

1 Simulation-based approaches to characterize the effect of sequencing depth on the quantity and quality  
2 of metagenome-assembled genomes

3

4 Taylor Royalty<sup>a</sup>, Andrew D. Steen<sup>a\*</sup>

5 <sup>a</sup>Department of Earth and Planetary Sciences, University of Tennessee – Knoxville

6 Running Title: Modeling metagenome-assembled genome properties

7

8 \* Corresponding: [asteen1@utk.edu](mailto:asteen1@utk.edu)

9

10 Abstract Word Count: 230

11 Importance Word Count: 108

12 Text Word Count: 4999

13

## 14 **Abstract**

15 We applied simulation-based approaches to characterize how sequencing depth influences the  
16 properties of genomes identified in metagenomes assembled from short read sequences. An initial  
17 analysis evaluated the quantity, completion, and contamination of metagenome-assembled genomes  
18 (MAGs) as a function of sequencing depth on four preexisting sequence read datasets taken from four  
19 environments: a maize soil, an estuarine sediment, the surface ocean, and the human gut. These were  
20 subsampled to varying degrees in order to simulate the effect of sequencing depth on MAG binning.  
21 The property, MAG quantity fit the Gompertz curve, which has been used to describe microbial growth  
22 curves. A second analysis explored the relationship between sequencing depth and the proportion of  
23 available metagenomic DNA sequenced during a sequencing experiment as a function of community  
24 richness, evenness, and genome size. Typical sequencing depths in published experiments (1 to 10 Gb)  
25 reached the point of diminishing returns for MAG creation. Simulations from the second analysis  
26 demonstrated that both community richness and evenness influenced the amount of sequencing  
27 required to sequence a metagenome to a target fraction of exhaustion. The most abundant genomes  
28 required comparable quantities of bases sequenced regardless of community evenness, while more  
29 uneven communities required considerably more sequences to fully sequence rarer members. Future  
30 whole-genome shotgun sequencing studies can use an approach comparable to the one described here  
31 to estimate the quantity of sequences required to achieve scientific objectives.

## 32 **Importance**

33 Short read sequencing with Illumina sequencing technology provides an accurate, high-throughput  
34 method for characterizing the metabolic potential of microbial communities. Short read sequences are  
35 assembled into metagenome-assembled genomes which allow metabolic processes influencing health,

36 agriculture, and biogeochemical cycles to be assigned to microbial clades. At present, no reliable  
37 guidelines exist to select sequencing depth as a function of experimental goals in metagenome-  
38 assembled genomes creation projects. The work presented here provides a framework for obtaining a  
39 constrained estimate on the number of short read sequences needed for sequencing microbial  
40 communities. Results suggested that both the microbe community richness and evenness influence the  
41 amount of sequencing in a predictable matter.

## 42 **Introduction**

43 The assembly of high-accuracy short read sequences into metagenome-assembled genomes (MAGs) is  
44 a recent approach to characterize microbial metabolisms within complex communities (1). The recent  
45 creation of ~8,000 MAGs from largely uncultured organisms across the tree of life (2), the spatial  
46 characterization of microbial metabolisms and ecology across Earth's oceans (3), and the  
47 characterization of the potential impact that fermentation-based microbial metabolisms have on  
48 biogeochemical cycling in subsurface sediment environments (4) provide a few examples of how  
49 MAGs helped constrain the relationships between microbial ecology, microbial metabolisms, and  
50 biogeochemistry. At present, there is little information to guide how much sequencing is appropriate  
51 for metagenomic shotgun sequencing experiments (5). For the year 2017, estimates compiled by  
52 Quince et al. (5) suggest that up till now, metagenomic shotgun sequencing experiments usually  
53 sequence between 1 Gb and 10 Gb DNA nucleotides. Nonetheless, more guidance is necessary for  
54 selecting an appropriate metagenomic shotgun sequencing depth for one's experimental question which  
55 balances the maximization of information and minimization of cost.

