 Science Center at Houston, Houston, TX, USA

16 e. Division of Disease Prevention and Control, Houston Health Department, Houston, TX, USA

17 f. Georgia Genomics and Bioinformatics Center, University of Georgia, Athens, GA, USA

18

19 **Email:** [justin.bahl@uga.edu](mailto:justin.bahl@uga.edu)

20 **Author Contributions:** Leke Lyu and Justin Bahl conceptualized and designed research. Leke  
21 Lyu, Gabriella Veytsel, Guppy Stott, Spencer Fox, Cody Dailey, Lambodhar Damodaran and Kayo  
22 Fujimoto performed research. Pamela Brown, Roger Sealy, Armand Brown and Magdy Alabady  
23 contributed new data. Leke Lyu, Gabriella Veytsel and Guppy Stott Analyzed data. Leke Lyu,  
24 Gabriella Veytsel, Guppy Stott, Spencer Fox, Cody Dailey, Lambodhar Damodaran, Kayo Fujimoto

1

25 and Justin Bahl wrote and reviewed the paper. Justin Bahl acquired funding, supervised work, and  
26 coordinated communication among team members.

27 **Competing Interest Statement:** The authors declare that they have no conflict of interest.

28 **Classification:** Biological sciences / Ecology

29 **Keywords:** Viral evolution, Genomic epidemiology, Pandemic control

30 **This PDF file includes:**

31 Main Text

32 Figures 1 to 4

33

## 34 **Abstract**

### 35 **Background**

36 Viral genomes contain records of geographic movements and cross-scale transmission dynamics.  
37 However, the impact of regional heterogeneity, particularly among rural and urban centers, on viral  
38 spread and epidemic trajectory has been less explored due to limited data availability. Intensive  
39 and widespread efforts to collect and sequence SARS-CoV-2 viral samples have enabled the  
40 development of comparative genomic approaches to reconstruct spatial transmission history and  
41 understand viral transmission across different scales.

42

### 43 **Methods**

44 We proposed a novel spatial transmission count statistic that efficiently summarizes the geographic  
45 transmission patterns imprinted in viral phylogenies. Guided by a time-scaled tree with ancestral  
46 trait states, we identified spatial transmission linkages and categorize them as imports, local  
47 transmissions, and exports. These linkages were then summarized to represent the epidemic  
48 profile of the focal area.

49

### 50 **Results**

51 We demonstrated the utility of this approach for near real-time outbreak analysis using over 12,000  
52 full genomes and linked epidemiological data to investigate the spread of the SARS-CoV-2 in  
53 Texas. Our study showed (1) highly populated urban centers were the main sources of the epidemic  
54 in Texas; (2) the outbreaks in urban centers were connected to the global epidemic; and (3)

55 outbreaks in urban centers were locally maintained, while epidemics in rural areas were driven by  
56 repeated introductions.

57

## 58 **Conclusions**

59 In this study, we introduce the Source Sink Score, which allows us to determine whether a localized  
60 outbreak may be the source or sink to other regions, and the Local Import Score, which assesses  
61 whether the outbreak has transitioned to local transmission rather than being maintained by  
62 continued introductions. These epidemiological statistics provide actionable information for  
63 developing public health interventions tailored to the needs of affected areas.

64

## 65 **Plain Language Summary**

66 This study examined how COVID-19 spread through urban and rural areas in Texas by analyzing  
67 the virus's genomes from over 12,000 samples. Our goal was to understand how the virus travels  
68 and impacts different regions. Our findings revealed that densely populated urban centers were the  
69 primary sources of the virus in Texas. In contrast, outbreaks in rural areas were often fueled by  
70 new introductions of the virus from external sources. To conduct this analysis, we employed new  
71 computational methods that track where the virus originates and where it spreads. These methods  
72 provide detailed information crucial for public health officials, particularly in regions where the virus  
73 has more severe impacts or exhibits unique spread patterns.

74

## 75 **Introduction**

76 Genomic epidemiology is a field that utilizes pathogen genomes to study the spread of infectious  
77 diseases through populations <sup>1</sup>. This approach has become increasingly popular due to the  
78 decreasing cost of genomic sequencing combined with increasing computational power. During the  
79 COVID-19 pandemic, increased number of countries started generating genomic data to inform

80 public health responses<sup>2</sup>. The Global Initiative on Sharing All Influenza Data (GISAID)<sup>3</sup> expanded  
81 to accommodate these novel data and now maintains the world's largest database of SARS-CoV-  
82 2 sequences. As of December 2023, over 16 million sequences, sampled from over 200  
83 countries/regions, have been submitted and archived. Such a vast and diverse dataset enables  
84 researchers and public health officials to identify key mutations<sup>4,5</sup> and track the emergence of  
85 variants of interest (VOIs) or variants of concern (VOCs). Additionally, this wealth of genomic  
86 information creates opportunities to uncover the hidden characteristics of the local-scale outbreak,  
87 such as the spatial dispersal of transmission and the demographic characteristics contributing to  
88 transmission patterns. However, effectively handling the complexity of the SARS-CoV-2 genomic  
89 dataset requires addressing key challenges, such as establishing robust sampling frameworks to  
90 draw reliable conclusions and developing efficient computational algorithms/pipelines.

91

92 In genomic epidemiology, analyzing sampling biases and develop an appropriate sampling strategy  
93 are crucial steps<sup>6</sup>. Recent studies have shown that differences in epidemiology and sampling can  
94 impact our ability to identify genomic clusters<sup>7</sup>. For instance, decreased sampling fraction can lead  
95 to the identification of multiple, separate clusters. Sampling biases can also impact  
96 phylogeographic analyses. When investigating diffusion in discrete spaces, if a specific area is  
97 overrepresented in the dataset, it may lead to an overrepresentation of the same area at inferred  
98 internal nodes<sup>1</sup>. Similarly, when investigating diffusion in continuous space, extreme sampling bias  
99 might cause the posterior distribution to exclude the true origin location of the root<sup>8</sup>.

100

101 Viral transmission happens at different spatial scales, encompassing international pandemics,  
102 domestic dispersal, and local outbreaks such as those in jails, nursing homes, hospitals, or schools.  
103 By mapping how pathogens spread through space and time, evidence-based interventions can be  
104 better developed and applied across various scales<sup>9</sup>. The well-established software package,  
105 Bayesian Evolutionary Analysis Sampling Trees (BEAST)<sup>10</sup>, implements discrete<sup>11</sup> and continuous  
106<sup>12</sup> phylogeographic models. Previous studies have used the discrete model to identify the

107 transmission clusters of SARS-CoV-2 introduced in Europe <sup>13</sup>, United States <sup>14</sup>, Denmark <sup>15</sup> and  
108 England <sup>16</sup>. Additionally, the continuous model has been applied to elucidate the spatial expansion  
109 of SARS-CoV-2 in Belgium <sup>17</sup> and New York City <sup>18</sup>. Moreover, the BEAST module can  
110 accommodate individual travel history <sup>19</sup> to yield high-accuracy prediction regarding the location of  
111 ancestral nodes. Apart from Bayesian analysis, TreeTime <sup>20</sup> applies a maximum likelihood  
112 approach to infer the transitions between discrete characters. As a component of the Nextstrain <sup>21</sup>  
113 pipeline, this fast analysis enables real-time tracking of pathogens. With the rapid growth in SARS-  
114 CoV-2 genome data, we are now facing extensive phylogenies with thousands of tips. This raises  
115 the question: How can we translate the evolutionary changes of geographic traits from such  
116 expansive trees into clear epidemiological insights?

