# What we talk about when we talk about species

*Running title: Compressing genomes to understand their evolution*

Apurva Narechania[1,2], Shyam Gopalakrishnan[2], M Thomas P Gilbert[2,3]

[1]Institute for Comparative Genomics, American Museum of Natural History, New York, NY, USA
[2]Center for Evolutionary Hologenomics, The Globe Institute, University of Copenhagen, Copenhagen, Denmark
[3]University Museum, NTNU, Trondheim, Norway
Corresponding authors: anarechania@amnh.org; tgilbert@sund.ku.dk

1    *Abstract*

2    Genome annotation, alignment, and phylogenetics are at the center of most work in

3    evolutionary genomics. These techniques function best when rooted in prior work. Genes are mined

4    from new genomes using evidence from old gene models. These genomes are aligned to well-worn

5    references to create matrices for tree reconstruction. And trees are often populated with well

6    characterized genomes to add context to the newly sequenced. Genome inference traces a line back

7    to model organisms, yoking the analysis of new genomes to layers of previous knowledge. We

8    instead highlight methods that use unannotated and unaligned sequence to understand the

9    information diversity of sequence ensembles. Any set of genomes can comprise our sequence

10   ensemble. In a pandemic context, a sequence ensemble might be clinically isolated strains from one

11   day. In a systematic context, a sequence ensemble could be the pangenome available for a clade.

12   The normal bioinformatics playbook would have us align. But we instead compress. A sequence

13   ensemble that compresses easily contains lower information diversity. For pandemics, we can use

14   curves of information diversity to trace genomic novelty and monitor selective sweeps in new

15   strains. For systematics, we can calculate compressibility quickly across all known bacterial taxa,

16   leveling the criteria for species across clades.  If we tolerate data loss, we can go one step further

17   and capture structural evolution as we compress. Our approach sacrifices a lot. We skip many of the

18   products of modern bioinformatics like variation anchored to known genes or genome alignment to

19   prescribed references or pangenome graphs. But we gain speed, breadth, and the ability to respond

20   to novelty.

21

22   *Introduction (The problem)*

23   Compression encodes information into reduced representations. Whether bits are

24   eliminated through statistical redundancy (lossless compression), or shed entirely (lossy

25   compression), compressed data always has a smaller footprint than the original. The act of

26   compression – its difficulty or ease – communicates information about the original data source.

27    Highly redundant data with many common patterns will compress easily. In contrast, novelty or

28    surprise with little repeated context is difficult to compress. Evolution creates ensembles of

29    sequence. These ensembles can be represented as pangenomes. Pangenomes are compressible

30    entities, but how compressible depends on evolutionary strategy.

31        Genomics is a retrospective field. Existing bioinformatic techniques often model new

32    genomes on sequences annotated in the past[1]. Alignment to these reference genomes

33    circumscribes our knowledge of diversity. Large swaths of the tree of life are presumably

34    unknown[2]. For example, much of the sequence from environmental samples passes through

35    annotation filters as undefined[3]. In a read streaming era[4], we need forward looking techniques

36    that flag genomic novelty by dispensing with references, annotation or alignment. Standard

37    methods are ill-equipped for these volumes. New species are not easily caught in the sparse web

38    of the known.

39        As genomics has swept through biology, systematics has come to favor molecular

40    character sets to help delimit species boundaries[5,6]. While morphology is still important, and

41    holdouts have been more than vocal[7], phylogenomics has more recently carried the day.

42    Phylogenomics extends the handful of marker genes that were the foundation of early molecular

43    systematics to matrices that concatenate thousands of orthologous genes[8]. This character

44    explosion has been a boon to systematics, but annotation is still anchored to the known.

45        The thousands of orthologous genes found in phylogenomic datasets are rarely evenly

46    distributed among the genomes that describe a species[9]. The complete set of these genes is one

47    definition of the pangenome, and its complexity was originally defined as the rate of gene

48    accumulation with newly sequenced genomes[10]. Genes found universally comprise a genomic

49    core and are considered indispensable for core species functions. Genes found sporadically may

50    contribute to strain success in particular niches but may not be essential to their overall biology.

51    The ratio of core genomes to accessory genomes informs genome fluidity[11]. Species whose

52     genomes are mostly core have closed, less fluid pangenomes. Species with a large fraction of

53     accessory genes are considered open and more fluid.

54          This gene-centric view of orthologs is blind to the diversity in the non-coding genome[12].

55     Whole genome alignment to annotated, chromosomal references[13,14] makes variation in non-

56     coding genome accessible but again circumscribes its characterization. If all we know is a linear

57     reference on a single coordinate system, our understanding of the non-coding genome will be

58     limited to what will stick.

