

1 **What do we gain when tolerating loss? The information bottleneck, lossy compression, and**  
2 **detecting horizontal gene transfer**

3  
4 Apurva Narechania<sup>1\*</sup>, Rob DeSalle<sup>1</sup>, Barun Mathema<sup>2</sup>, Barry Kreiswirth, and Paul J. Planet<sup>1,4,5\*</sup>

5  
6 <sup>1</sup>Institute for Comparative Genomics, American Museum of Natural History, New York, NY

7 <sup>2</sup>Department of Epidemiology, Mailman School of Public Health, Columbia University, New  
8 York, NY

9 <sup>3</sup>Center for Discovery and Innovation, Hackensack Meridien Health, Nutley, NJ

10 <sup>4</sup>Division of Infectious Diseases, Children's Hospital of Philadelphia, Philadelphia, PA

11 <sup>5</sup>Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania, United  
12 States of America

13

14 \*Correspondence to: Apurva Narechania (anarechania@amnh.org) & Paul J. Planet  
15 (planetp@email.chop.edu)

16

17 **Word Count**

18 Abstract: 263

19 Body: 4824

20

21 Running Title: Information bottleneck for detecting recombination

22

## 23 **Abstract**

24 Most microbes have the capacity to acquire genetic material from their environment.  
25 Recombination of foreign DNA yields genomes that are, at least in part, incongruent with the  
26 vertical history of their species. Dominant approaches for detecting such horizontal gene  
27 transfer (HGT) and recombination are phylogenetic, requiring a painstaking series of analyses  
28 including sequence-based clustering, alignment, and phylogenetic tree reconstruction. Given  
29 the breakneck pace of genome sequencing, these traditional pan-genomic methods do not  
30 scale. Here we propose an alignment-free and tree-free technique based on the sequential  
31 information bottleneck (SIB), an optimization procedure designed to extract some portion of  
32 relevant information from one random variable conditioned on another. In our case, this joint  
33 probability distribution tabulates occurrence counts of k-mers with respect to their genomes of  
34 origin (the relevance information) with the expectation that HGT and recombination will create  
35 a strong signal that distinguishes certain sets of co-occurring k-mers. The technique is  
36 conceptualized as a rate-distortion problem. We measure distortion in the relevance  
37 information as k-mers are compressed into clusters based on their co-occurrence in the source  
38 genomes. This approach is similar to topic mining in the Natural Language Processing (NLP)  
39 literature. The result is model-free, unsupervised compression of k-mers into genomic topics  
40 that trace tracts of shared genome sequence whether vertically or horizontally acquired. We  
41 examine the performance of SIB on simulated data and on the known large-scale  
42 recombination event that formed the *Staphylococcus aureus* ST239 clade. We use this  
43 technique to detect recombined regions and recover the vertically inherited core genome with  
44 a fraction of the computing power required of current phylogenetic methods.

45

## 46 Introduction

47 Whole microbial genomes are being sequenced at an unprecedented rate.<sup>1</sup> Focused  
48 sequencing of key organisms and broad sequencing of microbial environments have expanded  
49 our knowledge of evolution and the microbiosphere<sup>234</sup>. However, the production of data is  
50 outstripping our ability to analyze it<sup>5</sup>. Most work in molecular evolution is grounded in  
51 sequence alignment and phylogenetic tree reconstruction. However, whole genome alignment  
52 breaks down with increasing diversity, and tree-based techniques suffer from an exponential  
53 increase in compute time with broader taxon sampling. The evolution of microbes is particularly  
54 challenging because horizontally transferred elements contribute historical signal that is  
55 unrelated to vertical descent. Most dominant techniques for capturing horizontal gene transfer  
56 (HGT) and recombination require either alignment of reads across a reference genome (eg.,  
57 single nucleotide polymorphism (SNP) based analysis or whole genome alignment<sup>67</sup>. Where  
58 global alignment is impossible, phylogenomic tools require all-against-all analyses designed to  
59 fix genes into aligned orthologous groups<sup>8910</sup>. All of these approaches require careful curation,  
60 tree-building, HGT/Recombination detection analysis, and deliberate sampling to limit data to  
61 reasonable scales. For larger, unbiased datasets that include as much natural variation as  
62 possible, these approaches are not sustainable. To handle the onslaught of genomes, we need  
63 tools that can tolerate information loss without sacrificing knowledge of key evolutionary  
64 events.

65 Lossy compression, where an individual or algorithm makes decisions about which data  
66 are important (or relevant) from a large body of information<sup>11</sup>, may offer a solution. To do this  
67 in a principled way, the relevance of a given dataset can be measured as information retained

68 about some other correlated variable. For example, in unsupervised natural language  
69 processing (NLP) large corpora of texts are distilled to a few topics that reflect overall themes  
70 by comparing patterns of co-occurring words in the source texts. In topic modeling of this sort,  
71 the texts themselves are the relevance variable. The goal is to cluster the overall word  
72 distribution with respect to the documents from which they arise. If  $X$  is the original data  
73 distribution,  $T$  its compressed representation, and  $Y$  the relevance variable, the challenge is to  
74 pack  $X$  into as few clusters,  $T$ , as possible without sacrificing too much information,  $Y$ . This idea  
75 was first described by Tishby, Pereira and Bialek as the information bottleneck (IB)<sup>12</sup>. It was  
76 premised on rate distortion, Shannon's original theory of lossy compression which yoked signal  
77 distortion to the rate at which that signal can be encoded<sup>13</sup>. Distortion is severe if the signal is  
78 forced through a small communication channel and gets cleaner as the channel widens. The IB's  
79 primary innovation was the use of a relevance variable to quantify this distortion. Topic  
80 modeling was one of this technique's first applications.

81 Topic modeling has become an important part of the NLP literature with a number of  
82 wider applications to unsupervised machine learning. The dominant technique in the field is  
83 Latent Dirichlet Allocation (LDA)<sup>14</sup>, a probabilistic method, that like the IB, considers each  
84 document as a mixture of topics. Some groups have applied this idea to whole genomes<sup>15,16,17</sup>,  
85 and since the publication of STRUCTURE, LDA has become foundational in the genetics  
86 literature where populations are inferred by the distribution of alleles at measured loci<sup>18</sup>.  
87 Despite LDA's popularity and success, a number of authors have shown that unbalanced  
88 sampling can lead to erroneous or missed population assignments<sup>19</sup>. LDA also makes a number  
89 of statistical assumptions including the assignment of hyperparameters and a Dirichlet prior<sup>20</sup>.

90 In contrast, the IB is model free and less likely to suffer from size sample bias. The distortion  
91 measure emerges from the analysis of the relevance variable, revealing underlying topics  
92 without having to set any distributional parameters other than the number of clusters  
93 expected.

94 Because it is model free, the IB is a powerful approach for microbial genomics where  
95 very little is known about the diversity of the organisms in nature or their distribution.  
96 Genomes are living documents that can be sliced into words of arbitrary size. This metaphor is  
97 straightforward and has been explored with respect to other NLP techniques elsewhere<sup>212223</sup>. In  
98 a genomic context, where words are k-mers (X) and documents (Y) are their genomes of origin  
99 we hypothesized that IB derived topics (T) may represent co-occurring groups of k-mers that  
100 highlight shared ancestry. These topics might include k-mers arranged in co-linear blocks  
101 corresponding to a single element, or k-mers distributed across the genome that were inherited  
102 in concert. In either case, compression of these k-mers into topics is guided by how often they  
103 co-occur with respect to their genomes of origin. This mechanism will tend to group adjacent k-  
104 mers in a recombined region because the recombination event is likely restricted to just a  
105 subset of taxa. Additionally, shared tracts of co-occurring k-mers common to all genomes, offer  
106 a simple, operational definition of a genomic “core”.<sup>76</sup> For microbial genomes where HGT is  
107 rampant<sup>2425</sup> we can therefore use the technique to learn which portions of the genome form  
108 the vertically inherited core, and which portions have been recombined, or inherited  
109 horizontally. In the NLP topic modeling analogy, the core genome of a species could be  
110 considered the set of meaningful words across every book in a specialized library, while  
111 recombined regions are like themes or ideas restricted to only certain shelves.

112 Here we apply the IB to microbial genomes. Remarkably, our approach identifies  
113 recombination tracts without making any attempt to model evolution, annotate genes,  
114 reconstruct trees, or build alignments. In addition, the IB treats genic and intergenic portions of  
115 the genome equally, obviating the need for gene-based pangenomic analysis<sup>26</sup>. Applying the  
116 information bottleneck to a k-mer occurrence matrix identifies genome segments with shared  
117 vertical or horizontal evolutionary history in a fraction of the time used by other approaches.

118

## 119 Theory and Implementation

120

121 Consider a set of genomes each of which is chopped into overlapping k-mers. One way  
122 to measure the overall relatedness of two of these genomes is to compare their k-mer  
123 conditional distributions. To do this we can define

124

$$125 \quad p(x|y) = \frac{n(x|y)}{\sum_y n(x|y)}$$

126

127 where X is the set of all k-mers, Y the set of all genomes, and  $n(x|y)$  is the occurrence count of  
128 the k-mer, x, in genome y. The exercise would then be to group genomes with similar k-mer  
129 distributions across all k-mers. In the natural language processing literature, this idea was  
130 formalized as distributional clustering<sup>27</sup>.

131 However, finding the right distance or distortion measure between these distributions is  
132 non-trivial. It is especially difficult when the important features of the signal are unknown.  
133 Imagine compressing music into MP3s without data on which frequencies are most important

134 for human perception, or determining themes from a body of literature if words were  
135 decoupled from their books. Even when important components of the signals are known, most  
136 clustering algorithms will resort to domain specific, pairwise distances or quantization to find a  
137 compressed set of classes with either high levels of internal connectivity or low levels of  
138 internal distortion. However, domain specific distortions reduce the usefulness of these  
139 clustering techniques. For example, in bioinformatics, clustering based on sequence alignment  
140 is subject to all the vagaries of the alignment procedure and parameters therein.