56 Illumina sequencing technology is currently the most popular platform to generate  
57 metagenomic shotgun sequences (5). Here we present two distinct analyses which constrain the  
58 relationship between the quantity of Illumina metagenomic shotgun sequences and the quantity and  
59 quality of retrieved MAGs. First, we performed *in silico* experiments simulating the effect of how  
60 sequencing depth on Illumina sequence read datasets impacted the retrieved MAG properties for these  
61 datasets. Second, we applied a theoretical model and numerical simulations to estimate the minimum  
62 sequencing depth needed to sequence a metagenome to a target fraction of exhaustion. The work  
63 presented here illustrates how community evenness and richness control the sequencing depth  
64 necessary to sequence a metagenome to a target fraction of exhaustion. These patterns can be used to

65 guide sequencing depth decisions for future sequencing efforts in which MAG creation is a primary  
66 goal.

## 67 **Results**

### 68 MAG ASSEMBLY AS A FUNCTION OF SEQUENCING DEPTH IN EXISTING METAGENOMIC DATASETS

69 The number of “effective MAGs” (equivalents to 100%-complete MAGs, as defined in the  
70 Methods section) as a function of high quality bases empirically fit the Gompertz equation (equation 1;  
71 Fig 1B; parameters in Table 1). For each environment, the data fit the Gompertz equation better than a  
72 linear least-squares fit based on Akaike Information Criterion (AIC) (6). This equation is formulated  
73 for applications with microbial growth curves, such that the parameters  $A$ ,  $\mu$ , and  $\lambda$  correspond to  
74 maximum cell density, growth rate, and lag time (Fig 1A). Here,  $A$ ,  $\mu$ , and  $\lambda$  correspond to the  
75 maximum number of effective MAGs assembled with the pipeline, the maximum rate which effective  
76 MAGs form as with more sequencing, and the “lag bases,” or the bases which must be sequenced prior  
77 to rapid retrieval in effective MAGs. For the estuary, maize, and human gut datasets, MAG yield began  
78 to asymptote at higher sequencing depths, which indicates that further sequencing would yield  
79 diminishing returns with our pipeline. The Tara Ocean dataset followed a similar pattern at <25 Gb.  
80 However, when the number of sequenced bases was >25Gb, the number of effective MAGs decreased  
81 and became insensitive to sequencing depth. Since we have expressed MAG creation in terms of  
82 effective MAGs, the actual number of MAGs created in each example was considerably higher.

83 Mean MAG completeness also increased towards an asymptote with increasing sequencing  
84 depth (Fig 1C). Completeness was highest for the human gut dataset, with a maximum of 23.9%, and  
85 increased continuously as sequencing depth increased. The mean MAG completeness reached an  
86 asymptote of ~10-15% for the other three datasets with sequenced bases >10 Gb. Note that when >10  
87 Gb were sequenced, the number of effective MAGs created still increased as new sequences were

88 added. For all datasets, mean MAG contamination was <2% (Fig 1D) and did not depend strongly on  
89 sequencing depth.

## 90 SIMULATION EXPERIMENTS

91 Using equation 7, we calculated the number of  $k$ -length sequence reads required to sequence all  
92 unique DNA sequences of length,  $k$  ( $k$ -mers), in four hypothetical metagenomes. Three of the  
93 community structures are ecologically unrealistic but represented a community in which taxa are  
94 distributed perfectly evenly, highly unevenly, and at an intermediate level of evenness (Fig 2A-C). The  
95 fourth community structure, which is lognormally distributed, is ecologically realistic (Fig 2D; (7, 8)).  
96 The expectation value of the log number of sequences required to fully sequence metagenomes of those  
97 hypothetical communities was linear with respect to log-transformed size of the metagenome (i.e.,  
98 number of unique  $k$ -mers in the population, approximate number of unique base pairs in a  
99 metagenome); this suggests a power-law relationship between metagenome size and expectation value  
100 of sequence reads required to sequence the metagenome to exhaustion (Fig 2E). For all community  
101 structures, the slope of the relationship between log-transformed sequenced reads and log-transformed  
102 unique number of sequenced reads was within 1% of 1.06. The structure of the population strongly  
103 influenced the number of reads required such that more even community structures required far fewer  
104 reads than less even structures.