117

118 The transmission dynamics of SARS-CoV-2 are shaped by host immunity, host movement patterns,  
119 and other demographic characteristics <sup>22</sup>. For instance, in Chile, people aged under 40 in  
120 municipalities with the lowest socioeconomic status had an infection fatality rate 3.1 times higher  
121 than those with the highest socioeconomic status <sup>23</sup>. The severity of SARS-CoV-2 infection and the  
122 risk of mortality increased significantly with age <sup>24</sup>. Accordingly, identifying at-risk populations is  
123 crucial for determining the potential burden on public health. In the US, rural populations have been  
124 particularly vulnerable to COVID-19 complications <sup>25</sup>, experiencing higher incidences of disease  
125 and mortality <sup>26</sup>. This vulnerability is largely attributed to limited access to healthcare and social  
126 services <sup>27</sup>, as well as reduced access to and utilization of health information sources <sup>28</sup> compared  
127 to urban residents. Previous phylodynamic analyses have shown that frequent bi-directional  
128 transmission occurs between rural and urban communities <sup>29</sup>. However, few studies have  
129 investigated the differences in transmission patterns between these areas.

130

131 In this study, we developed a pipeline to understand local-scale epidemic trends. Our approach  
132 includes proportional genome sampling based on case counts <sup>14</sup>, followed by phylogeographic  
133 analysis using the Nextstrain framework <sup>21</sup>. Lastly, we summarize and compare transmission

134 patterns across subregions to identify viral sources and sinks. To demonstrate the utility of this  
135 method, we focused on the Delta wave in Texas, aiming to characterize viral diffusion within the  
136 state and compare epidemic trends between urban and rural areas.

137

## 138 **Methods**

### 139 **Surveillance and genetic dataset**

140 The United States Office of Management and Budget (OMB) defines Texas as having 25  
141 metropolitan areas (Table S1). Any population, housing, or territory not included in these  
142 metropolitan areas is classified as rural. The Rural-Urban Continuum Codes (RUCC) further  
143 categorize metropolitan areas based on population size. Dallas–Fort Worth, Houston, San Antonio,  
144 and Austin, all classified as RUCC-1<sup>30</sup>, represent the most populous metropolitan areas in Texas.

145

146 We obtained historical COVID-19 data for confirmed cases in Texas from the Texas Department of  
147 State Health Services (DSHS) website<sup>31</sup>. These weekly case counts, organized by county, were  
148 aggregated into metropolitan areas to inform our genome sampling strategy. Following the  
149 approach of Anderson F. Brito<sup>14</sup>, we developed R scripts, later consolidated into an R package  
150 called *Subsampler*. This package processes case count tables and genome metadata, enabling  
151 visual exploration of sampling heterogeneity and the implementation of proportional sampling  
152 schemes.

153

154 With support from the Houston Health Department (HHD), we accessed a large dataset of SARS-  
155 CoV-2 genomes sampled in Texas: 51,229 genomes with linked metadata. We focused on the  
156 Delta variant for our analysis, as its outbreak caused severe illness, spread rapidly before  
157 widespread immunity was established, and was intensely sampled at multiple scales<sup>32</sup>. Of the  
158 available genomes, 24,593 were identified as Delta variant, and 5,899 were subsampled  
159 proportionally to the case counts. Additionally, we sampled 6,386 Delta genomes from 49 countries

160 to provide global context and estimate viral migration to and from Texas. Our final dataset  
161 comprises 12,285 epidemiologically linked SARS-CoV-2 genomes.

162

### 163 **Phylogeographic analysis pipeline**

164 The pipeline comprises two major components: (1) Phylogenetic Reconstruction and (2)  
165 Characterization of Spatial Transmission Linkages.

166

167 Phylogenetic Reconstruction: This component utilizes the Nextstrain pipeline <sup>21</sup> to generate a time-  
168 labeled phylogeny with inferred ancestral trait states. Sequence alignment was conducted using  
169 Nextalign <sup>21</sup>, while the maximum likelihood tree construction was achieved with IQ-TREE <sup>33</sup>,  
170 applying a GTR substitution model. TreeTime <sup>20</sup> was employed to produce a time-scaled phylogeny  
171 and infer ancestral node states. The phylogeny was rooted using early samples from Wuhan  
172 (Wuhan-Hu-1/2019). Its temporal resolution was set based on an assumed nucleotide substitution  
173 rate of  $8 * 10^{-4}$  substitutions per site per year (default setting of Nextstrain build for SARS-CoV-2).  
174 Migration patterns between distinct geographic regions were inferred through time-reversible  
175 models <sup>20</sup>.

176

177 Characterization of Spatial Transmission Linkages: This component used custom scripts to identify  
178 spatial transmission linkages from the phylogeny and summarize epidemic trends in the focal  
179 region. The tree file was imported and read using the 'treeio' <sup>34</sup> package in R. The tree was then  
180 converted into a structured data frame for further analysis, facilitated by the 'tidytree' package <sup>34</sup>.  
181 Branches with durations exceeding 15 days were excluded, and the shorter branches in the  
182 phylogeny were designated as spatial transmission linkages. By analyzing trait states, we identified  
183 whether transmissions occurred within the focal area, involved imports from another location, or



184 resulted in exports to another area. The time series of spatial transmission counts, categorized by  
185 type, provides an overview of the epidemic trends in the focal region.

186

### 187 **Metrics that describe transmission pattern**

188 Different areas possess varying population sizes, levels of population mobility, and immunological  
189 characteristics, all of which can contribute to differences in the size and dynamics of the epidemic.

190 We introduced two metrics to quantitatively compare the characteristics of epidemics in different  
191 areas.

192

193 We define the Local Import Score to estimate the proportion of new cases due to introductions:

194 
$$Local\ Import\ Score = \frac{C_t(Import)}{C_t(Import) + C_t(Local\ Trans)}$$

195  $C_t(Import)$  represents the count of viral imports over a specific period  $t$  and  $C_t(Local\ Trans)$   
196 represents the count of local transmissions during the same period. The choice of the time window  
197 for calculation is contingent on the research objective. It can encompass the entire duration of the  
198 epidemic wave to assess cumulative effects, or it might focus on shorter intervals, such as

199 epidemiological weeks, for real-time surveillance. The Local Import Score ranges between 0 and  
200 1.