141 An antidote to these narrow clustering applications is to operate in an information  
142 theoretic space where the primary measurement is relevant quantization<sup>12</sup>. The IB extends  
143 Shannon's rate distortion theory by guiding it with an additional, orienting variable. Tishby et  
144 al<sup>12</sup> enriched a theory about transmission efficiency with the concept of relevance (Y), or the  
145 value of the information transmitted. The choice of Y defines relevant features in the signal. If X  
146 and Y are tabulated as a joint probability distribution, the information that X provides about Y is  
147 squeezed through a simpler representation, T. For the technique to work, the two variables in  
148 our joint distribution  $p(x,y)$  must be non-independent, or more precisely, must have positive  
149 mutual information,  $I(X,Y)$ :

150

$$151 \quad I(X,Y) = \sum_x \sum_y p(x)p(y|x) \log \frac{p(y|x)}{p(y)}$$

152

153 T is now a meaningful compression of the data, maximizing the mutual information between  
154 the clusters and documents,  $I(T;Y)$ , while minimizing the mutual information between the  
155 words and the clusters,  $I(T;X)$ . The IB is a classic optimization problem.

156           With the distribution in hand and implemented as a k-mer occurrence matrix, we can  
157   quantize the set of all k-mers directly by minimizing information lost about their source  
158   genomes. If  $X$  is compressed into  $T$  then we can find the optimal assignments for  $X$  by  
159   minimizing the following Lagrangian with respect to  $Y$ :

160

$$161 \quad \mathcal{L}[p(t|x)] = I(X;T) - \beta I(X;Y)$$

162

163   This formulation balances the compactness of  $X$ , with the erosion of information about  $Y$ .  $\beta$  is a  
164   multiplier that slides through the optimization landscape. As beta approaches 0, k-mers are  
165   clumped into fewer and fewer clusters, emphasizing compression. As beta approaches infinity,  
166   every k-mer is its own cluster, preserving all relevant information. Of course, collapsing all k-  
167   mers into one cluster is overly reductive, and assigning each k-mer to its own cluster is  
168   meaningless. The IB negotiates these two extremes (Figure 1). In NLP, the result is a set of  
169   clusters that coalesce into topics over a body of literature<sup>28</sup>. In genomics, these same clusters  
170   might yield co-occurring and/or spatially co-located k-mers with distinct biological and/or  
171   evolutionary meaning.

172           Remarkably, minimizing the Lagrangian above has an exact, optimal solution<sup>12</sup>. The most  
173   surprising outcome of this solution is that the relative entropy, or Kullback Liebler divergence<sup>29</sup>,  
174   emerges as the distortion measure for the information bottleneck. The relative entropy is a  
175   fundamental quantity in information theory, and in the IB context, it measures the distortion  
176   between the points,  $x$  (k-mers), as they are quantized into their clusters,  $t$ , with respect to the  
177   relevance variable,  $y$  (genomes):



178

$$179 \quad D_{KL} = \sum_y p(y|x) \log \frac{P(y|x)}{P(y|t)}$$

180

181 Calculation of the optimal solution requires soft clustering, that is, any given k-mer can exist in  
182 more than one cluster. But soft clustering can be slow and difficult to devise. Early  
183 implementations of the information bottleneck therefore settled on hard clustering  
184 approximations. In hard or deterministic clustering, each k-mer is assigned to only one cluster,  
185 an assumption that eases computational burden but does not generally arrive at globally  
186 optimal solutions.

187 The most obvious hard clustering algorithm is agglomerative, or bottom-up<sup>30</sup>. Consider  
188 again the set of all genomes,  $X$ , and their compressed representation,  $T$ . If we start with a  
189 scenario where every k-mer in  $X$  occupies its own singleton cluster, we can systematically  
190 reduce the dimensionality by merging clusters that minimize some distortion score. This greedy  
191 merging procedure produces a tree. But agglomerative clustering does not yield stable cluster  
192 membership. The tree varies every time the process is reinitialized. Worse, its computation is  
193 expensive, requiring cubic time complexity and quadratic memory complexity. In a genomic  
194 context where we routinely deal with billions of k-mers, this approach is a nonstarter.

195 Instead, we implemented a sequential clustering procedure where the number of  
196 clusters is defined at the outset and remains consistent throughout the calculation. From an  
197 initial random distribution of all k-mers across this set of clusters, we draw one k-mer out, and  
198 represent it as a singleton. Now using greedy optimization, we merge this singleton into one of

199 the existing bulk clusters. Slonim's sequential information bottleneck (SIB)<sup>31</sup> employs the  
200 Jensen-Shannon divergence<sup>32,21</sup> in the cost of merging a k-mer,  $x$ , into a cluster,  $t$ :

201

$$202 \quad d(x, t) = (p(x) + p(t)) * D_{JS}(p(y|x), p(y|t))$$

203

204 A k-mer will join a new cluster only if its new address reduces the total distortion. Otherwise it  
205 will remain in its existing cluster. With respect to our initial random conditions, this algorithm is  
206 guaranteed to converge to a local optimum. We mitigate the risk of getting trapped in local  
207 optima by testing several random initializations.

208         Once the clusters stabilize, we quantify the information captured by calculating the  
209 normalized mutual information,  $NMI = I(T;X) / I(X;Y)$ . Trivially,  $NMI = 1$  when each k-mer  
210 occupies its own cluster. The curve traced between  $T = 1$  ( $NMI = 0$ ) and  $T = x$  is called the  
211 relevance compression curve<sup>33</sup>. This is analogous to the optimization of  $\beta$  in the Lagrangian  
212 above, but for the deterministic case involving hard clustering. As with  $\beta$ , the shape of this  
213 curve describes the compressibility of the data.

214         The most important aspect of the SIB, and the reason we chose it for this work, is that it  
215 makes the concept of the information bottleneck accessible to modern genomics. The time  
216 complexity is linear in the number of k-mers and the number of clusters. This improvement  
217 makes information theoretic NLP a useful tool to discover genomic topics encoded as clusters  
218 of co-occurring k-mers.

219

220 **Results and Discussion**

221

222 *The bottleneck in test: one large, simulated HGT event*

223           The simple example in Figure 2 illustrates how the bottleneck works in practice. In  
224 SimBac<sup>34</sup>, we simulated four 1 megabase genomes with a single 200 kilobase recombination  
225 event. The event is common to genomes 0, 2 and 3, but is not found in strain 1. We initialized  
226 the simulation with a random distribution of 19-mers across five clusters. To learn the true  
227 distribution, we leveraged information in our relevance variable, the source genomes. The inset  
228 table shows how this distribution evolves as we iterate through the sequential information  
229 bottleneck (SIB). Since the relevance variable is expected to drive the unsupervised  
230 compression of these k-mers, we also included the genomes in this table. Counts across each  
231 row therefore reflect how many times a k-mer in that cluster is found in a particular genome.

232           The SIB starts by randomly distributing the k-mers, destroying all information available  
233 in the original occurrence matrix. At the outset, the normalized mutual information is therefore  
234 zero. With each SIB loop, we attempt to reclaim as much of this information as possible given  
235 the number of clusters we choose to model. Because the technique is inherently lossy, the SIB  
236 will never recover all of the information originally encoded, but aims to extract the most salient  
237 themes, or topics.

238           In the example shown here, after the first loop, cluster 3 (the cluster designations are  
239 arbitrary) has attracted the most k-mers in roughly even proportion across the genomes. The  
240 normalized mutual information has also jumped to 0.69, indicating that just one pass of sorting  
241 k-mers into five bins effectively captures 70% of the information available in the original  
242 occurrence matrix. The second and third loops refine the other clusters into mutually exclusive

243 sets and add to cluster 3, which strengthens into a genomic “core” defined here as the cluster  
244 of k-mers with the highest average representation across all genomes and the lowest index of  
245 dispersion.

246 By the third pass through the k-mers, the SIB reaches a plateau in the normalized  
247 mutual information, and the counts of k-mers across clusters and genomes have stabilized. For  
248 this particular set of starting conditions, the SIB reclaims nearly 91% of the information in the  
249 original matrix. To put this in perspective, we have effectively reduced the outsized,  
250 uninterpretable dimensions of our original data – 1.25 million unique k-mers – into the 5  
251 clusters we set out to model, while sacrificing only 9% of the original information present in the  
252 relevance variable.

253 In a genomic context, we hypothesized that the spatial organization of k-mer clusters  
254 would correspond to areas of common ancestry. In Figure 2, we mapped k-mers from various  
255 clusters to the genome backbones of strain 1 and strain 2. Cluster 3 occupies the outer tracks of  
256 both strains. This cluster emerges as a dense block of shared genome sequence and  
257 corresponds to our definition of a bottleneck-defined core. But the block is interrupted by our  
258 simulated recombination event. Since this event is restricted to only genomes 0, 2 and 3, the  
259 region is absent from the core. Its k-mers are instead captured by cluster 4 while cluster 1  
260 serves as a counterpoint, containing the ancestral state prior to the simulated event.

261

### 262 Several smaller, simulated HGT events

263 Though large hybridization events like the one we simulated here do occur (see our  
264 analysis of ST239 *S. aureus* below), smaller and more abundant events typify most microbial

265 evolution<sup>35</sup>. To see how the bottleneck performs in this more challenging case, we simulated  
266 ten 1 megabase genomes with a background mutation rate of 0.01 and a recombination rate of  
267 0.0001, resulting in 57 discrete events averaging 500 basepairs in size (from 6 to 2884 bases). In  
268 **Figure 3**, the innermost track marks the locations of these events.

269         The ability to detect horizontally transferred sequence is strongly dependent on its  
270 evolutionary distance from the genome background<sup>7</sup>. To visualize this dependence, we  
271 modulated the divergence of our 57 recombination events (an arbitrary number derived from  
272 the first simulation) and measured the effect on the core cluster, one of 60 modeled for this  
273 simulation. The innermost histogram in Figure 3 shows the core pattern with an external  
274 (between species) divergence rate of 0.1, an order of magnitude higher than the background.  
275 We observe clear “valleys” in the k-mer distribution of the core that are coincident with the  
276 positions of our 57 events. But this pattern steadily disappears as we sweep through lower  
277 rates of divergence (0.05, 0.03, and 0.01). The outermost track models the same mutation rate  
278 as the background, resulting in dulled or partially filled valleys in the core genome. Plots of core  
279 k-mers function almost as a photographic negative, highlighting blank spaces as regions of  
280 potential evolutionary interest.