105 As equation 7 only estimates the number of reads to sequence a metagenome to exhaustion, we  
106 used a numerical simulation to estimate the number of  $k$ -sized reads to sequence a metagenome to a  
107 target fraction of exhaustion. Numerical simulation results predicted the same number of sequences  
108 reads to sequence 100% of a given metagenome as the numerically integrated expected sequences from  
109 equation 7 (Fig 3); this supported the use of this simulation. The log-transformation of both total  
110 unique  $k$ -sized reads ( $|K_{MG}|$  and sequenced reads showed a linear response for all target fractions and

111 all community structures. The amount of sequences required to achieve a given target of  $|K_{MG}|$  was  
112 variable for the different communities shown in Fig 2A. For instance, the lognormally-distributed  
113 community required the most amount of sequencing to sequence a metagenome to a target fraction of  
114 exhaustion but required similar amount of sequencing to sequence the metagenome to a target fraction  
115 of 50% as the other communities.

116 We applied the simulation to semi-quantitatively demonstrate the effect that community  
117 evenness has on the number of reads required to sequence a community to a target fraction of  
118 completion. These communities ranged from perfectly even ( $a=0$ , eq. 9) to more uneven ( $a = 0.02$ , Fig  
119 4A). Evenness was quantified using the Pielou evenness index, which expresses Shannon diversity  
120 relative to the diversity of a perfectly even community (9). Computational limits precluded simulating  
121 communities with Pielou evenness less than 0.977 given the richness and size of genomes within the  
122 communities. The number of sequence reads required to sequence genomes to a target fraction of  
123 completion depended strongly on both the evenness and the target fraction of completion (Fig 4B).  
124 Again, more even communities required more sequence reads than less even communities. The strength  
125 of this relationship also depended on the target fraction of completion. A community with Pielou  
126 evenness of 0.97 required 3 orders of magnitude more sequence reads to sequence a metagenome to a  
127 target fraction of exhaustion than a perfectly even community while the same community only required  
128 about 42% more reads to sequence 50% of the metagenome.

129 The minimum number of sequence reads required to sequence a microbe genome given a  
130 combination of target fraction, genome size, and fraction of the metagenome community was modeled  
131 with a generalized additive model. The smooth dimensions for target fraction, genome size, and  
132 fraction of the metagenome community was 7, 3, and 9, respectively, to achieve a normal distribution  
133 of residuals. To normalize for different sequence read length, sequence reads were converted to bases

134 and ranged from  $1 \times 10^7$  to  $1 \times 10^{13}$ . More bases were required to sequence microorganisms when 1) the  
135 genome was relatively rarer in the community, 2) to achieve better coverage of the genome, and 3)  
136 when the genome increased in size.

## 137 **Discussion**

138 We sought to establish evidence-based guidelines for selecting a sequencing depth during  
139 shotgun metagenomic sequencing experiments with the goal of creating MAGs of a given quantity and  
140 quality. Random subsamples of existing short read datasets, which were each individually assembled  
141 and binned, simulated the effect of creating MAGs from datasets of different sizes and environments.  
142 The datasets analyzed here are argued to be representative of both the order of magnitude of  
143 sequencing depth (1 to 10 Gb) (5) and the types of target environments microbial ecologists often  
144 investigate (10). A variety of software is available for all steps of MAG creation pipelines, and the  
145 quantity/quality of MAGs will depend on software selection, software configuration, and sequenced  
146 environment (5). Furthermore, it is best-practice to manually curate algorithmically-created MAG bins  
147 (11). We do not argue that the pipeline used here is objectively optimal for generating “true” MAGs  
148 (i.e., represent true genomes). Thus, MAG quantity was not directly reported but expressed as effective  
149 MAGs. The metric, effective MAGs, represents the integrated completeness (12) divided by 100 for  
150 MAGs retrieved with a taxonomic rank of at least phylum. In effect, effective MAGs represents  
151 phylogenetic signal, as defined by the presence of marker genes in assembled contigs (necessary for  
152 constructing MAGs). Thus, increases in effective MAGs should scale proportionally with increases in  
153 the quantity of true MAGs.

154 As sequencing depth increased, there was at first a “lag time” (more precisely a lag depth, or  
155 number of bases before effective MAGs began to increase) followed by a rapid increase in effective  
156 MAG quantity, and then diminishing returns at higher sequencing depths. Previous investigators