201

202 We introduce the Source Sink Score to identify whether a region acts primarily as a viral source or  
203 sink:

$$204 \quad \text{Source Sink Score} = \frac{C_t(\text{Export}) - C_t(\text{Import})}{C_t(\text{Export}) + C_t(\text{Import})}$$

205  $C_t(\text{Export})$  represents the count of exports over a specific period  $t$ . The Source-Sink Score ranges  
206 from -1 to 1. A score close to 1 indicates that the region primarily acts as a viral source, while a  
207 score near -1 suggests that the region mainly functions as a viral sink.

208

## 209 **Phylogenetic-based spatial network**

210 We constructed a weighted, undirected network to capture the viral flow between metropolitan  
211 areas in Texas. Each metropolitan area is represented as a node, and the edge carries weight  
212 corresponding to the spatial transmission counts. After establishing the network, we conducted the  
213 centrality analysis to rank the metropolitan areas based on their betweenness, closeness, and  
214 degree centrality. We processed the various network data objects using the 'igraph' package <sup>35</sup> in  
215 R. Visualizations were generated with the 'ggplot2' package <sup>36</sup>. We utilized the 'qgraph' package <sup>37</sup>  
216 to compute the centrality statistics of nodes.

217

## 218 **Sensitivity analysis of Source Sink Score and Local Import Scores**

219 To evaluate the robustness of the Source Sink Score and Local Import Score, we conducted a  
220 sensitivity analysis by generating nine additional genome datasets for Texas. These datasets were  
221 created using the same proportional sampling scheme as the original dataset. We then ran the  
222 same phylogeographic workflows on each dataset. By comparing the results across these

223 replicates, we assessed how uncertainties in sampling and phylogenetic inference affected the  
224 calculated scores.

225

## 226 **Results**

### 227 **Genome sampling bias and subsampling scheme adjustments**

228 With support from the Houston Health Department (HHD), we collected 24,593 Delta samples  
229 (B.1.167.2 and AY\*) with high-coverage complete genomes (>29,000 bp) and linked sampling site  
230 ZIP codes. Our genome database contained over a thousand distinct ZIP code records, which we  
231 translated into their affiliated metropolitan areas. We calculated the sampling ratio by dividing the  
232 number of available genomes by the number of reported cases to explore sampling biases.  
233 Significant heterogeneity in sampling ratios was observed across different metropolitan areas from  
234 Epi-Week 14 to Epi-Week 43 (Figure S1A). Victoria, Wichita Falls, and Bryan-College Station were  
235 identified as the top three under-sampled metropolitan areas, while Houston, San Angelo, and  
236 Abilene were the most over-sampled. To mitigate potential sampling biases, we applied a  
237 proportional sampling scheme (Figure S1B), thereby enhancing the accuracy of our  
238 phylogeographic analysis<sup>9,10</sup>. We adopted a consistent sampling ratio of 0.006 as a baseline for all  
239 regions. In regions that were under-sampled (sampling ratio below the baseline), all available  
240 genomes were retained. Conversely, over-sampled regions (with a sampling ratio exceeding the  
241 baseline) were down-sampled to match the baseline rate. As a result, we selected 5,899 Texas  
242 genomes, and the variance in sampling ratios across all metropolitan areas dropped substantially  
243 from 5.74e-05 to 7.56e-07.

244

### 245 **The transmission dynamics in Texas**

246 We conducted a comprehensive phylogeographic analysis of 12,048 SARS-CoV-2 Delta genomes  
247 sampled from March 27, 2021, to October 24, 2021, to investigate the timing of virus introduction  
248 into Texas and the dynamics of the resulting local transmission lineages. These genomes were

249 selected to ensure a roughly 1:1 ratio between Texas sequences (Table S2) and globally contextual  
250 sequences (Table S3). The Nextstrain <sup>21</sup> phylogenetic workflow was applied, in which a  
251 phylogenetic tree was estimated using IQ-TREE <sup>33</sup>, and a time-adjusted phylogeny was inferred  
252 with TreeTime <sup>20</sup>. The trait states of ancestral nodes were reconstructed as either 'Texas' or 'non-  
253 Texas' using the 'mugration' model implemented in TreeTime.

254

255 By considering the branches connecting each node to its parent as spatial transmission links, the  
256 location trait assigned to the nodes helps us categorize these connections into imports, local  
257 transmissions, and exports (Figures 1A, 1C). We defined a time series for these links as spatial  
258 transmission counts, providing a comprehensive summary of the epidemic's trends over time  
259 (Figures 1B, 1D). Given that the infectious period for SARS-CoV-2 typically ranges from day 2 to  
260 day 15 post-infection <sup>22</sup>, longer branches in the phylogeny likely indicate multiple transmission  
261 events. To reduce uncertainty, we excluded branches with durations exceeding 15 days, removing  
262 9,995 out of 22,991 branches. Our findings reveal that the Delta variant was first introduced into  
263 Texas on April 5, 2021, with a confidence interval from March 18, 2021, to April 5, 2021, preceding  
264 the first documented case in Houston in mid-April 2021 <sup>39</sup>. The Texas epidemic featured at least  
265 311 viral imports and 433 viral exports, linking statewide cases to the global pandemic. The  
266 outbreak in Texas was predominantly driven by local transmission, with 6,584 branches classified  
267 as local transmission.

268

## 269 **Characterizing spatial transmission heterogeneity**

270 To understand the spatial transmission of SARS-CoV-2 in Texas, we estimated ancestral location  
271 states on the phylogeny described above, incorporating 27 location traits: one contextual trait and  
272 26 subregions of Texas (25 metropolitan areas and one combined rural area) (Figure S2). We then  
273 constructed a network of metropolitan areas in Texas based on phylogeographic signals (Figure  
274 2A). The inferred network consisted of 25 nodes and 88 edges. Centrality analysis, detailed in  
275 Table S4, highlighted four pivotal nodes: Dallas–Fort Worth, Houston, San Antonio, and Austin.

276 These subregions were consistently identified as key hubs based on degree, betweenness, and  
277 connectedness <sup>40</sup>. Notably, all four of these metropolitan areas are classified as RUCC-1,  
278 suggesting populated urban centers played a crucial role in the viral spread across Texas.

279

## 280 **Community source-sink dynamics**

281 We introduced the Source Sink Score to classify populations as either viral sources or sinks. This  
282 score ranges from -1 to 1, with a score near 1 indicating a population is predominantly a viral  
283 source—where the number of exports greatly exceeds imports—and a score near -1 indicating a  
284 population is primarily a viral sink, where imports dominate over exports.

285

286 Subregions of Texas were categorized as sources or sinks based on their cumulative Source Sink  
287 Score, with the full list provided in Table S5. Our analysis showed that the RUCC-1 group, which  
288 represents densely populated urban centers, had the highest Source Sink Scores, emphasizing its  
289 role as a major source during the outbreak in Texas (Figure 2B). Within the RUCC-1 group, Dallas-  
290 Fort Worth had the highest score at 0.092, followed by Houston (0.063), San Antonio (0.000), and  
291 Austin (-0.444). In contrast, rural areas, with a score of -0.717, primarily acted as viral sinks.