59          More recent pangenome methods attempt to enhance the reference by conveying it as a

60     graph[15]. For example, a species graph through elements of the genomic core would collapse into

61     a single consensus, punctuated by bubbles that code small scale variation like single nucleotide

62     polymorphism and small insertion/deletion elements. In contrast, the accessory genome forks the

63     pangenome graph along entirely disparate paths. Graph-based methods attempt to incorporate

64     nuance and novelty into a more complex reference structure. But the game is still the same: new

65     data is aligned to a set of old genomes bound together into a complex, branching network.

66          Is there another way? Can we measure some other property of whole genomes that isn't

67     contingent on their alignment? Can we de-center the gene so we aren't limited to the protein

68     coding genome? Can we dispense with phylogenomics so we aren't spending CPU years

69     deciphering a bifurcating set of species relationships that convey a mere shadow of a more

70     reticulate truth[16]?

71          Here, we propose several new information theoretic techniques that reimagine genomes

72     as ensembles of information, containers subject to compression. This view of genomic

73     information does not require annotation. Because we aren't concerned with genes or the

74     contiguous arrangements of genomic elements, we also forgo alignment. We instead describe

75     pangenomes with summary statistics of string-based intersections. In this article, we argue that

76     compression can enhance existing comparative genomic strategies, highlight structural evolution

77    through controlled information loss, and democratize the bacterial species question by applying a

78    uniform mathematics across the Linnean taxonomy.

79

80    *The toolkit (entropy)*

81        Our approach is guided by two foundational concepts at the very root of information

82    theory: entropy and relative entropy. Both ideas rely heavily on Claude Shannon's seminal ideas

83    on information introduced in "A Mathematical Theory of Communication", the founding

84    document of information theory[17]. Information is data that reduces uncertainty. Shannon's

85    original formulation resembled the thermodynamic construction of entropy devised for statistical

86    mechanics[18]. We measure information entropy as

87
$$H = -\sum_{i=1}^{N} p_i \ln p_i$$

88    where $N$ is the set of all possible states, $i$, and $p_i$ is the probability of the $i$th state. This expression

89    quantifies data into bits (base two logarithm) or nats (natural log). The bit is the most irreducible

90    unit of information. A bit is gained when a binary variable is assigned either a 1 or a 0.

91        In genomics, our data comes in sequences. We can measure the entropy of sequences by

92    digesting into substrings of specific size. In the bioinformatics literature, substrings of biological

93    readouts (DNA, RNA, protein) are called k-mers. In a comparative setting, we're most interested

94    in the entropy of a group of sequences, or sequence ensemble[19,20]. For genome sequence

95    ensembles, alignment has been the tool of choice. But alignment is computationally arduous and

96    breaks down with evolutionary distance. Fields as diverse as linguistics[21], neurobiology[22], and

97    statistical mechanics[23] have successfully employed entropy to quantify ensemble complexity. In

98    each of these fields, researchers code a linear string of observations and divide into

99    subsequences, calculating the entropy  of each set across the ensemble. In Figure 1, we show

100    how the entropy of genome sequence (e.g. DNA/RNA) typically increases with increasing

101    subsequence size. This is a block entropy curve.

102          Block entropy curves contain information about the complexity of the ensemble[19].

103    Systems with more ensemble structure – repeated elements across sequences – will peak at lower

104    entropy. More novelty across sequences yields

105    higher entropy. In genomics, closed pangenomes,

106    with a large core shared across all species

107    genomes, have low entropy. Auxiliary genes

108    unique to subsets of genomes add entropy to the

109    ensemble system. The uneven distribution of

110    these elements is the hallmark of an open

111    pangenome. But to measure complexity we don't

112    need the annotated and aligned genes. Signal is
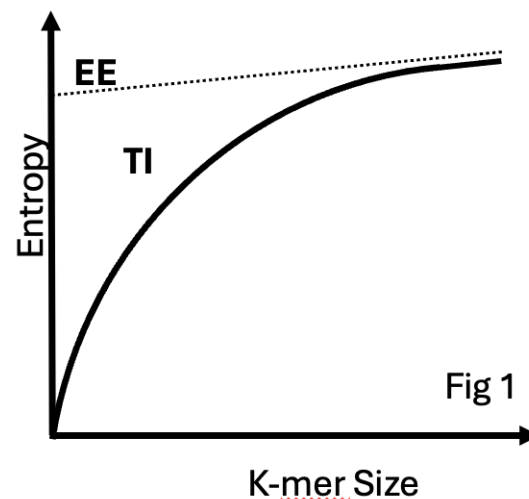
113    preserved in unaligned and unannotated k-mers.