281         The k-mers that would otherwise occupy these gaps, are sorted into other clusters  
282 because they are unique to only a subset of the genomes, and carry the recombination signal.  
283 As we have shown in our first simulation, k-mers corresponding to the ancestral state should  
284 fall into a different cluster. Note that this does not necessarily mean that each side (donor and  
285 recipient) of an HGT event has its *own* cluster. Recall that compression is driven by genome  
286 origin. If a single common ancestor sustains multiple transfer events, all k-mers from those

287 events will merge into a single cluster because they are shared by the same subset of  
288 descendants.

289         The accounting becomes increasingly complicated when events overlap. Overlapping  
290 events might mix across clusters depending on their arrangement and how frequently they  
291 have been overwritten. When detection becomes difficult, we instead rely on an evolutionary  
292 event's imprint on the core cluster. This approach exploits the idea of the core as a  
293 photographic negative or a clonal frame. The pattern of HGT events in this negative is evident  
294 by eye, but if the number of input genomes and the number of modeled clusters is large, visual  
295 inspection is a burden, and subject to error in interpretation. Instead we introduce a method  
296 based in change point detection to automatically detect changes in k-mer frequency<sup>36</sup>. We  
297 specifically employ Bayesian change point detection<sup>37</sup> to model probabilities of change in the k-  
298 mer frequency stream. As shown in Figure 4 change point probabilities spike at the start and  
299 end of HGT events.

300         In addition to change point detection, we note that if counts of k-mers in an HGT region  
301 are significantly lower than the rest of the core's background (Wilcoxon,  $p < 0.05$ ), these  
302 depletions can qualify as a simple signal marking some combination of HGT events. With these  
303 criteria, at a divergence rate of 0.1, the bottleneck captures 56 of the 57 simulated events,  
304 missing only the smallest.

305

### 306 *The k-mer skim*

307         Accounting for every overlapping k-mer in each strain is an unnecessarily close reading  
308 of our genomic text. We can save on both memory and computation by selecting fewer k-mers

309 (skimming) from our source genomes with some set space between each sample. In Figure 5 we  
310 show that even when sampling every 25<sup>th</sup> 19-mer in our ten 1 Mbase simulated genomes, we  
311 still detect 55 of our 57 recombination events. Because the bottleneck relies on the signal  
312 inherent in k-mer co-occurrence, as we reduce the density of our k-mer sampling, we lose  
313 detection of the smallest events first. However, the compute time savings more than  
314 compensate for this loss in sensitivity. While analyzing every 19-mer requires nearly 12  
315 minutes, skimming every 25<sup>th</sup> reduces the runtime to 30 seconds. This compares favorably with  
316 the efficiency of both ClonalFrameML<sup>7</sup> and Gubbins<sup>6</sup>, the two dominant HGT detection  
317 methods in the literature. ClonalFrameML requires 110 seconds and captures only 47 of our 57  
318 events. Gubbins finds 54 in 21 seconds. However, both ClonalFrameML and Gubbins require  
319 alignment and phylogenetic tree reconstruction, which both add massive prior computational  
320 cost and time.

321       Because the IB is alignment-free and tree-free, it is theoretically capable of handling  
322 larger datasets than any existing technology in reasonable amounts of time. To test this, we  
323 simulated 1000 1 Mb genomes with the same parameters as the smaller dataset shown in  
324 Figure 3. The simulation generated 620 unique recombination events. ClonalFrameML detected  
325 564 (91%). Including time required to build a guide tree, this calculation consumed 32.5 CPU  
326 hours. Gubbins was slightly more accurate and significantly faster: 583 (95%) events over 16.3  
327 CPU hours. Using Figure 5 as a guide, we ran the 1000 genome dataset through the SIB using a  
328 25 base-pair skim. We detected an HGT imprint at 92% of sites in 1.5 CPU hours.

329

330 *How well does the IB hold up under extreme evolutionary pressure?*

331 To evaluate the performance of our technique with respect to recombination size and  
332 divergence rate, we simulated sets of ten 1 megabase (Mb) genomes for each variable. We set  
333 default parameters to 0.01 for background rate, 0.001 for recombination rate, 0.1 for HGT  
334 divergence rate, and 500 base pairs for average recombination tract size. We performed 100  
335 replicates at each size and rate, and measured the imprint of the simulated events on the core  
336 cluster without the skim feature. Figure 6A shows this sweep for recombination tract length,  
337 and Figure 6B, for recombination tract divergence. In both cases, we observe saturating  
338 behavior. We see recombination imprints at 90% accuracy when events are larger than 100  
339 base pairs with divergence rates of at least 0.02. Notably, our procedure can detect HGT in at  
340 least half of events that diverge at the very low rate of 0.005, well below the background. And  
341 only the very smallest recombination events (less than 7 basepairs) elude our technique  
342 completely.

343 Recombination tract length and divergence have direct and measureable effects on the  
344 efficacy of detection. As long as the total length of all recombination events is less than half the  
345 size of the genome, the core remains intact, and we can easily isolate HGT events of sufficient  
346 size and divergence. But recombination and background mutation rates are problematic  
347 because they redefine the core. For example, at high rates of recombination, every base of a 1  
348 Mb genome is likely scrambled. Under such flux, some sites recombine several times. A high  
349 background mutation rate also disrupts stretches of common sequence that mark the core. As  
350 these rates increase, the core genome itself erodes. To measure this phenomenon, we again  
351 simulated 100 sets of ten 1 Mb genomes across a variety of recombination and background  
352 mutation rates. All three curves in Figure 7 show a steep decline in the size of the core with



353 increasing recombination rate. At rates of 0.01 and 0.1, we see no shared core at all. Each  
354 genome has essentially rewritten itself into something distinct from all others. Core genome  
355 signal grows stronger with lower background mutation, but even with background mutation set  
356 to essentially zero, a high recombination rate destroys the core.

357

### 358 *The bottleneck in action: one large, real world hybridization event*

359 We used genomes from ST239 *Staphylococcus aureus* to illustrate that our method can  
360 corroborate known, large scale recombination events found in nature. The ST239 strain is a  
361 hybrid: a segment from a CC30 (clonal complex 30) donor replaced nearly 20% of the  
362 homologous region in a CC8 strain<sup>38</sup>. The evolutionary histories of genes across these segments  
363 are incongruent. Previous studies compared the histories of thousands of genes to reach this  
364 conclusion<sup>39</sup>. Here, we attempt to localize this same phenomenon using the co-occurrence  
365 pattern of k-mers alone. We chose 10 genomes (GCA\_000146385.1, GCA\_000012045.1,  
366 GCA\_000011505.1, GCA\_000011265.1, GCA\_000013425.1, GCA\_000204665.1,  
367 GCA\_000159535.2, GCA\_000027045.1, GCA\_000017085.1, and SA21300), sampled from both  
368 the donor clade (CC30), the recipient clade (CC8), and genomes outside of the evolutionary  
369 event. When cut into overlapping 19-mers (no skim), these 10 genomes dissolve into 28.8  
370 million k-mers, 4.72 million of which are unique.

371 Figure 8 highlights two of these 10 genomes, and three of the 60 clusters we modeled  
372 for this analysis. Both *S. aureus* COL (CC8) and *S. aureus* T0131 (ST239) share a large, congruent  
373 core. The gap in this core characterizes the dimensions of the recombination event, whose k-  
374 mers are split into two other clusters, shown here as the second and third tracks. Like subtopics

375 in a vast library, the bottleneck learns the complete structural evolution of the clade as tracts,  
376 or topics, of co-occurring sequence. The clusters themselves comprise an evolutionary model  
377 for the structural event and the core genome. This evolutionary model is derived not from  
378 traditional character-based phylogenetic analysis, but from the presence/absence pattern of k-  
379 mers squeezed into a predefined number of groups. Genome origin guides the k-mer sort by  
380 forming the basis of the distortion measure. We lose information in a controlled and  
381 quantitative way, and we short circuit the long and arduous tasks phylogenomic analyses  
382 require<sup>39</sup> with an information theoretic procedure that runs for 2 hours on 1 CPU.

383 By definition, this sort of lossy compression is not perfect. In Figure 8, seemingly  
384 unrelated contaminants pollute the recombined region's clusters. This is equivalent to channel  
385 noise. It recalls Shannon's original formulation of the rate distortion problem<sup>13</sup>. When we force  
386 all the signal in our k-mer occurrence matrix through a narrow five cluster channel, portions of  
387 the original message emerge garbled. In this case, modeling more clusters increases the rate of  
388 transmission, and reduces the distortion of the message received.

389 With respect to the information bottleneck, we can quantify this effect using a  
390 relevance-compression curve<sup>28</sup>. Figure 9 shows curves for the ST239 genomes alongside 10  
391 genomes of *Mycobacterium tuberculosis* and *Helicobacter pylori*. In all three cases, as the  
392 number of clusters modeled increases, we capture more normalized mutual information. The  
393 theoretical extremes for this curve are intuitive. At the origin, all the relevant information is  
394 destroyed. At the other end, we retain too much relevant information to interpret. The curve  
395 traced between these two extremes is a fingerprint of the data. A convex shape suggests  
396 natural structure easily modeled with just a few clusters. We see this in *M. tuberculosis*, a

397 species thought to be largely clonal with little recombination. On the other hand, data that  
398 resists compression flattens this curve. Highly recombinogenic species like *H. pylori* suffer this  
399 sort of steep information loss. Theoretically, the space above the curve for each species is  
400 unachievable by any process, forming an upper bound. The relevance-compression curve  
401 therefore defines absolute limits on the quantity and quality of information communicated as  
402 we sweep through a dilating channel. This approach introduces a new type of comparative  
403 genomics based not on alignments and trees, but on compression. We interpret the shape of  
404 the relevance compression curve as a proxy for evolutionary mode. A convex curve implies  
405 fewer recombination events and more vertical signal, whereas a flattened curve may signal a  
406 species with a more open pangenome.