292

## 293 **Epidemic trends in populated urban centers compared to rural areas**

294 We introduced the Local Import Score to estimate the proportion of new cases due to introductions.  
295 This score ranges from 0 to 1, with values closer to 1 indicating that the outbreak is primarily driven  
296 by external introductions, and values closer to 0 suggesting that local transmission is well-  
297 sustained. Identifying when most new cases are locally acquired is crucial for informing public  
298 health resource allocation, contact tracing efforts, and control strategies during emergency  
299 situations.

300

301 Using Houston as a representative city, we compared epidemic trends in densely populated urban  
302 centers to those in rural areas (Figure 3). Epidemic trends for other subregions are shown in Figures

303 S3–S26. The accumulated Local Import Score for Houston during the entire Delta wave was 0.176,  
304 indicating that the outbreak was largely sustained by local transmission. In contrast, rural areas  
305 had an accumulated Local Import Score of 0.558, suggesting that the epidemic there was primarily  
306 driven by external introductions. Our results suggest that while an outbreak may initially rely on  
307 external introductions, once the epidemic becomes locally sustained, the region can evolve into a  
308 primary source of pathogen spread to other areas (Figures 3C and 3D).

309

310 We also analyzed viral flow between global contexts and urban centers (e.g., Houston) (Figure 4A),  
311 as well as between global contexts and rural areas (Figure 4B). Introductions from outside Texas  
312 accounted for 60% of all imports to Houston, while 25% of all exports from Houston were to  
313 locations outside Texas. By comparison, introductions from non-Texas sources accounted for 26%  
314 of all imports to rural areas, and 3% of rural exports were to locations outside Texas. These findings  
315 suggest that Houston, as a highly connected and large urban center, served as an important hub  
316 linking the outbreak across Texas to the broader global pandemic.

317

### 318 **Assessing the sensitivity of the new metrics**

319 Despite the uncertainties inherent in sampling and phylogenetic reconstruction, our previous  
320 conclusions remained consistent across replicates. All 10 replicates supported RUCC-1 regions as  
321 the predominant viral sources, as these regions consistently showed the highest Source Sink  
322 Scores (Figure 4). Houston and Dallas–Fort Worth displayed the most robust results, as reflected  
323 by their narrow score ranges. Specifically, the Source Sink Score for Dallas–Fort Worth ranged  
324 from 0.049 to 0.151, while Houston's score ranged from 0.063 to 0.190. The Local Import Score for  
325 Dallas–Fort Worth ranged from 0.142 to 0.169, while Houston's score ranged from 0.152 to 0.178.

326 A detailed record of the sensitivity analysis conducted across different metrics is provided in Table  
327 S6.

328

## 329 **Discussion**

330 In this study, we introduced a novel spatial transmission count statistic, which characterizes the  
331 weekly counts of local spread, viral inflow, and outflow, illustrating transmission trends over time.

332 The Source Sink Score and Local Import Score are heuristic metrics that allow for quantitative  
333 comparison of epidemic trends between regions. The Source Sink Score measures net viral

334 exports, weighted by the outbreak size, while the Local Import Score compares the significance of  
335 external introductions versus local transmission in shaping the epidemic. We investigated the

336 geographic diffusion pattern of SARS-CoV-2 in Texas to demonstrate the utility of this novel  
337 phylogeographic approach. At the state level, we characterized the timing and size of viral imports.

338 Within the state of Texas, we reconstructed regional dissemination and contrasted the epidemic  
339 trends between urban centers and rural areas.

340

341 The size of our genomic data offers unprecedented opportunities for high-resolution investigations  
342 of spatial transmission history. Our analysis revealed that cryptic transmissions began as early as

343 late March, 2 to 3 weeks before the identification of the first Delta case in Houston <sup>39</sup>. Additionally,  
344 we identified at least 311 imports and 433 exports, highlighting Texas's intensive connection to the

345 global pandemic. Our results indicated that the Delta variant invaded Texas through multiple  
346 introductions. These independent imports subsequently formed massive local transmission clusters

347 in Texas. This pattern aligns with observations from Connecticut's initial COVID-19 wave <sup>41</sup>, the  
348 UK's first wave <sup>42</sup>, the emergence of B.1.1.7 variant across the United States <sup>14</sup>, and the presence

349 of Omicron BA.1 in England <sup>16</sup>.

350

351 The spatial transmission count statistic represents the time-series of categorized transmission  
352 linkages related to the focal regions. Informed by the annotated viral phylogeny, it summarizes the

353 trends of local spread and viral flow at a minimal computational cost. Adopting a simplified model,  
354 we assume that transmission events take place along all the branches of the viral phylogeny.  
355 However, phylogenetic trees are not equivalent to transmission trees; they do not directly reveal  
356 who infected whom <sup>43,44</sup>. As a result, our model may introduce bias in the estimation of local  
357 transmission counts. Despite this limitation, it provides valuable insights into local-scale  
358 transmission and epidemic trajectories that can inform control efforts. The efficiency of this statistic  
359 enables real-time surveillance of tens of thousands of viral genomes, which is crucial for addressing  
360 the challenges posed by the current pandemic or potential future outbreaks.

361

362 The role of a population as a source or sink evolves dynamically as the outbreak progresses and  
363 host immunity develops. Therefore, the Source Sink Score should be interpreted as a comparative  
364 measure, emphasizing relative differences between regions rather than absolute values. In Texas,  
365 populated urban centers functioned as the primary viral sources during the outbreak. Among all  
366 subregions, the RUCC-1 group had the highest Source Sink Scores, with Dallas-Fort Worth had  
367 the score at 0.092, followed by Houston (0.063), San Antonio (0.000), and Austin (-0.444). The  
368 significant role of these urban centers in spreading the epidemic can be linked to their key locations  
369 in road and air travel networks. Houston, Dallas-Fort Worth, and San Antonio, connected by  
370 Interstates 10, 45, and 35, form the vertices of the Texas Triangle <sup>45</sup>, one of 11 megaregions in the  
371 US and home to the majority of the Texas's population. This complex connectivity, along with the  
372 presence of major airports such as George Bush Intercontinental Airport in Houston (a United  
373 Airlines hub), Dallas-Fort Worth International Airport (American Airlines' largest primary hub and  
374 headquarters), and San Antonio International Airport (a Southwest Airlines hub), highlights their  
375 pivotal role in airway travel. Our analysis underscored the crucial role of urban centers in driving  
376 the outbreak. This insight provides valuable information that can guide public health decision-



377 making. Increased control efforts in highly connected urban centers may have a disproportionate  
378 impact on connected rural areas <sup>46</sup>.

379

380 Rural areas exhibit a lower level of viral flow in relation to global contexts, with epidemics in these  
381 regions predominantly relying on external introductions, thus establishing them as viral sinks.  
382 Notably, urban centers and rural areas demonstrate distinct transmission patterns. It is important  
383 to note that our analysis assumes that virus transmission in each region is influenced only by  
384 population size and density, without accounting for the effects of community behavior and beliefs,  
385 healthcare disparities, environmental factors, and other influences on viral transmission. Future  
386 studies addressing these aspects will provide more comprehensive insights into the underlying  
387 drivers of transmission.