**Figure 1. Block Entropy Curve. We show that entropy increases with k-mer size. We use this curve to calculate Excess Entropy (EE) and Transient Information (TI).**

114          Block entropy curves asymptote at the

115    minimum block size required to efficiently capture information across the sequences. We use

116    three quantities calculated from these curves to describe the complexity of a pangenome: source

117    entropy, excess entropy and transient information[19]. The source entropy (Hmu) is the irreducible

118    randomness that remains even as larger block sizes capture most ensemble correlations. Hmu is a

119    direct measure of randomness. Random distributions are hard to compress. A high source

120    entropy is associated with the accumulation of unevenly distributed accessory genes, resulting in

121    a more complex pangenome. The excess entropy (EE) is the non-random fraction of the total

122    information in the system. It's the information we model from redundancies across the ensemble.

123    Alignment is anchored to these same redundancies. In fact, alignment only works if enough of

124    these redundancies are spread across the query genomes. Finally, the transient information (TI)

125  measures how much information we must invest to learn Hmu and EE. In Figure 1 we show it as

126  the area between the block entropy curve and the line defining Hmu. Species with closed

127  pangenomes typically have a lower TI than those open to accumulating gene diversity. Closed

128  pangenomes with a large core set of genes compress at lower k-mer sizes, approaching their

129  Hmu quickly.

130

131  *More tools in the toolkit (the information bottleneck)*

132        Entropy is the workhorse of lossless compression. In fact, it defines lossless

133  compression's limit. We cannot compress any further than the entropy of the source. In our

134  context, the block entropy curve follows compression limits along a k-mer spectrum. Lossless

135  compression preserves all data, but sacrifice can bring evolution into relief by isolating patterns

136  from genomic noise. Using lossy compression, we can identify the core genome of any species

137  without alignment or annotation. Along the way, we unlock the homologous and non-

138  homologous recombination events that violate vertical signal.

139        To understand how we can detect structural evolution without annotation or alignment,

140  we leverage Shannon's ideas on lossy compression. Shannon based his theory in communication.

141  A sender passes a message to a receiver through a channel. The fundamental problem of

142  communication is reconstructing that message. Communication channels suffer distortion. Data

143  rarely reaches the receiver whole. Information entropy represents the limit on how efficiently a

144  message can be compressed in the noise-less ideal.

145        No channel is noise-less. Still, the distortion introduced by noisy channels does not doom

146  message passing. A sender can compensate for noise by encoding more information into a

147  message, or a receiver can tolerate some level of distortion while ascertaining a sender's core

148  meaning. Shannon formalized this concept as rate-distortion theory[17]. On the sender side, the rate

149  is measured as bits of information per symbol. The sender's message is distorted as it passes

150 through a channel. The sender's rate and the receiver's distortion are inversely related. The

151 function describing the two variables for any given channel informs lossy compression. How

152 much information loss can we tolerate in reconstructing a sender's message?

153      This idea is central not only to information theory and lossy coding, but also to modern

154 machine learning methods that use variational autoencoders to populate the compressive layers

155 of a neural net[24,25]. The two key questions are 1. How well does a dataset compress, and 2. How

156 much data can we afford to lose?

157      We use these concepts to further understand pangenome complexity. Imagine the

158 compression regime in Figure 2. A set of genomes comprising a sequence ensemble are digested

159 into k-mers and compressed into a set number of clusters. This compression is analogous to a

160 communication channel. The more clusters we model, the higher the rate, and the lower the

distortion. With fewer clusters, we force the k-mers through a narrower channel and suffer more distortion.
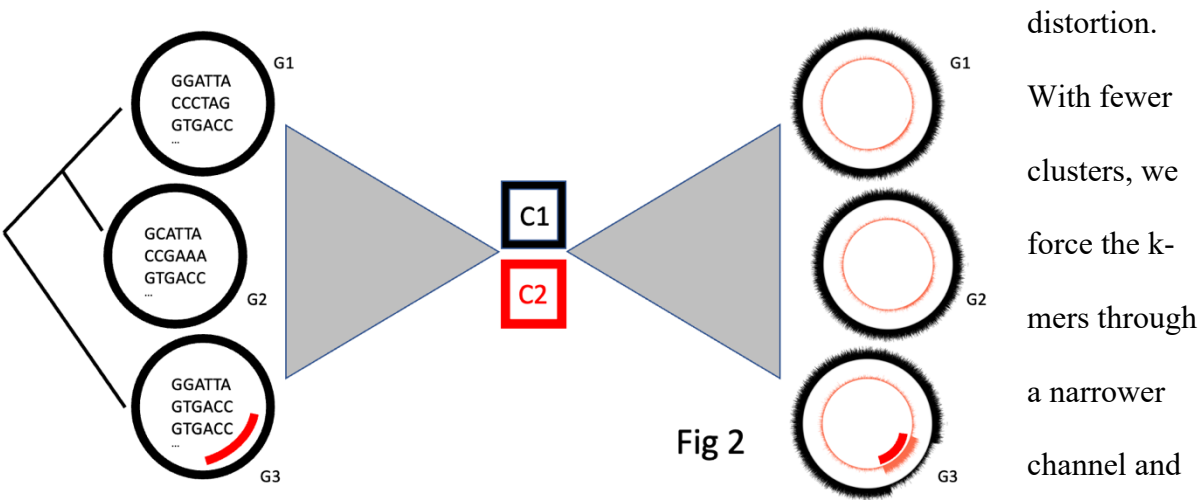