407         In the case of ST239, asking for just two clusters – a very narrow channel – captures  
408 more than 40% of the relevant information. Remarkably, these two clusters separate the core  
409 from the recombined region. Even the simplest model learns the most prominent evolutionary  
410 process. Further along the curve, fifteen clusters capture almost all of the information. Beyond  
411 fifteen, the curve elbows, and modeling gains are slight. In this way, the relevance-compression  
412 curve defines the optimal number of clusters.<sup>3340</sup> But in the light of evolution this bend may  
413 have a deeper meaning. Fifteen clusters are enough to adequately capture the complete set of  
414 k-mer aggregation patterns across our chosen genomes. This point of diminishing returns may  
415 signify an opportunity for interpretive balance: not so many clusters that we drown dominant  
416 evolutionary events, and not so few that we neglect to model subtle k-mer co-occurrence  
417 patterns. This particular use of the well-known elbow method in our information theoretic

418 context puts a crude limit on the dominant evolutionary paths taken by the genomic elements  
419 that comprise our species.

420  
421 *Conclusion (words=149)*

422 The information bottleneck, a lossy compression technique borrowed from the information  
423 theoretic and Natural Language Processing literature, is well suited to detecting evolutionary  
424 patterns in sets of co-occurring k-mers. Here we have shown that we can detect simulated and  
425 real recombination events while highlighting a core set of k-mers that comprise the vertically  
426 inherited portion of any set of genomes. Moreover, the compressibility of any given set of  
427 genomes, as embodied in their relevance compression curves, offers a new way to compare the  
428 pangenomes of very different clades in the microbial tree of life. In our application, the  
429 bottleneck is informed by genome origin, our relevance variable. But the technique is general.  
430 The information bottleneck can be used for any biological contingency matrix where the goal is  
431 to cluster a variable into interpretable groups by preserving as much information as possible in  
432 the variable to which it is linked.

433  
434 Software implementation: NECK (<https://github.com/narechan/neck>)

435  
436  
437 **Figure Legends**

438  
439 Figure 1. The information bottleneck. In the information bottleneck a distribution,  $X$ , is  
440 compressed into  $T$  while retaining as much information as possible about a correlated relevance  
441 variable,  $Y$ . The joint distribution,  $p(x,y)$ , has positive mutual information and the goal of the  
442 information bottleneck is to capture as much of that information as possible at interpretive

443 scale. The technique is a classic optimization problem wherein the mutual information between  
444  $T$  and  $X$  is minimized, while the mutual information between  $T$  and  $Y$  is maximized. At  
445 optimality,  $T$  is presumed to be a lossy but adequate model of  $X$ .

446

447 Figure 2. One simulated HGT event. A simple set of four simulated genomes with a single large  
448 transfer event is shown. The transfer occurs in the common ancestor to genomes 0, 2, and 3.  
449 The inset chart clearly shows that the k-mers corresponding to this event are captured by  
450 cluster 4, while the ancestral state is captured by cluster 1. K-mers from these clusters map to  
451 the location of the simulated event in genomes 0, 2 and 3 and genomes 1, respectively. Cluster  
452 3 is the core and contains only one gap corresponding to the HGT region.

453

454 Figure 3. Several simulated HGT events. The innermost ring of this circos plot shows the  
455 locations of 57 simulated HGT events across 10 1 Mbase genomes. The remaining concentric  
456 tracks plot the core set of k-mers as calculated by the information bottleneck. In the outermost  
457 frequency plot, the 57 HGT events diverge at the same mutation rate as the background, 0.01.  
458 Going in towards the center, we increase the HGT divergence rate of the events to 0.03, 0.05,  
459 and 0.1. Gaps in the core correspond with the simulated HGT events whose k-mers are sorted  
460 into other clusters.

461

462 Figure 4. Bayesian change point detection. The two innermost rings mirror those in Figure 3.  
463 The outermost ring plots the posterior probabilities of change in the k-mer frequencies.

464

465 Figure 5. The k-mer skim. Here we show the decrease in HGT detection sensitivity as a function  
466 of the density of k-mers sampled. The higher the k-mer skim factor (defined as the number of  
467 positions skipped before the next k-mer is sampled), the lower the density of k-mers subject to  
468 the information bottleneck. The inset shows the plateau behavior near the origin for k-mer  
469 skim factors of 1, 5, 10, 25, and 50.

470

471 Figure 6. Varying HGT length and divergence. HGT detection rates are shown with respect to  
472 increasing HGT length and divergence.

473

474 Figure 7. Varying recombination and background mutation rates. We measure the fraction of  
475 unique k-mers in each simulation captured by the core genome cluster as a function of  
476 recombination rate and background mutation rate. The core genome signal is strongest at low  
477 rates of recombination and background mutation. At higher recombination rates, there is no  
478 evidence for a core genome of any kind regardless of the background mutation rate.

479

480 Figure 8. Modelling ST239's hybridization event. We selected 10 *S. aureus* genomes to track the  
481 ST239 hybridization event with the information bottleneck. COL was chosen to represent the  
482 CC30 donor strain, and T0131 the CC8 acceptor. Of the 60 clusters we calculated, we show the  
483 three that capture the hybridization event. The innermost track is a frequency plot of k-mers  
484 that define the core. The second and third tracks are flipsides of the HGT event that created  
485 ST239.

486

487 Figure 9. Relevance compression curves. In an information bottleneck experiment, the  
488 relevance compression curve traces the increase in normalized mutual information with the  
489 number of clusters modeled. The curves quantify the amount of information lost at a given  
490 modeling threshold. We show how this type of relationship can function as a marker for  
491 evolutionary strategy by calculating curves for three very different groups of microbes: *M.*  
492 *tuberculosis*, a species thought to demonstrate little if any HGT; *S. aureus*, a species considered  
493 largely clonal with occasional HGT; and *H. pylori*, a species known to employ HGT as an engine  
494 for diversity.  
495

496

497 **References**

498

499 1. GenBank and WGS Statistics. <https://www.ncbi.nlm.nih.gov/genbank/statistics/>.

500 2. Hug, L. A. *et al.* A new view of the tree of life. *Nat. Microbiol.* **1**, (2016).

501 3. Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark matter.

502 *Nature* **499**, 431–437 (2013).

503 4. Sogin, M. L. *et al.* Microbial diversity in the deep sea and the underexplored “rare biosphere”.

504 *Proc. Natl. Acad. Sci.* **103**, 12115–12120 (2006).

505 5. Stephens, Z. D. *et al.* Big Data: Astronomical or Genomical? *PLOS Biol.* **13**, e1002195

506 (2015).

507 6. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant bacterial

508 whole genome sequences using Gubbins. *Nucleic Acids Res.* **43**, e15–e15 (2015).

509 7. Didelot, X. & Wilson, D. J. ClonalFrameML: Efficient Inference of Recombination in Whole

510 Bacterial Genomes. *PLOS Comput. Biol.* **11**, e1004041 (2015).

511 8. Chiu, J. C. *et al.* OrthologID: automation of genome-scale ortholog identification within a

512 parsimony framework. *Bioinformatics* **22**, 699–707 (2006).

513 9. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**,

514 3691–3693 (2015).

515 10. Zhao, Y. *et al.* PGAP: pan-genomes analysis pipeline. *Bioinformatics* **28**, 416–418

516 (2012).

517 11. Marzen Sarah E. & DeDeo Simon. The evolution of lossy compression. *J. R. Soc.*

518 *Interface* **14**, 20170166 (2017).

519 12. Tishby, N., Pereira, F. C. & Bialek, W. The information bottleneck method.

520 *arXiv:physics/0004057* (2000).



- 521 13. Shannon, C. E. A Mathematical Theory of Communication. *Bell Syst. Tech. J.* **27**, 379–  
522 423 (1948).
- 523 14. Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent Dirichlet Allocation. *J. Mach. Learn. Res.*  
524 **3**, 993–1022 (2003).
- 525 15. Liu, L., Tang, L., Dong, W., Yao, S. & Zhou, W. An overview of topic modeling and its  
526 current applications in bioinformatics. *SpringerPlus* **5**, 1608 (2016).
- 527 16. La Rosa, M., Fiannaca, A., Rizzo, R. & Urso, A. Probabilistic topic modeling for the  
528 analysis and classification of genomic sequences. *BMC Bioinformatics* **16**, S2 (2015).
- 529 17. Chen, X., Hu, X., Shen, X. & Rosen, G. Probabilistic topic modeling for genomic data  
530 interpretation. in *2010 IEEE International Conference on Bioinformatics and Biomedicine*  
531 *(BIBM)* 149–152 (2010). doi:10.1109/BIBM.2010.5706554.
- 532 18. Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of Population Structure Using  
533 Multilocus Genotype Data. *Genetics* **155**, 945–959 (2000).
- 534 19. Wang, J. The computer program structure for assigning individuals to populations: easy  
535 to use but easier to misuse. *Mol. Ecol. Resour.* **17**, 981–990 (2017).
- 536 20. Wallach, H. M., Mimno, D. M. & McCallum, A. Rethinking LDA: Why Priors Matter. in  
537 *Advances in Neural Information Processing Systems 22* (eds. Bengio, Y., Schuurmans, D.,  
538 Lafferty, J. D., Williams, C. K. I. & Culotta, A.) 1973–1981 (Curran Associates, Inc., 2009).
- 539 21. Sims, G. E., Jun, S.-R., Wu, G. A. & Kim, S.-H. Alignment-free genome comparison  
540 with feature frequency profiles (FFP) and optimal resolutions. *Proc. Natl. Acad. Sci.* **106**,  
541 2677–2682 (2009).
- 542 22. Cong, Y., Chan, Y. & Ragan, M. A. A novel alignment-free method for detection of  
543 lateral genetic transfer based on TF-IDF. *Sci. Rep.* **6**, 30308 (2016).