388

389 Despite uncertainties in sampling and phylogenetic reconstruction, all replicates from the sensitivity  
390 analysis supported RUCC-1 regions as the predominant viral sources. The robustness of both the  
391 Source Sink Score and the Local Import Score varied across regions. Houston and Dallas–Fort  
392 Worth exhibited more stable results, with narrower score ranges, likely due to the larger volume of  
393 data available (>1500 genomes). In contrast, regions such as Amarillo, Odessa, and San Angelo  
394 had fewer genomes (<50 genomes), leading to broader score ranges and making interpretation  
395 less reliable. We believe that data availability and volume significantly impact the robustness of  
396 these metrics. Therefore, future users must carefully inspect data disparities and be cautious when  
397 interpreting results from regions with limited genome data.

398

399 Former Bayesian phylodynamic analyses, such as those conducted in Washington State <sup>47,48</sup>,  
400 investigated the role of viral introductions in community spread. These studies use effective  
401 population sizes estimated from approximate structured coalescent models to determine the  
402 percentage of new cases resulting from introductions. Inspired by these studies, we propose  
403 integrating the Source Sink Score and Local Import Score into a Bayesian phylodynamic framework

404 as future direction. This integration would allow us to calculate Bayesian Credible Intervals for these  
405 scores, providing a reliable measure of their uncertainty. This approach is particularly valuable  
406 when testing whether the Source Sink Score in one region, such as region A, is significantly higher  
407 than in another, such as region B, thereby facilitating robust regional comparisons.

408

## 409 **Data availability**

410 The GISAID accession IDs of the genomes used in this study are provided on our GitHub repository  
411 (<https://github.com/leke-lyu/transmissionCount>). Additional data obtained during the study is  
412 available from the corresponding author upon reasonable request.

413

## 414 **Code availability**

415 The R package *Subsamplerr*, which enables visual exploration of sampling heterogeneity and the  
416 implementation of proportional sampling schemes, is publicly available at [https://github.com/leke-](https://github.com/leke-lyu/subsamplerr)  
417 [lyu/subsamplerr](https://github.com/leke-lyu/subsamplerr). For the pipeline setup and configurations used in the Nextstrain build, including  
418 Snakemake profiles, visit our GitHub repository at [https://github.com/leke-](https://github.com/leke-lyu/surveillanceInTexas)  
419 [lyu/surveillanceInTexas](https://github.com/leke-lyu/surveillanceInTexas). All scripts used to generate the results in the Texas case study are publicly  
420 available at <https://github.com/leke-lyu/transmissionCount>.

421

## 422 **Acknowledgments**

423 This work has been funded in part from the National Institute of Allergy and Infectious Diseases, a  
424 component of the NIH, Department of Health and Human Services, under contract no.  
425 75N93021C00018 (NIAID Centers of Excellence for Influenza Research and Response, CEIRR)  
426 and Centers for Disease Control and Prevention, Department of Health and Human Services, under  
427 contracts 75D30121C10133 and NU50CK000626. We acknowledge the GISAID contributors

428 (acknowledgment table of genomes used is provided on our GitHub repository) for sharing genomic  
429 data.

430

## 431 **References**

- 432 1. Hill, V., Ruis, C., Bajaj, S., Pybus, O. G. & Kraemer, M. U. G. Progress and challenges in  
433 virus genomic epidemiology. *Trends in Parasitology* **37**, 1038–1049 (2021).
- 434 2. Genomic sequencing of SARS-CoV-2: a guide to implementation for maximum impact on  
435 public health. <https://www.who.int/publications-detail-redirect/9789240018440>.
- 436 3. Elbe, S. & Buckland-Merrett, G. Data, disease and diplomacy: GISAID's innovative  
437 contribution to global health. *Glob Chall* **1**, 33–46 (2017).
- 438 4. Hodcroft, E. B. *et al.* Want to track pandemic variants faster? Fix the bioinformatics  
439 bottleneck. *Nature* **591**, 30–33 (2021).
- 440 5. Grubaugh, N. D. *et al.* Tracking virus outbreaks in the twenty-first century. *Nat Microbiol* **4**,  
441 10–19 (2019).
- 442 6. Inward, R. P. D., Parag, K. V. & Faria, N. R. Using multiple sampling strategies to estimate  
443 SARS-CoV-2 epidemiological parameters from genomic sequencing data. *Nat Commun* **13**,  
444 5587 (2022).
- 445 7. Sobkowiak, B. *et al.* The utility of SARS-CoV-2 genomic data for informative clustering under  
446 different epidemiological scenarios and sampling. *Infection, Genetics and Evolution* **113**,  
447 105484 (2023).
- 448 8. Kalkauskas, A. *et al.* Sampling bias and model choice in continuous phylogeography: Getting  
449 lost on a random walk. *PLOS Computational Biology* **17**, e1008561 (2021).
- 450 9. Attwood, S. W., Hill, S. C., Aanensen, D. M., Connor, T. R. & Pybus, O. G. Phylogenetic and  
451 phylodynamic approaches to understanding and combating the early SARS-CoV-2  
452 pandemic. *Nat Rev Genet* **23**, 547–562 (2022).
- 453 10. Suchard, M. A. *et al.* Bayesian phylogenetic and phylodynamic data integration using BEAST  
454 1.10. *Virus Evol* **4**, vey016 (2018).