**Figure 2. The information bottleneck. As k-mers from our input genomes are compressed into a narrow channel, patterns of structural evolution emerge from the resulting clusters.**

170 If we hold the channel constant and model the same number of clusters across species

171 ensembles, open pangenomes will suffer more data loss than their closed counterparts. Open

172 pangenomes have more complex information to communicate. We employ the information

173 bottleneck[26], an idea first proposed in the Natural Language Processing literature, to measure

174 loss in our compression framework. Moreover, the clusters we glean from the information

175   bottleneck comprise a model of structural evolution. The largest cluster usually represents the

176   core. K-mers from recombination regions populate the others. The act of compression therefore

177   deconstructs real biological events without the need to align genomes, build sets of orthologs, or

178   calculate any trees.

179

180   *Even more tools in the toolkit (relative entropy)*

181          Block entropy curves measure the compressibility of any sequence ensemble. The

182   bottleneck compresses information into clusters that communicate only the most salient bits.

183   Compressibility is directly related to evolutionary strategy. Some species are open to genomic

184   input, others have narrower, closed pangenomes. But as we've described it here, entropy treats

185   the entire pangenome distribution as a single entity. This allows us to measure overall

186   complexity, but doesn't account for each genome's departure from that distribution. Relative

187   entropy, a measure of how one distribution (any given single genome) diverges from the overall

188   distribution (our pangenome) adds nuance to our approach. Summed across all genomes, the

189   relative entropy gives another, complementary angle on pangenome complexity and

190   compression.

191          To formalize this concept, we turn to bedrock principals in ecology. Ecology has an

192   extensive history of incorporating ideas from information theory and compression[27]. The

193   Shannon Index has long been used to combine the effects of species richness, the absolute

194   number of unique species in an environment, and species evenness, the relative abundances of

195   those species[28]. But for ecologists the core equation is more general than Shannon entropy.

196   Ecological datasets span many types of environments. Comparing diversity across those

197   environments is crucial. Hill introduced the effective number of species as an intuitive solution[29].

198   The effective number of species of order $q$ is given as

199
$$D_q = \left( \sum_{i=1}^{N} p_i^q \right)^{1/(1-q)}$$

200 where $p_i$ is the frequency of a particular species $i$, and $N$, the total number of unique species.

201 Sweeping through the parameter $q$ controls the metric's responsiveness to rare ($q = 0$) or

202 common species ($q = 2$ or more). At $q = 0$, the expression reduces to species richness, and at $q =$

203 2, the expression expands into the Simpson Index. But the sweet spot is at $q = 1$. The limit of this

204 equation as q approaches 1 is Shannon's information entropy (or the Shannon Index if overheard

205 in an ecology department). This transformation connects ideas from mathematical ecology to

206 information theory. The exponent of Shannon entropy yields the Hill number at q = 1, or the

207 effective number of species:

208
$$D_1 = exp\left( -\sum_{i=1}^{N} p_i ln\, p_i \right)$$

209 More diverse samples have higher Hill numbers. Hill numbers convey species diversity as an

210 intuitive number. Because of its connection to information theory, Hill numbers are not the

211 exclusive domain of ecologists. In Natural Language Processing, perplexity[30] is used to measure

212 how well a language model can predict a string of text. Perplexity is the effective number of

213 words in a library. Perplexity and Hill numbers draw from the same mathematical toolkit. This

214 toolkit's simplicity allows for easy comparisons between entirely different experiments. But the

215 expression collapses each experiment's observations into a single distribution.

216     We can enrich Hill numbers by extending beyond species measured as single variable

217 distributions. To this point, we've defined what an ecologist would term alpha diversity[31], or the

218 diversity of species in any one sample. One sample usually doesn't cut it. Ecologists sample

219 multiple transects from their environment of interest. Sampling introduces several opportunities.