- 544 23. Cong, Y., Chan, Y. & Ragan, M. A. Exploring lateral genetic transfer among microbial  
545 genomes using TF-IDF. *Sci. Rep.* **6**, 29319 (2016).
- 546 24. Polz, M. F., Alm, E. J. & Hanage, W. P. Horizontal gene transfer and the evolution of  
547 bacterial and archaeal population structure. *Trends Genet.* **29**, 170–175 (2013).
- 548 25. Planet, P. J. Reexamining microbial evolution through the lens of horizontal transfer. *EXS*  
549 247–303 (2002).
- 550 26. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of *Streptococcus*  
551 *agalactiae*: Implications for the microbial “pan-genome”. *Proc. Natl. Acad. Sci.* **102**, 13950–  
552 13955 (2005).
- 553 27. Pereira, F., Tishby, N. & Lee, L. Distributional Clustering of English Words. in  
554 *Proceedings of the 31st Annual Meeting on Association for Computational Linguistics* 183–  
555 190 (Association for Computational Linguistics, 1993). doi:10.3115/981574.981598.
- 556 28. Slonim, N. The Information Bottleneck: Theory and Applications. *Dr. Diss. Hebr. Univ.*  
557 *Jerus. Isr.* 2003 157.
- 558 29. Kullback, S. & Leibler, R. A. On Information and Sufficiency. *Ann. Math. Stat.* **22**, 79–  
559 86 (1951).
- 560 30. Slonim, N. & Tishby, N. Agglomerative Information Bottleneck. in *Proceedings of the*  
561 *12th International Conference on Neural Information Processing Systems* 617–623 (MIT  
562 Press, 1999).
- 563 31. Slonim, N., Friedman, N. & Tishby, N. Unsupervised Document Classification Using  
564 Sequential Information Maximization. in *Proceedings of the 25th Annual International ACM*  
565 *SIGIR Conference on Research and Development in Information Retrieval* 129–136 (ACM,  
566 2002). doi:10.1145/564376.564401.

- 567 32. Lin, J. Divergence measures based on the Shannon entropy. *IEEE Trans. Inf. Theory* **37**,  
568 145–151 (1991).
- 569 33. Still, S. & Bialek, W. How Many Clusters? An Information-Theoretic Perspective.  
570 *Neural Comput* **16**, 2483–2506 (2004).
- 571 34. Brown, T., Didelot, X., Wilson, D. J. & Maio, N. D. SimBac: simulation of whole  
572 bacterial genomes with homologous recombination. *Microb. Genomics* **2**, (2016).
- 573 35. Didelot, X. & Maiden, M. C. J. Impact of recombination on bacterial evolution. *Trends*  
574 *Microbiol.* **18**, 315–322 (2010).
- 575 36. Truong, C., Oudre, L. & Vayatis, N. Selective review of offline change point detection  
576 methods. *Signal Process.* **167**, 107299 (2020).
- 577 37. Barry, D. & Hartigan, J. A. A Bayesian Analysis for Change Point Problems. *J. Am. Stat.*  
578 *Assoc.* **88**, 309 (1993).
- 579 38. Robinson, D. A. & Enright, M. C. Evolution of *Staphylococcus aureus* by Large  
580 Chromosomal Replacements. *J. Bacteriol.* **186**, 1060–1064 (2004).
- 581 39. Narechania, A. *et al.* Clusterflock: a flocking algorithm for isolating congruent  
582 phylogenomic datasets. *GigaScience* **5**, (2016).
- 583 40. Slonim, N., Atwal, G. S., Tkačik, G. & Bialek, W. Information-based clustering. *Proc.*  
584 *Natl. Acad. Sci.* **102**, 18297–18302 (2005).

585  
586  
587  
588

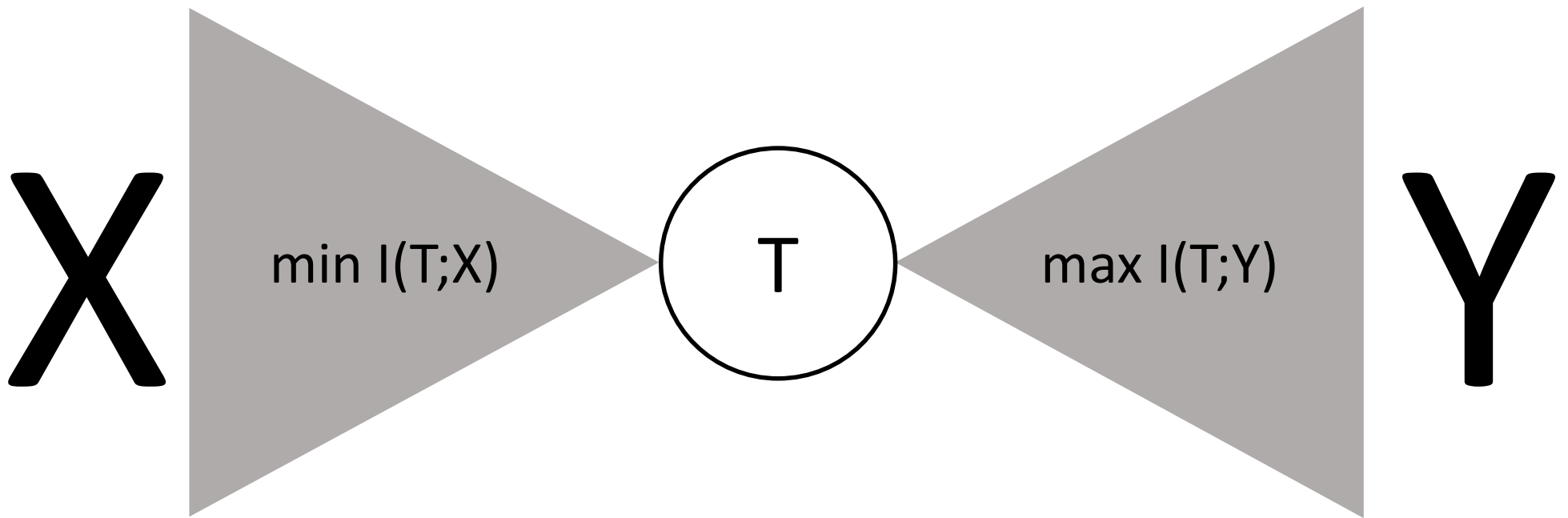
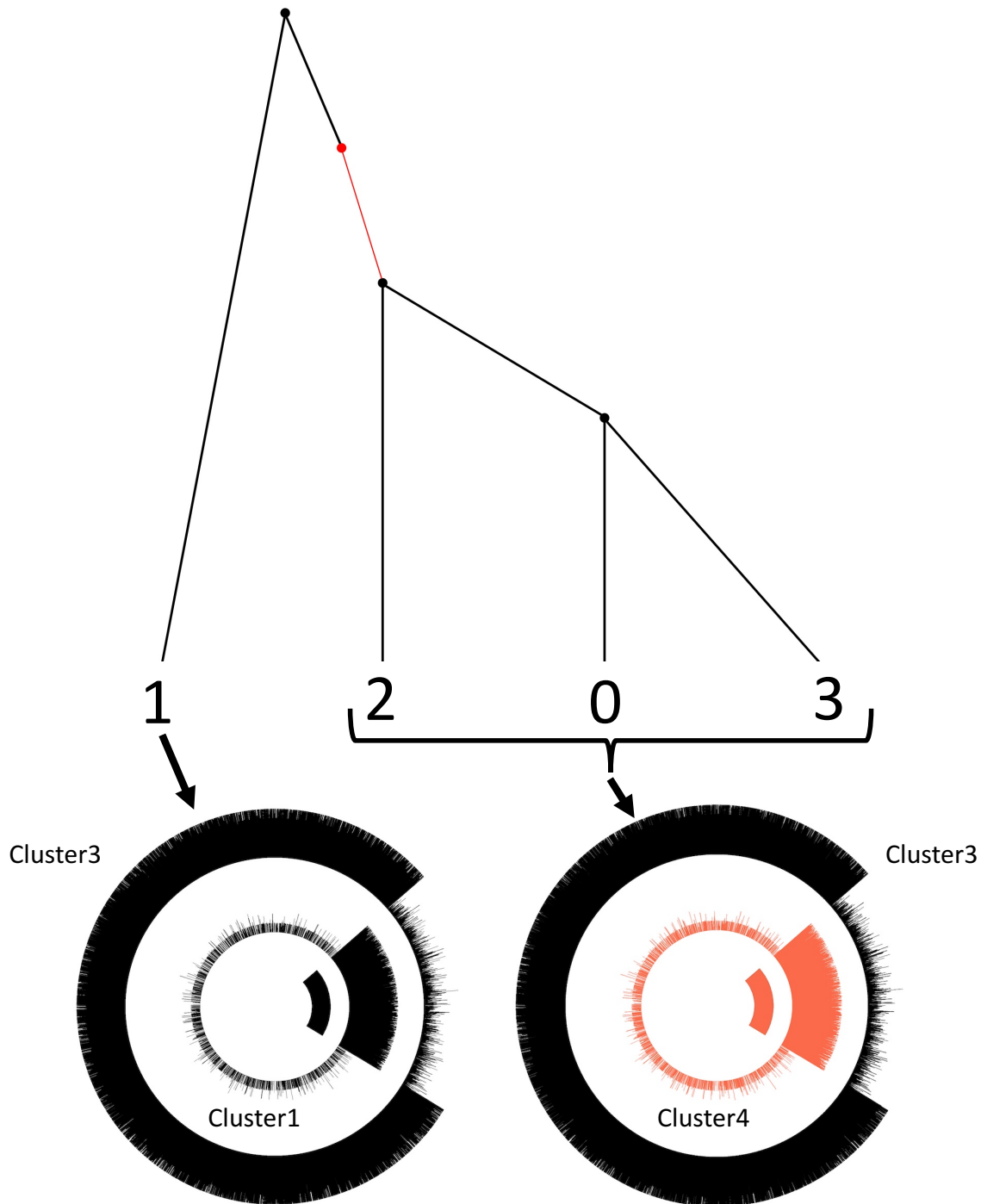


Figure 1



<b>Initialize</b>	genome0	genome1	genome2	genome3
CLUST0	199366	199350	199437	199282
CLUST1	200184	200439	200134	200194
CLUST2	199693	199696	199591	199808
CLUST3	200765	200718	200857	200785
<b>CLUST4</b>	<b>199974</b>	<b>199779</b>	<b>199963</b>	<b>199913</b>
<b>NMI = 0</b>				
<b>loop 1</b>				
CLUST0	3270	3954	4079	0
CLUST1	15808	237766	19369	15327
CLUST2	9348	10535	1175	10552
CLUST3	746003	747727	747728	747728
<b>CLUST4</b>	<b>225553</b>	<b>0</b>	<b>227631</b>	<b>226375</b>
<b>NMI = 0.69</b>				
<b>loop 2</b>				
CLUST0	7093	19812	50361	0
CLUST1	0	206335	0	0
CLUST2	28277	21421	0	37257
CLUST3	746103	752414	752413	752314
<b>CLUST4</b>	<b>218509</b>	<b>0</b>	<b>197208</b>	<b>210411</b>
<b>NMI = 0.89</b>				
<b>loop 3</b>				
CLUST0	0	12719	43268	0
CLUST1	0	206335	0	0
CLUST2	48715	21421	0	51754
CLUST3	753196	759507	759506	752314
<b>CLUST4</b>	<b>198071</b>	<b>0</b>	<b>197208</b>	<b>195914</b>
<b>NMI = 0.91</b>				