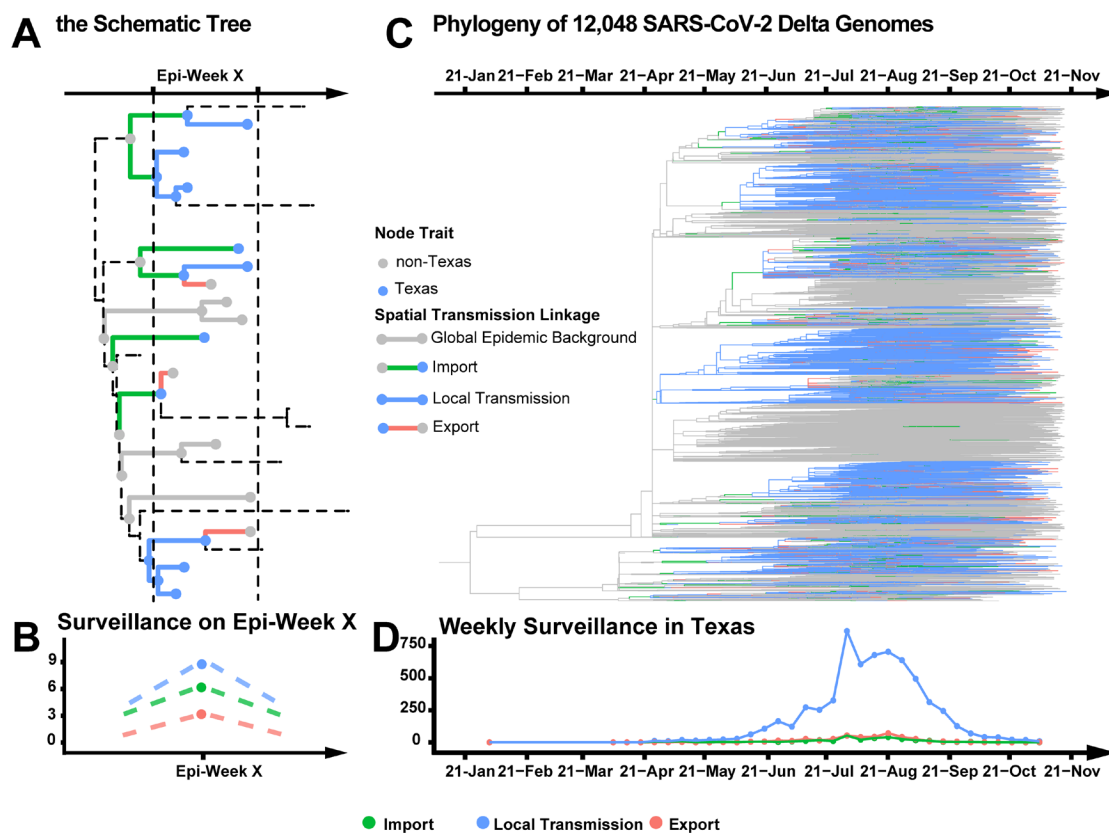
- 455 11. Lemey, P., Rambaut, A., Drummond, A. J. & Suchard, M. A. Bayesian Phylogeography Finds  
456 Its Roots. *PLOS Computational Biology* **5**, e1000520 (2009).
- 457 12. Lemey, P., Rambaut, A., Welch, J. J. & Suchard, M. A. Phylogeography Takes a Relaxed  
458 Random Walk in Continuous Space and Time. *Mol Biol Evol* **27**, 1877–1885 (2010).
- 459 13. Lemey, P. *et al.* Untangling introductions and persistence in COVID-19 resurgence in  
460 Europe. *Nature* **595**, 713–717 (2021).
- 461 14. Alpert, T. *et al.* Early introductions and transmission of SARS-CoV-2 variant B.1.1.7 in the  
462 United States. *Cell* **184**, 2595-2604.e13 (2021).
- 463 15. Michaelsen, T. Y. *et al.* Introduction and transmission of SARS-CoV-2 lineage B.1.1.7, Alpha  
464 variant, in Denmark. *Genome Medicine* **14**, 47 (2022).
- 465 16. Tsui, J. L.-H. *et al.* Genomic assessment of invasion dynamics of SARS-CoV-2 Omicron  
466 BA.1. *Science* **381**, 336–343 (2023).
- 467 17. Dellicour, S. *et al.* A Phylodynamic Workflow to Rapidly Gain Insights into the Dispersal  
468 History and Dynamics of SARS-CoV-2 Lineages. *Molecular Biology and Evolution* **38**, 1608–  
469 1613 (2021).
- 470 18. Dellicour, S. *et al.* Variant-specific introduction and dispersal dynamics of SARS-CoV-2 in  
471 New York City – from Alpha to Omicron. *PLOS Pathogens* **19**, e1011348 (2023).
- 472 19. Lemey, P. *et al.* Accommodating individual travel history and unsampled diversity in Bayesian  
473 phylogeographic inference of SARS-CoV-2. *Nat Commun* **11**, 5110 (2020).
- 474 20. Sagulenko, P., Puller, V. & Neher, R. A. TreeTime: Maximum-likelihood phylodynamic  
475 analysis. *Virus Evolution* **4**, vex042 (2018).
- 476 21. Hadfield, J. *et al.* Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* **34**,  
477 4121–4123 (2018).
- 478 22. Markov, P. V. *et al.* The evolution of SARS-CoV-2. *Nat Rev Microbiol* **21**, 361–379 (2023).
- 479 23. Mena, G. E. *et al.* Socioeconomic status determines COVID-19 incidence and related  
480 mortality in Santiago, Chile. *Science* **372**, eabg5298 (2021).

- 481 24. O’Driscoll, M. *et al.* Age-specific mortality and immunity patterns of SARS-CoV-2. *Nature*  
482 **590**, 140–145 (2021).
- 483 25. USDA ERS - Rural Residents Appear to be More Vulnerable to Serious Infection or Death  
484 From Coronavirus COVID-19. [https://www.ers.usda.gov/amber-waves/2021/february/rural-](https://www.ers.usda.gov/amber-waves/2021/february/rural-residents-appear-to-be-more-vulnerable-to-serious-infection-or-death-from-coronavirus-covid-19/)  
485 [residents-appear-to-be-more-vulnerable-to-serious-infection-or-death-from-coronavirus-](https://www.ers.usda.gov/amber-waves/2021/february/rural-residents-appear-to-be-more-vulnerable-to-serious-infection-or-death-from-coronavirus-covid-19/)  
486 [covid-19/](https://www.ers.usda.gov/amber-waves/2021/february/rural-residents-appear-to-be-more-vulnerable-to-serious-infection-or-death-from-coronavirus-covid-19/).
- 487 26. Cuadros, D. F., Branscum, A. J., Mukandavire, Z., Miller, F. D. & MacKinnon, N. Dynamics of  
488 the COVID-19 epidemic in urban and rural areas in the United States. *Ann Epidemiol* **59**, 16–  
489 20 (2021).
- 490 27. Mueller, J. T. *et al.* Impacts of the COVID-19 pandemic on rural America. *Proc Natl Acad Sci*  
491 *U S A* **118**, 2019378118 (2021).
- 492 28. Chen, X. *et al.* Differences in Rural and Urban Health Information Access and Use. *J Rural*  
493 *Health* **35**, 405–417 (2019).
- 494 29. Tang, C. Y. *et al.* Rural populations facilitated early SARS-CoV-2 evolution and transmission  
495 in Missouri, USA. *npj Viruses* **1**, 1–11 (2023).
- 496 30. USDA ERS - Rural-Urban Continuum Codes. [https://www.ers.usda.gov/data-products/rural-](https://www.ers.usda.gov/data-products/rural-urban-continuum-codes/)  
497 [urban-continuum-codes/](https://www.ers.usda.gov/data-products/rural-urban-continuum-codes/).
- 498 31. COVID-19 (Coronavirus Disease 2019) | Texas DSHS. [https://www.dshs.texas.gov/covid-19-](https://www.dshs.texas.gov/covid-19-home)  
499 [home](https://www.dshs.texas.gov/covid-19-home).
- 500 32. den Hartog, G. *et al.* Assessment of hybrid population immunity to SARS-CoV-2 following  
501 breakthrough infections of distinct SARS-CoV-2 variants by the detection of antibodies to  
502 nucleoprotein. *Sci Rep* **13**, 18394 (2023).
- 503 33. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective  
504 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology and*  
505 *Evolution* **32**, 268–274 (2015).
- 506 34. Yu, G. *Data Integration, Manipulation and Visualization of Phylogenetic Trees*. (CRC Press,  
507 2022).