220 First, the degree of sample overlap is a potential gauge of efficacy. Second, sample diversity

221 yields insight into the overall, hypothetical, unapproachable diversity of the system, or the

222   gamma diversity. If samples are highly diverse, ascertaining the diversity of the target

223   environment may require more samples to be taken. If gamma diversity is too high, no sampling

224   scheme may be enough. Beta diversity measures the degree of overlap between samples[32,33].

225   Grounding the concept in information theory, we extend the Hill number of species into a Hill

226   number of samples. The following expression yields the effective number of samples:

227
$$D_\beta = exp \left( \sum_{S=1}^{M} w_s \sum_{i=1}^{N} p_{si} ln \frac{p_{si}}{p_i} \right)$$

228   This equation incorporates the Kullback-Leibler divergence or relative entropy, a formulation as

229   frequently used as entropy in the information theory literature[34]. The relative entropy measures

230   the divergence of any one genome's k-mer distribution against the k-mer distribution of the

231   entire pangenome. Here, $N$ is the number of unique species, $M$, the number of samples, $p_{si}$ is the

232   frequency of species $i$ in sample $s$, $p_i$ is the frequency of species $i$ across all samples, and $w_s$

233   weighs all observations in sample $s$ relative to all individuals collected in the experiment.

234          The effective number of samples is another measure of compression. If species richness

235   and evenness is the same across all samples, the effective number of samples reduces to 1. If the

236   samples contain no species in common, or if species have wildly different occurrence counts, the

237   effective number of samples approaches the number of samples taken. In the first case, we have

238   perfect compression. In the latter, no compression at all.

239          We take this ecological concept and adapt it to genomics. Our goal is to calculate the

240   information diversity embedded in sequence ensembles. This requires a complete reframe.

241   Rather than species in a community (alpha diversity), we think k-mers in a genome. Rather than

242   transects in an environment (beta diversity), we think genomes in a pangenome. The shift is in

243   the container. Employed in this way, we recast Hill numbers as the effective number of genomes

244   or genome equivalents. We coin KHILL, an intuitive metric that quantifies the information space

245   of a pangenome, or the degree to which it will compress. We calculate KHILLs in a fraction of

246    the time it takes to annotate genomes, run alignments, and build the orthologs required to

247    compute pangenome fluidity.

248

249    *The toolkit applied*

250        Biological datasets are large and growing. Other fields also contend with large datasets,

251    and some have been grappling with them for decades longer.  For example, astronomers have big

252    data, perhaps the biggest data in the sciences[35]. Processing and saving all astronomic data is

253    impossible. Astronomers have known for years the importance of sensing data as it shines onto

254    their mirrors. Compression normally happens at the point of collection. We are quickly reaching

255    this point in biology.

256        Organized, collaborative genome sequencing projects began in earnest in the 1990s.

257    Starting then and through the first two decades of this century, genomic datasets were sacrosanct.

258    Groups held onto their data until every angle was exhausted. Though genomic data has always

259    been big data, generating it back then was costly. This is no longer the case. The price of genome

260    sequencing has seen steep decline. Storing this accumulating data has become nearly impossible.

261    Perhaps it is time to let go. With the information bottleneck, we tolerate controlled data loss.

262    New sequencing platforms emit data in nearly unending streams. Sensors are designed to glean

263    information from data streams in real time. There are sensors that detect change in acceleration

264    (engineers), in light (astronomers), in brain activity (doctors). Perhaps streams of biological

265    sequence can also be processed and discarded. Can sequence become a sensor?

266        Take for example SARS-CoV-2. Fifteen million SARS-CoV-2 genomes are now

267    available in various repositories around the world[36]. The state of the art in surveilling these

268    genomes as they accumulate in time and space is phylodynamic[37,38]. But phylodynamics is

269    retrospective. Investigators curate a fraction of the genomes available, compare them against an

270    even more rarefied set of references, and embed the new alongside the old either in phylogenetic

271 trees or networks. Alignment is the linchpin in this arrangement. Genome alignments feed tools

272 like Nextstrain[39], which employ Bayesian and likelihood phylogenetic approaches – some of the

273 most computationally costly algorithms in bioinformatics – to extend our view of SARS-CoV-2

274 biology slightly beyond the anointed references in a database.

275      We find this limiting. We can use KHILL to look forward, analyzing all the sequence

276 available to us outright[40]. Whether it's 15 million clinical genomes or streams of wastewater,

277 KHILL is capable of processing terabytes of streaming sequence and flagging the emergence of

278 new variants without relying on the references that confine biological novelty. KHILL can also

279 achieve rapid community analysis as exemplified in our study of the microbial shifts in the

280 making of cheese (ref), and the microbiome perturbations caused by broad spectrum antibiotics

281 (unpublished data). Whether it's a life threatening virus or the cheese you spread on crackers, we

282 use *all* sequence, not just the bits that will stick to existing references.