Figure 2

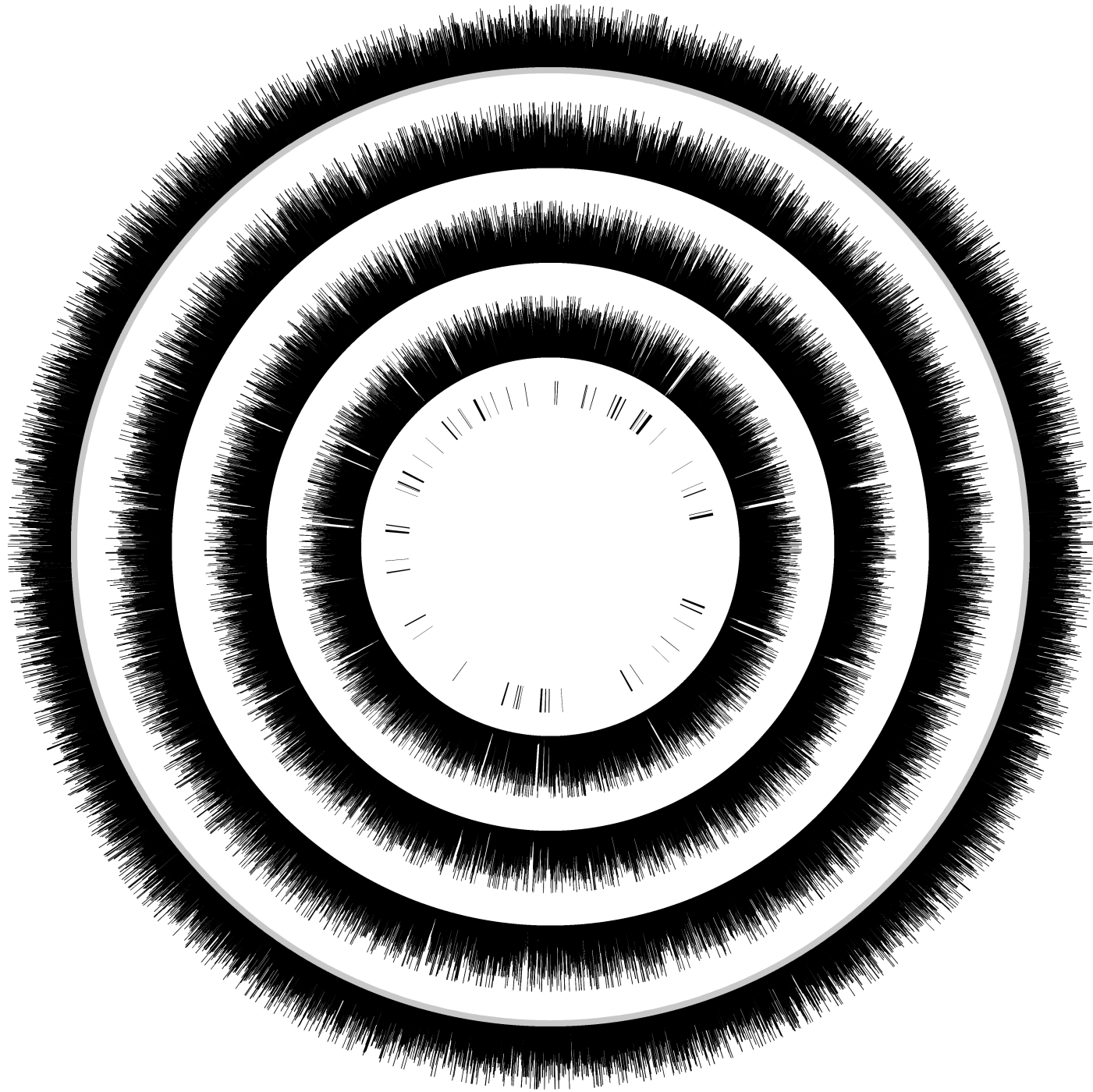


Figure 3

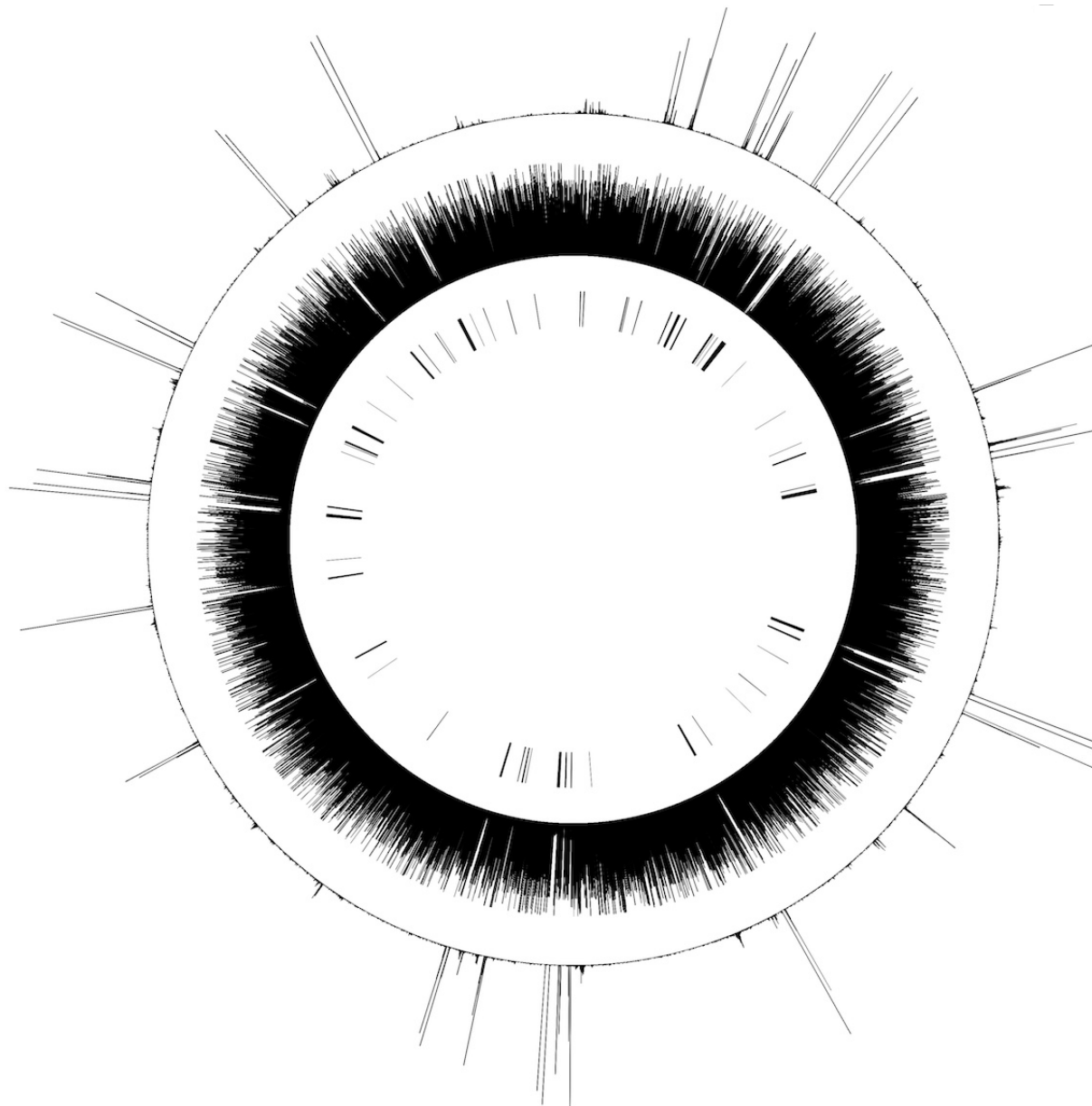


Figure 4