- 508 35. Csardi, G. & Nepusz, T. The Igraph Software Package for Complex Network Research.  
509 *InterJournal Complex Systems*, 1695 (2005).
- 510 36. Valero-Mora, P. M. ggplot2: Elegant Graphics for Data Analysis. *Journal of Statistical*  
511 *Software* **35**, 1–3 (2010).
- 512 37. Epskamp, S., Cramer, A. O. J., Waldorp, L. J., Schmittmann, V. D. & Borsboom, D. qgraph:  
513 Network Visualizations of Relationships in Psychometric Data. *Journal of Statistical Software*  
514 **48**, 1–18 (2012).
- 515 38. Frost, S. D. W. *et al.* Eight challenges in phylodynamic inference. *Epidemics* **10**, 88–92  
516 (2015).
- 517 39. Christensen, P. A. *et al.* Delta Variants of SARS-CoV-2 Cause Significantly Increased  
518 Vaccine Breakthrough COVID-19 Cases in Houston, Texas. *Am J Pathol* **192**, 320–331  
519 (2022).
- 520 40. Freeman, L. C. Centrality in social networks conceptual clarification. *Social Networks* **1**, 215–  
521 239 (1978).
- 522 41. Fauver, J. R. *et al.* Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the  
523 United States. *Cell* **181**, 990-996.e5 (2020).
- 524 42. Establishment and lineage dynamics of the SARS-CoV-2 epidemic in the UK | Science.  
525 <https://www.science.org/doi/10.1126/science.abf2946>.
- 526 43. Hall, M. D. & Colijn, C. Transmission Trees on a Known Pathogen Phylogeny: Enumeration  
527 and Sampling. *Mol Biol Evol* **36**, 1333–1343 (2019).
- 528 44. Didelot, X., Fraser, C., Gardy, J. & Colijn, C. Genomic infectious disease epidemiology in  
529 partially sampled and ongoing outbreaks. *Mol Biol Evol* msw075 (2017)  
530 doi:10.1093/molbev/msw275.
- 531 45. Hagler, Y. Defining U.S. Megaregions.
- 532 46. Polo, G., Soler-Tovar, D., Villamil Jimenez, L. C., Benavides-Ortiz, E. & Mera Acosta, C.  
533 SARS-CoV-2 transmission dynamics in the urban-rural interface. *Public Health* **206**, 1–4  
534 (2022).

- 535 47. Paredes, M. I. *et al.* Local-scale phylodynamics reveal differential community impact of  
536 SARS-CoV-2 in a metropolitan US county. *PLOS Pathogens* **20**, e1012117 (2024).
- 537 48. Müller, N. F. *et al.* Viral genomes reveal patterns of the SARS-CoV-2 outbreak in Washington  
538 State. *Science Translational Medicine* **13**, eabf0202 (2021).
- 539
- 540
- 541

542 **Figures and Tables**



543

544 **Figure 1. The spatial transmission count statistic investigates the transmission dynamic. A.**

545 Conceptual figure showing that transmissions can be classified into three categories: import, local

546 transmission, and export. **B.** The schematic tree depicts a total of 18 spatial transmission linkages

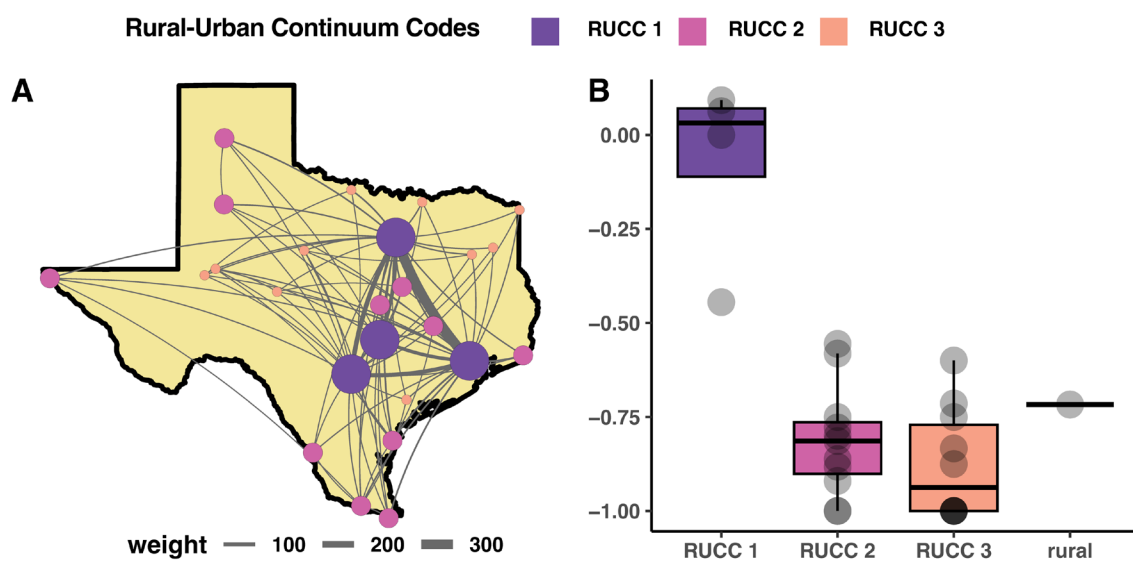
547 in Epi-Week X: 6 imports, 9 local transmissions, and 3 exports. **C.** In the time-adjusted phylogeny,

548 branches are colored based on the categories of the corresponding spatial transmission linkages.

549 **D.** The time series of spatial transmission counts summarizes the epidemic trend in Texas.