283      For SARS-CoV-2, we calculate one KHILL number per day along a pandemic time

284 course. Compiling these genome equivalents yields an information diversity curve through time.

285 KHILL increases as variants of concern ascend in a population mixing with a prior background.

286 KHILL decreases once these variants grow dominant and sweep away all other genomic

287 heterogeneity. In this way, we detect the emergence of concerning strains well before annotation

288 clearinghouses have blessed new database entries.

289      As a genomic measure of compression, KHILL also naturally lends itself to the analysis

290 of pangenomes. In fact, with SARS-CoV-2, we used KHILL as a rolling measure of pangenome

291 complexity. Because of their contracted timeline, pandemic genomes occupy a small information

292 space. The KHILL of all the millions of sequenced SARS-CoV-2 compresses to about 1.15

293 effective genomes. But KHILL is not restricted to any one biological scale. We can measure the

294 complexity of strains, species, genus, and collections at even higher taxonomic levels.

For example, we have used KHILL to calculate the pangenome complexity of all known bacterial species[41]. An analysis at this scale is impossible with current alignment-based bioinformatic techniques. But because KHILL is fast, we can compute genome equivalents for every species in the database. We couple this with metrics derived from block entropy curves (Hmu, EE and TI) to calculate the information space occupied by all known bacterial species. This information theoretic approach democratizes species classification, labeling each pangenome with a single number.

As we've defined it, KHILL species complexity mixes two separate phenomena. First, species definitions vary. The Linnaean taxonomy imposes a hierarchy on life, but this hierarchy is not uniformly applied. Species in one part of the taxonomic tree may not mean the same thing to its experts as species in another part of the tree. This is cultural. But it does influence the relative breadth of species buckets. We expect some variation in KHILL based just on these very human inconsistencies.

More interesting, however, is our second observation. Pangenome fluidity[11] has been shown to track with some gross aspects of bacteria phenotype[42]. For example, host-bound species accustomed to a uniform environment typically have less complex pangenomes. Cosmopolitan species occupying diverse niches tend towards more pangenome diversity. Obligate bacteria are less complex than their facultative counterparts. Non-motile organisms, less complex than those on the move. Complexity, in this case, was measured as pangenome fluidity. Pangenome fluidity is as near to measuring information-theoretic complexity as alignment-based techniques can get. We find that KHILL, a more direct, swifter measure of complexity, also corresponds to bacterial lifestyle. We see this borne out in KHILLs. For example, pathogens have significantly lower KHILL than mutualists. Challenging environments presumably encourage the accretion of pangenome complexity as species contend with

319  instability. Our compression based techniques squeeze this information from genomes without

320  the normal bioinformatics playbook.

321

322  *Challenges? In sacrifice there is clarity!*

323     Metrics based in compression can distort mechanism. KHILL increases with population

324  heterogeneity, as in the case of our SARS-CoV-2 populations. But it also increases with genetic

325  distance. This genetic diversity could be the result of environmental pressure, or it could simply

326  be lazy, inconsistent categorization. Because block entropy curves and KHILL dispense with

327  alignment, we also lose the ability to pinpoint change in genomic space.  In criticism of this

328  work, we've heard over and over how obscuring mechanism, sacrificing location, or conflating

329  biological forces is a weakness. But in a field saturated with sequence data, our approach allows

330  researchers to skim data streams without resorting to the heaviest, most cumbersome algorithms

331  in bioinformatics.

332     The idea of conflating signal is a hallmark of information-based approaches. Shannon's

333  communication problem is emblematic of this compromise. Distortion is inevitable as

334  information is relayed from sender to receiver. This concept has been used in everything[43] from

335  telecommunications, to thermodynamics, to data encoding in Natural Language Processing.

336  More complex data requires a broader channel to communicate. But sometimes we must

337  sacrifice nuance for meaning. In fact, compressing away the noise can sometimes distill signal.

338  In other words, conflation sometimes yields clarity.

339     We take this concept to genomics[44]. Mutation, homologous recombination, and

340  horizontal gene transfer all distort genomic signal. We can capture the degree of distortion by

341  measuring how difficult it is to compress strings (k-mers) from a set of genomes into a set

342  number of clusters. If the compression is easy, we need fewer clusters – a narrower channel – to

343  achieve communication at an acceptable level of distortion. But if the genomes are labile, we

344 need more clusters to communicate the added information diversity. The information

345 bottleneck[26] quantifies complexity.