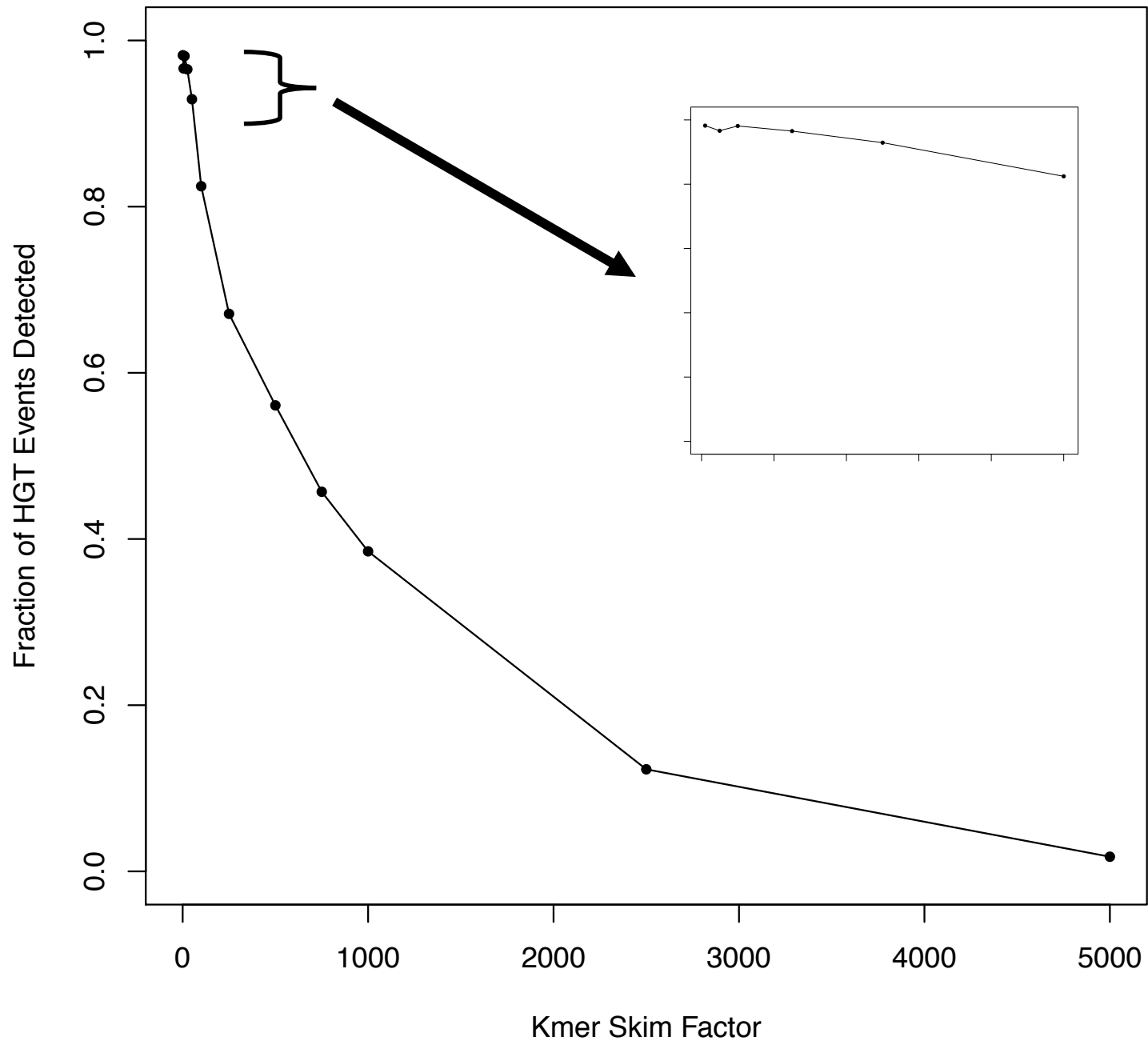


Figure 5



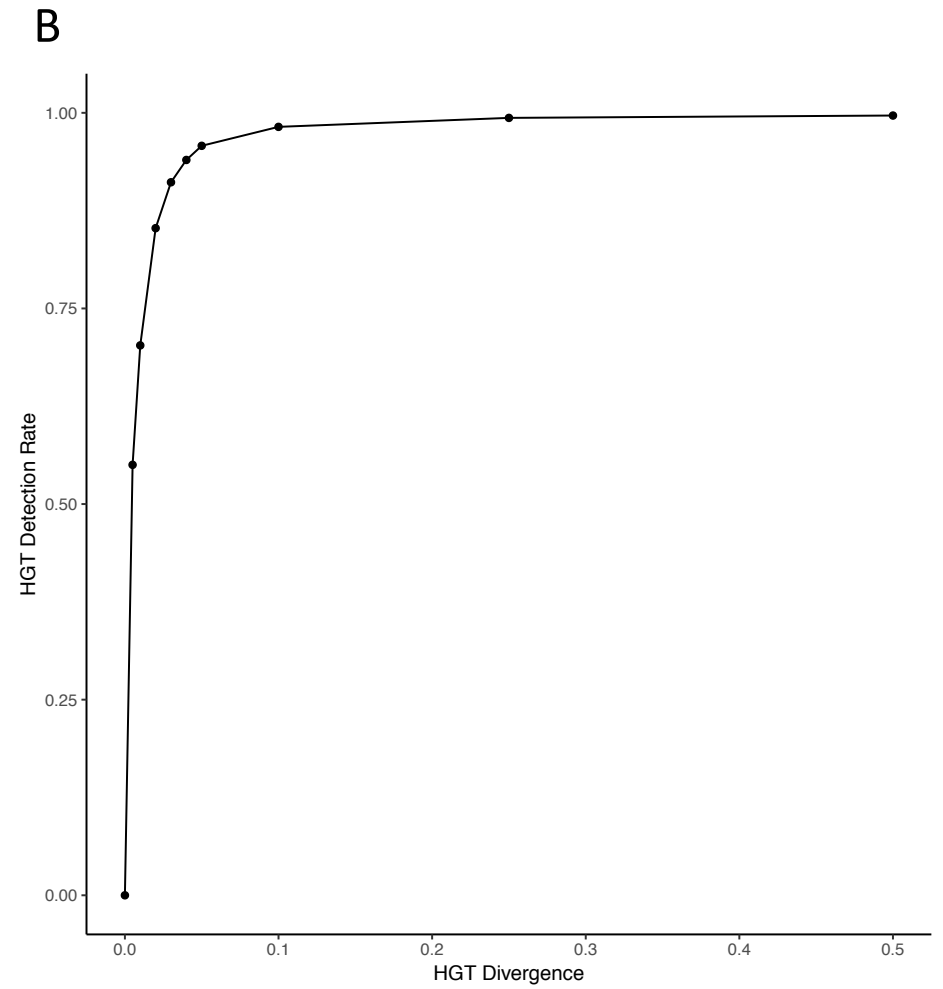
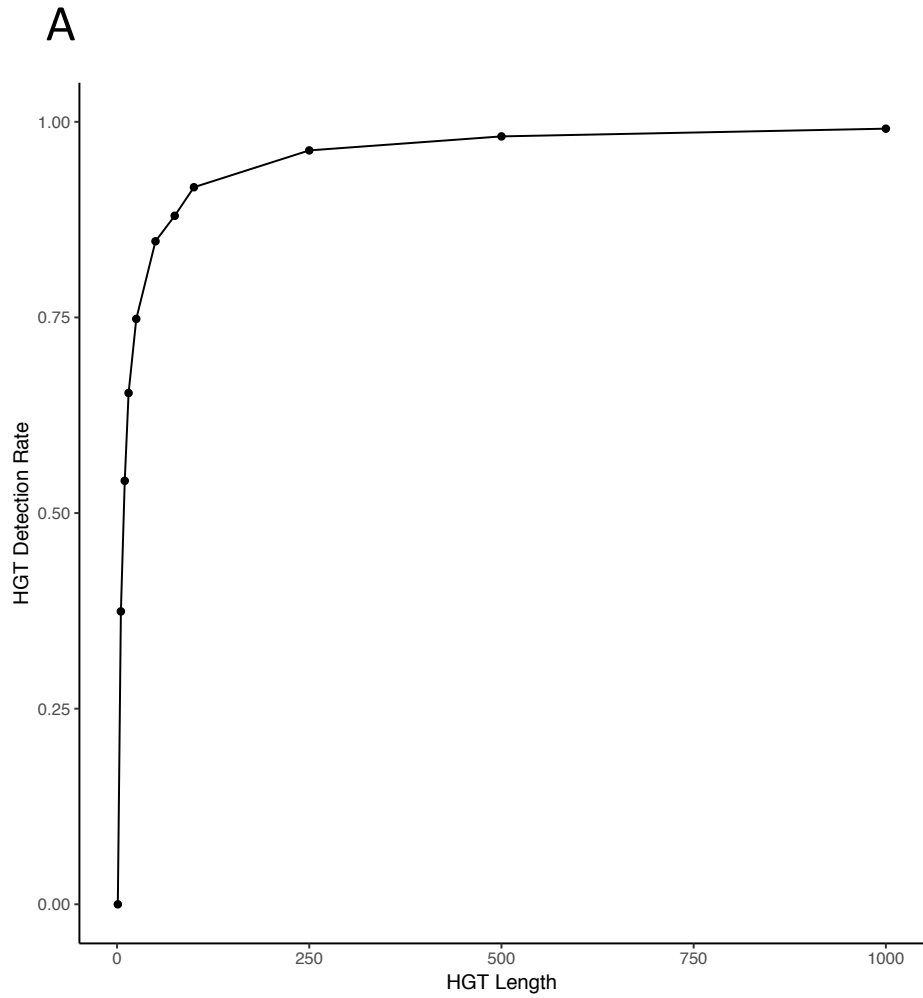