550

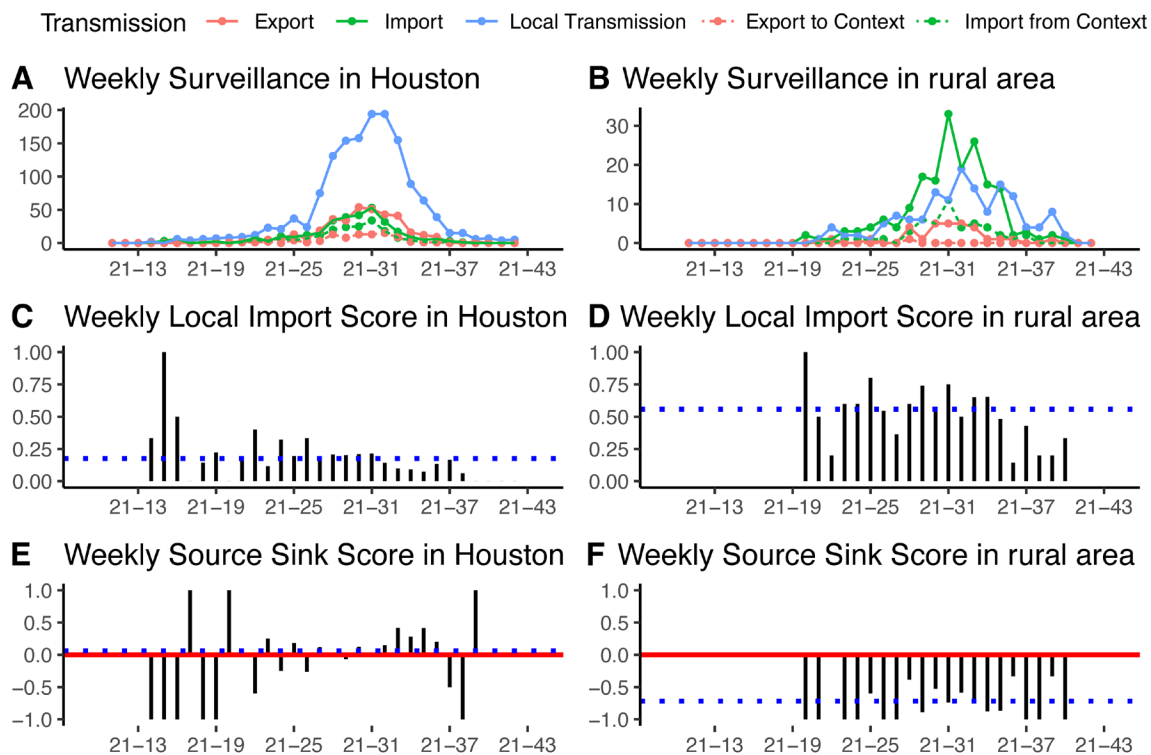




551

552 **Figure 2. Characterizing Spatial Transmission Heterogeneity.** Subregions across Texas are  
553 categorized by their Rural-Urban Continuum Codes (RUCC). RUCC-1 includes metropolitan areas  
554 with over 1 million residents, RUCC-2 includes areas with populations between 250,000 and 1  
555 million, and RUCC-3 represents areas with fewer than 250,000 residents. The four major urban  
556 centers—Dallas-Fort Worth, Houston, San Antonio, and Austin—are classified as RUCC-1. **A.**  
557 Phylogeographic network of Texas metropolitan areas. In this network, each node represents a  
558 metropolitan area, and the width of the edges is proportional to the spatial transmission counts. **B.**  
559 The Source Sink Score identifies key source hubs of SARS-CoV-2 spread in Texas. Dots in the  
560 box plot represent subregions of Texas.

561



562

563 **Figure 3. Epidemic trends on the Delta outbreak in populated urban centers vs the rural**

564 **areas. A.** The epidemic trend of Houston. **B.** The epidemic trend of the rural areas. The top of the

565 panel shows the time series of spatial transmission counts by week. The dashed pink line

566 represents exports from the analyzed regions to non-Texas. The dashed green line represents

567 imports from non-Texas into the analyzed regions. **C.** The trend of Local Import Score in Houston.

568 **D.** The trend of Local Import Score in rural areas. The black bars in the middle of the panel depict

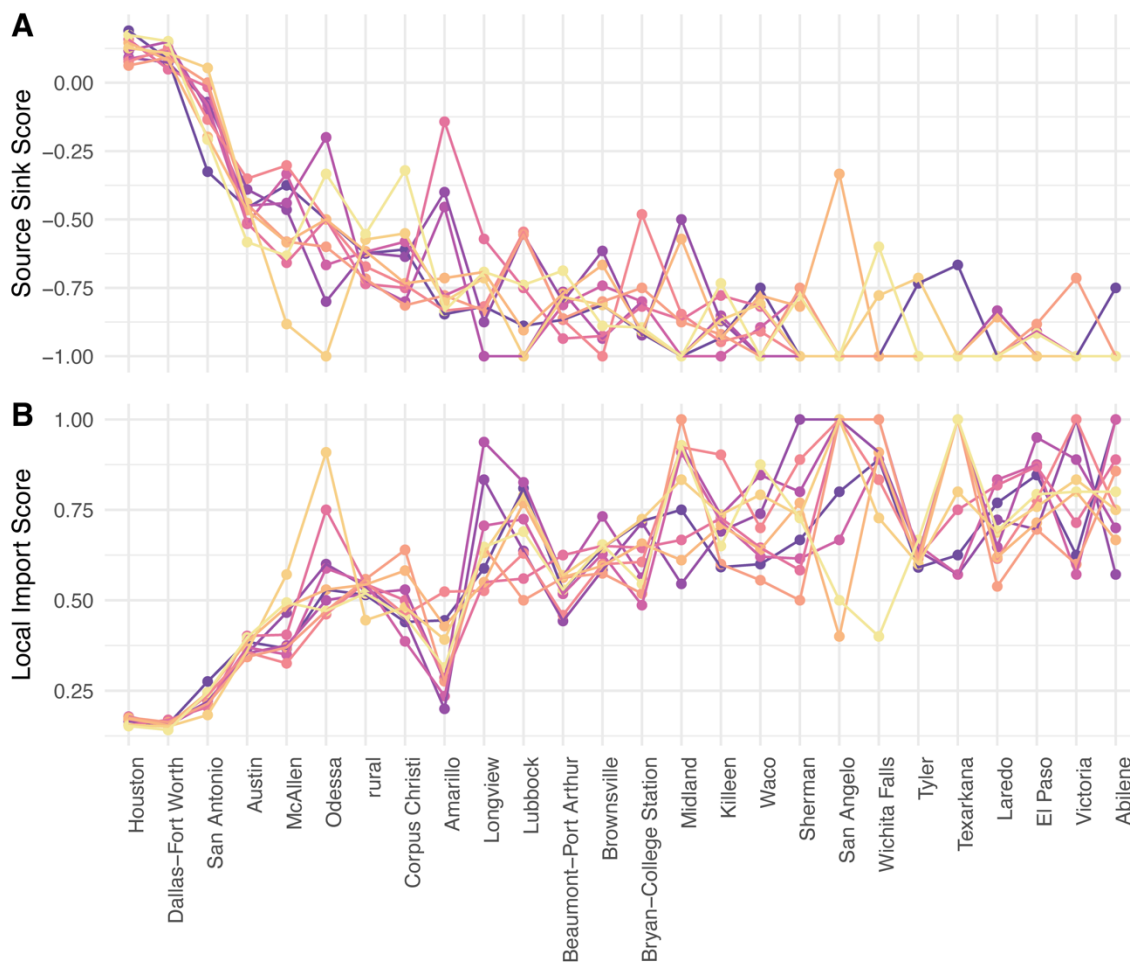
569 the weekly dynamics of Local Import Score. The dashed blue line indicates the accumulated Local

570 Import Score. **E.** The trend of Source Sink Score in Houston. **F.** The trend of Source Sink Score in

571 rural areas. The solid red line represents the benchmark of 0, indicating a balance between imports

572 and exports. The dashed blue line marks the accumulated Source Sink Score.

573



574

575 **Figure 4. Sensitivity Analysis of 10 Replicates. A.** Source Sink Score for subregions in Texas.

576 **B.** Local Import Score for subregions in Texas. Both plots share the same x-axis, where regions

577 are ranked from highest to lowest mean Source Sink Score. Each colored line connects the

578 statistics estimated from the same replicate.