346      The clusters that comprise our information channel, are datasets that sort meaning. Where

347 KHILL is a mark of compression, these clusters are actual compressed representations. We can

348 measure the fidelity of the original 'message' carried by the genomes relative to these

349 compressed representations. Clonal, tree-like, bifurcating species generally require fewer clusters

350 to model modes of genomic change. Recombinogenic species require more clusters to achieve

351 the same signal clarity.

352      Like KHILL, this approach conflates biological phenomena. Lossy compression through

353 the bottleneck does not distinguish between mutation and recombination. But for both KHILL

354 and the bottleneck, the compressibility of a set of genomes becomes a metric that can be used to

355 compare sets of species.

356      We began this essay bemoaning genomics as a retrospective enterprise. We believe

357 information theory allows us to shift our gaze forward. Eliminating references opens us to

358 novelty. De-centering the gene offers a new view of pangenome complexity. And eliminating

359 alignment boosts speed. Together these efficiencies recast sequencing as a sensor delimiting

360 change. We can sense change along a pandemic trajectory. We can predict bacterial lifestyle

361 from compression. And we can probe the unbalanced hierarchies of bacterial taxonomy.

362

363 **Data Accessibility**

364 All data discussed in this paper is freely available online at NCBI

365 (https://www.ncbi.nlm.nih.gov).

366

367 **Author Contributions**

368 AN conceived the project and drafted the original and final versions; SG co-wrote the

369 manuscript; MTPG co-wrote the manuscript.

370

371 **References**

372 1. Guigó, R. Genome annotation: From human genetics to biodiversity genomics. *Cell*

373 *Genomics* **3**, 100375 (2023).

374 2. Wu, D. *et al.* A phylogeny-driven genomic encyclopaedia of Bacteria and Archaea. *Nature*

375 **462**, 1056–1060 (2009).

376 3. Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome research. *BMC*

377 *Biol.* **17**, (2019).

378 4. Erlich, Y. A vision for ubiquitous sequencing. *Genome Res.* **25**, 1411–1416 (2015).

379 5. Baker, R. H. & DeSalle, R. Multiple Sources of Character Information and the Phylogeny of

380 Hawaiian Drosophilids. *Syst. Biol.* **46**, 654–673 (1997).

381 6. Rokas, A., Williams, B. L., King, N. & Carroll, S. B. Genome-scale approaches to resolving

382 incongruence in molecular phylogenies. *Nature* **425**, 798–804 (2003).

383 7. Neumann, J. S., Desalle, R., Narechania, A., Schierwater, B. & Tessler, M. Morphological

384 Characters Can Strongly Influence Early Animal Relationships Inferred from Phylogenomic

385 Data Sets. *Syst. Biol.* **70**, 360–375 (2021).

386 8. Zhu, Q. *et al.* Phylogenomics of 10,575 genomes reveals evolutionary proximity between

387 domains Bacteria and Archaea. *Nat. Commun.* **10**, (2019).

388 9. Shi, T. & Falkowski, P. G. Genome evolution in cyanobacteria: The stable core and the

389 variable shell. *Proc. Natl. Acad. Sci.* **105**, 2510–2515 (2008).

390 10. Tettelin, H. *et al.* Genome analysis of multiple pathogenic isolates of Streptococcus

391 agalactiae: Implications for the microbial "pan-genome". *Proc. Natl. Acad. Sci.* **102**, 13950–

392 13955 (2005).

393    11. Kislyuk, A. O., Haegeman, B., Bergman, N. H. & Weitz, J. S. Genomic fluidity: an

394         integrative view of gene diversity within microbial populations. *BMC Genomics* **12**, 32

395         (2011).

396    12. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the

397         human genome. *Nature* **489**, 57–74 (2012).

398    13. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler

399         transform. *Bioinformatics* **25**, 1754–1760 (2009).

400    14. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–

401         2079 (2009).

402    15. Eizenga, J. M. *et al.* Pangenome Graphs. *Annu. Rev. Genomics Hum. Genet.* **21**, 139–162

403         (2020).

404    16. Doolittle, W. F. Phylogenetic Classification and the Universal Tree. *Science* **284**, 2124–2128

405         (1999).

406    17. Shannon, C. E. A Mathematical Theory of Communication. *Bell Syst. Tech. J.* **27**, 379–423

407         (1948).

408    18. Boltzmann, L. On the Relation between the Second Law of the Mechanical Theory of Heat

409         and the Probability Calculus with respect to the Theorems of Thermal Equilibrium. *Wien.*

410         *Berichte* **76**, 373–435 (1877).