Figure 6

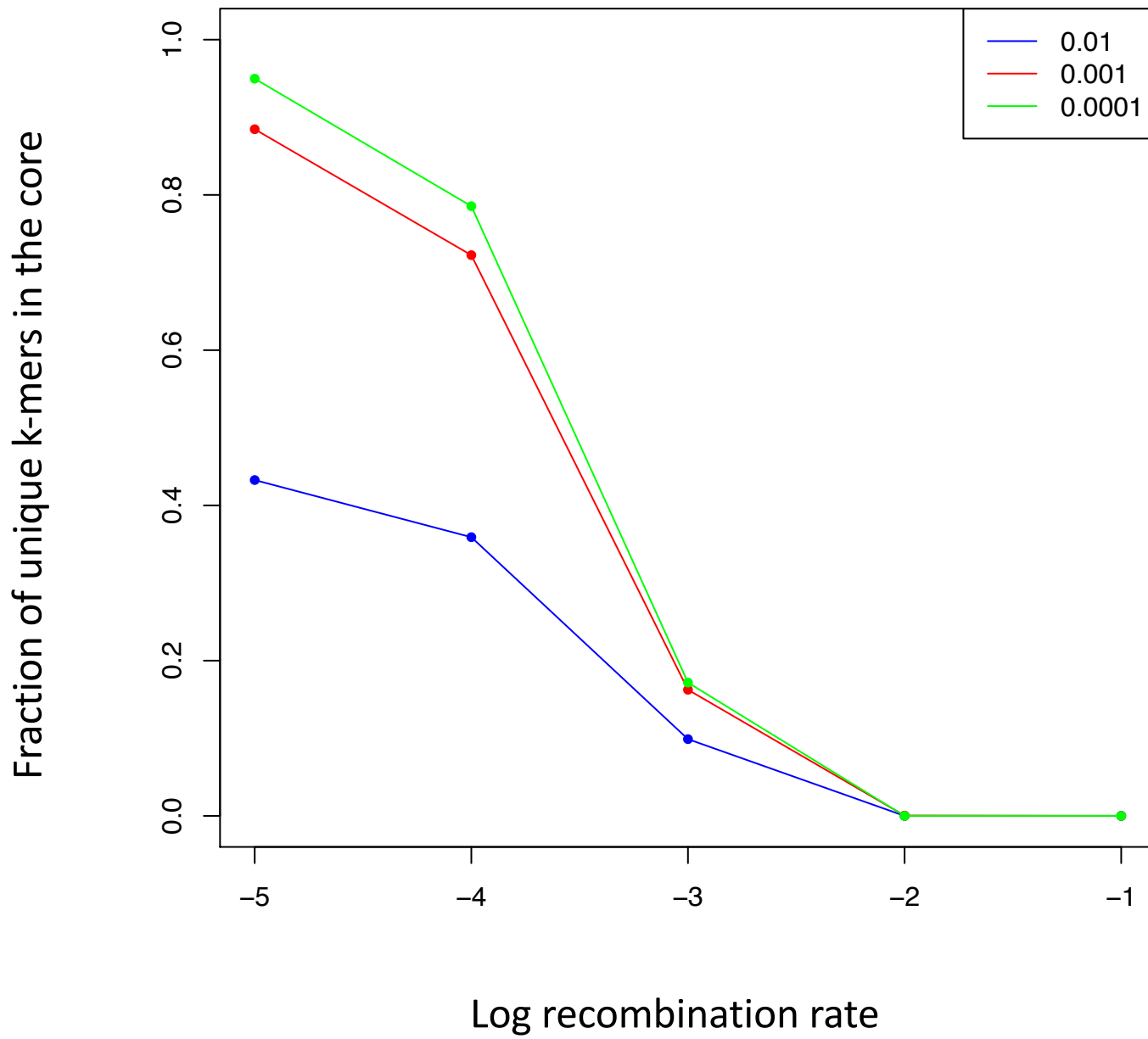
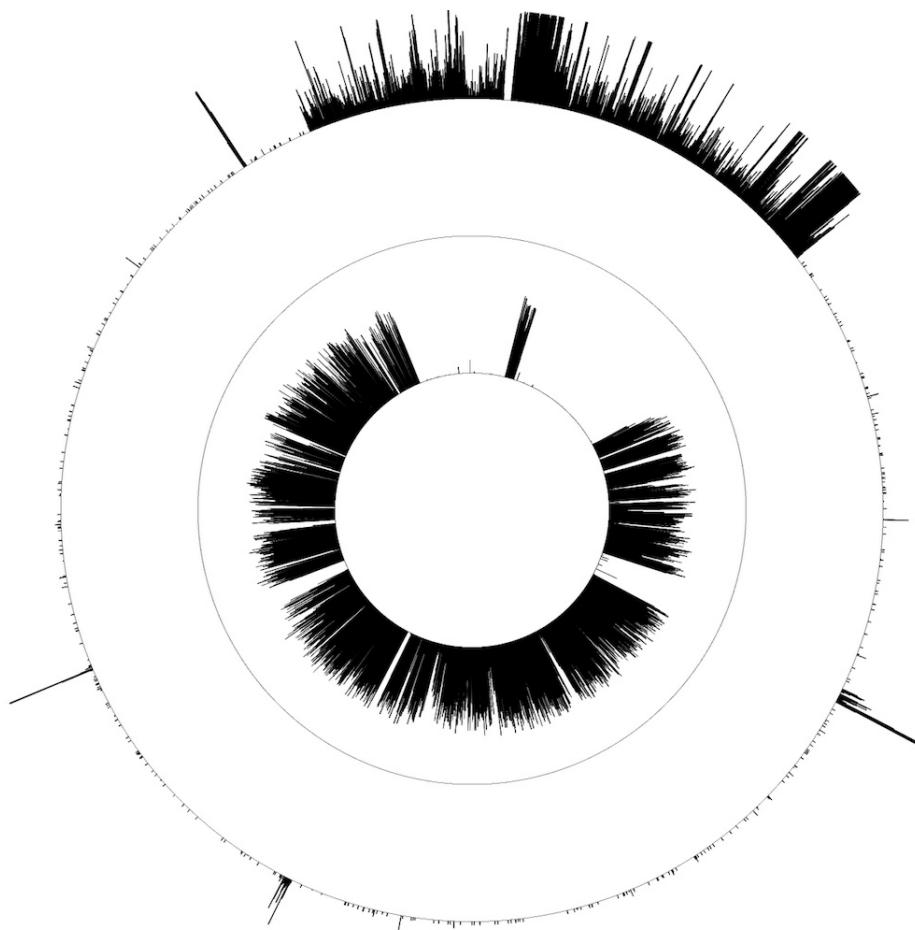
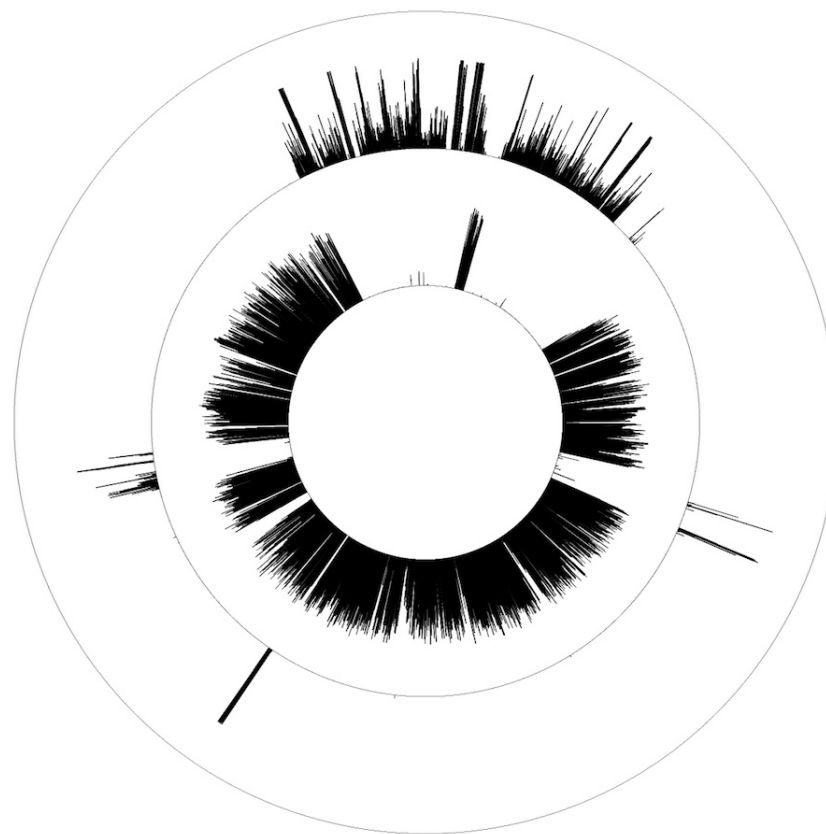


Figure 7



SaCOL



SaT0131

Figure 8

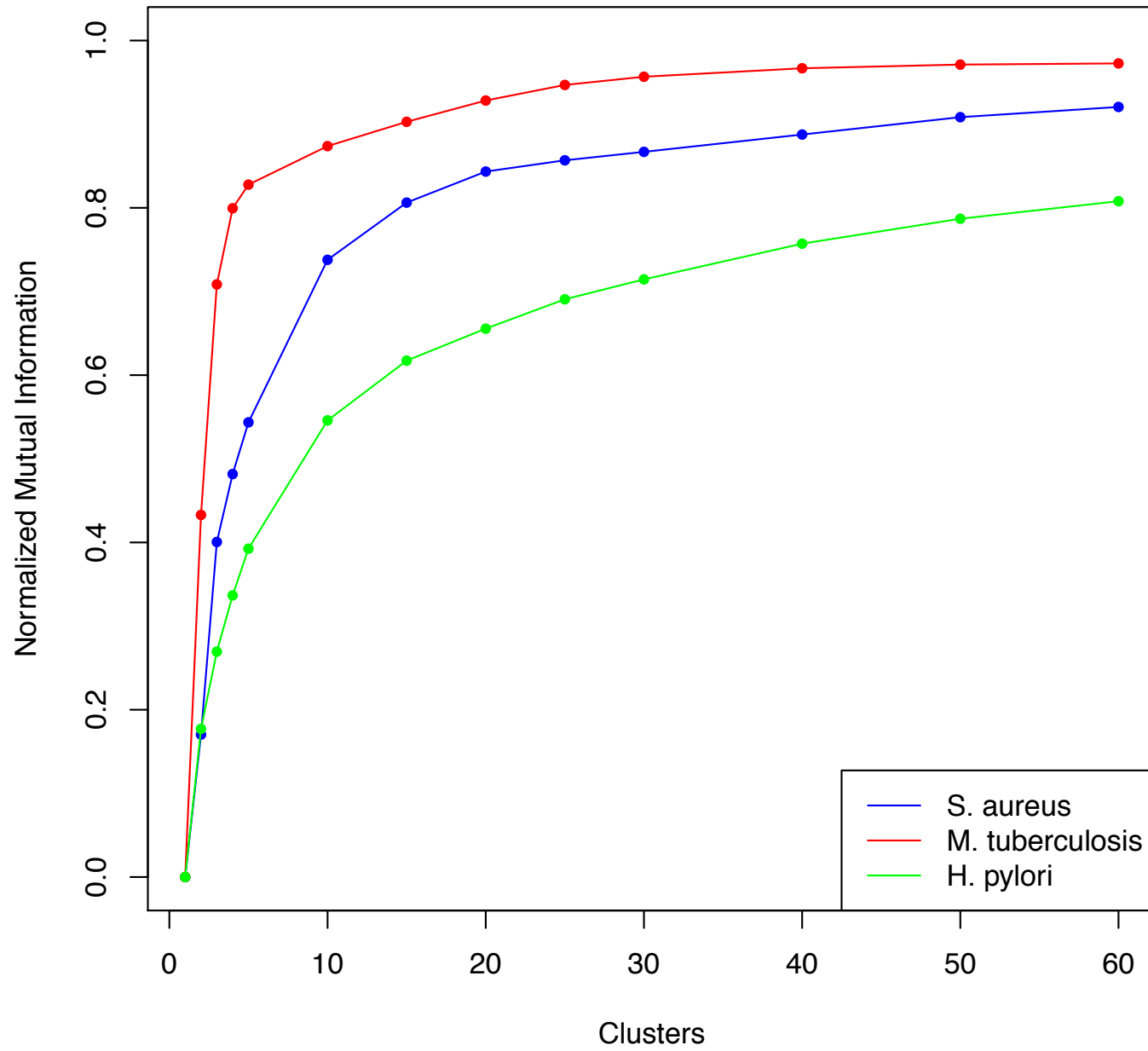


Figure 9