411    19. Crutchfield, J. P. & Feldman, D. P. Regularities unseen, randomness observed: Levels of

412         entropy convergence. *Chaos Interdiscip. J. Nonlinear Sci.* **13**, 25–54 (2003).

413    20. Bialek, W., Nemenman, I. & Tishby, N. Predictability, Complexity, and Learning. *Neural*

414         *Comput.* **13**, 2409–2463 (2001).

415    21. Bentz, C. & Alikaniotis, D. The word entropy of natural languages. Preprint at

416         https://doi.org/10.48550/ARXIV.1606.06996 (2016).

417    22. Bialek, W., Rieke, F., De Ruyter Van Steveninck, R. R. & Warland, D. Reading a Neural

418          Code. *Science* **252**, 1854–1857 (1991).

419    23. Jaynes, E. T. Information Theory and Statistical Mechanics. *Phys. Rev.* **106**, 620–630

420          (1957).

421    24. Kingma, D. P. & Welling, M. Auto-Encoding Variational Bayes. Preprint at

422          https://doi.org/10.48550/ARXIV.1312.6114 (2013).

423    25. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep

424          convolutional neural networks. *Commun. ACM* **60**, 84–90 (2017).

425    26. Tishby, N., Pereira, F. C. & Bialek, W. The information bottleneck method. (2000)

426          doi:10.48550/ARXIV.PHYSICS/0004057.

427    27. Chao, A., Wang, Y. T. & Jost, L. Entropy and the species accumulation curve: a novel

428          entropy estimator via discovery rates of new species. *Methods Ecol. Evol.* **4**, 1091–1100

429          (2013).

430    28. Jost, L. The Relation between Evenness and Diversity. *Diversity* **2**, 207–232 (2010).

431    29. Hill, M. O. Diversity and Evenness: A Unifying Notation and Its Consequences. *Ecology* **54**,

432          427–432 (1973).

433    30. Brown, P., Pietra, S. D., Pietra, V. D., Lai, J. & Mercer, R. An Estimate of an Upper Bound

434          for the Entropy of English. *CL* (1992).

435    31. Jost, L. PARTITIONING DIVERSITY INTO INDEPENDENT ALPHA AND BETA

436          COMPONENTS. *Ecology* **88**, 2427–2439 (2007).

437    32. Marcon, E., Hérault, B., Baraloto, C. & Lang, G. The decomposition of Shannon's entropy

438          and a confidence interval for beta diversity. *Oikos* **121**, 516–522 (2012).

439    33. Marcon, E., Scotti, I., Hérault, B., Rossi, V. & Lang, G. Generalization of the Partitioning of

440          Shannon Diversity. *PLoS ONE* **9**, e90289 (2014).

441  34. Kullback, S. & Leibler, R. A. On Information and Sufficiency. *Ann. Math. Stat.* **22**, 79–86

442      (1951).

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

444      (2015).

445  36. Shu, Y. & McCauley, J. GISAID: Global initiative on sharing all influenza data – from

446      vision to reality. *Eurosurveillance* **22**, (2017).

447  37. Grenfell, B. T. *et al.* Unifying the Epidemiological and Evolutionary Dynamics of

448      Pathogens. *Science* **303**, 327–332 (2004).

449  38. Volz, E. M., Koelle, K. & Bedford, T. Viral Phylodynamics. *PLoS Comput. Biol.* **9**,

450      e1002947 (2013).

451  39. Hadfield, J. *et al.* Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* **34**,

452      4121–4123 (2018).

453  40. Narechania, A. *et al.* Rapid SARS-CoV-2 surveillance using clinical, pooled, or wastewater

454      sequence as a sensor for population change. *Genome Res.* **34**, 1651–1660 (2024).

455  41. Narechania, A., Bobo, D., Gilbert, M. T. P. & Gopalkrishnan, S. The information content of

456      species: formal definitions of pangenome complexity track with bacterial lifestyle. Preprint

457      at https://doi.org/10.1101/2025.03.28.645969 (2025).

458  42. Dewar, A. E., Hao, C., Belcher, L. J., Ghoul, M. & West, S. A. Bacterial lifestyle shapes

459      pangenomes. *Proc. Natl. Acad. Sci.* **121**, e2320170121 (2024).

460  43. Noise and Distortion. in *Advanced Digital Signal Processing and Noise Reduction* 29–43

461      (Wiley, 2001). doi:10.1002/0470841621.ch2.

462  44. Narechania, A. *et al.* What Do We Gain When Tolerating Loss? The Information Bottleneck

463      Wrings Out Recombination. *Mol. Biol. Evol.* **42**, msaf029 (2025).

464