157 modeled the response of 16S RNA gene (13–15), Hill’s number diversity (16), taxon-resolved  
158 abundance (17), and gene abundance (17) as a function of sequencing depth using rarefaction curves,  
159 or collectors curves. The effective number of MAGs created did not match a traditional collector’s  
160 curve, which does not contain any initial lag. The Gompertz function, conversely, fit the data well,  
161 suggesting that MAG construction as a function of sequencing depth behaves similarly to microbial  
162 growth in a constrained medium, in concept if not in precise mechanism. The Gompertz function is  
163 defined in terms of three parameters,  $A$ ,  $\mu$ , and  $\lambda$ . These parameters correspond to the maximum  
164 effective MAGs at infinite sequencing depth ( $A$ ), maximum rate that effective MAGs increased with  
165 increases in sequencing depth ( $\mu$ ), and a minimum threshold of sequencing necessary prior to rapid  
166 effective MAGs retrieval ( $\lambda$ ) (Fig 1A). The Gompertz equation achieves the same asymptotic behavior  
167 of conventional rarefaction models while also modeling the apparent lag ( $\lambda$ ) in effective MAGs  
168 observed during this work (Fig 1B).

169 The four environments analyzed demonstrated different responses to increases in sequencing  
170 depth. Specifically, the predicted maximum effective MAGs varied from ~17 to ~97, the predicted  
171 maximum rate that effective MAGs increased varied from ~1.4 to ~5.8, and the minimum threshold of  
172 sequencing necessary prior to seeing effective MAGs varied from ~0.6 to ~6.7. The Tara Ocean  
173 dataset, where effective MAGs decreased at sequencing depth >20 Gbp, was an exception. We  
174 speculate that our choice of pipeline, and specifically the fact that we discarded contigs <3kb, caused  
175 poor performance at higher sequencing depth for the Tara Ocean dataset.

176 As mean MAG completeness converged to an asymptote considerably less than 100% (Fig 1B),  
177 MAG yields (Table 1) were close to 100%. This suggests the maximum effective MAGs ( $A$ ) likely  
178 represents sequence reads associated with abundant MAGs. Thus, we asked how much sequencing was  
179 necessary to sequence a community to exhaustion. The expected number of sequence reads required to

180 sequence an entire metagenome was estimated using equation 7 for four hypothetical communities (Fig  
181 2A-D). The total unique  $k$ -sized reads (i.e., richness) and community structure influenced how much  
182 sequencing is necessary to sequence an entire metagenome (Fig 2E). For a given community structure,  
183 increases in community richness lead to linear increases the sequencing depth necessary to exhaust the  
184 metagenome. All regressions had similar slopes, indicating that community structure did not exert a  
185 major influence on that relationship. Interestingly, the sequencing depth necessary to sequence an  
186 entire metagenome depended strongly on the structure of the target microbial community (Fig 2E). As  
187 sequencing depth was log-transformed in Fig 2E, the differences in model intercepts indicate orders of  
188 magnitude differences in the necessary sequencing depth. The primary implication of Fig 2 is that the  
189 sequencing depth increased in a predictable trend in response to richness, regardless of the community  
190 structure.

191 One limitation to equation 7 is that it only provides an estimate of the sequencing depth  
192 required to sequence a metagenome to exhaustion. For practical applications, a continuous increase in  
193 sequencing depth eventually leads to diminishing returns in identifying unique sequence reads while  
194 also leading to a disproportional increase in monetary resources needed to find these unique sequence  
195 reads (18). Thus, it is desirable to constrain the fraction of unique sequence reads (e.g., 50%, 70%,  
196 90%, etc.) sequenced from a metagenome in relation to monetary investment necessary to achieve that  
197 fraction of a metagenome. Simulations show that as target metagenome completeness increases, the  
198 sequencing depth required increases dramatically (Fig 3). Simulation results were validated by  
199 comparing the sequencing depth necessary to sequence 100% of a metagenome with predictions from  
200 equation 7. While the numerical approach successfully reproduced and extended equation 7,  
201 communities with large values of richness ( $|K_{MG}| > 1 \times 10^8$ ) became computationally burdensome.  
202 Nonetheless, when the target fraction and community structures were held constant, the linear increase

203 in sequencing depth as a function of increased richness suggests linear regression may be sufficient to  
204 estimate sequencing depth for communities with large values of richness.

205 One observation from the numerical simulations was the impact that community structure had  
206 on the required depth of sequencing (Fig 2E and 3). Even communities required less sequencing to  
207 achieve a fraction of  $|K_{MG}|$ . Conceptually this makes sense, as abundant taxa (i.e., large  $n$  values in  
208 equation 3) should be sequenced more deeply compared to rarer taxa. To further explore the influence  
209 that community evenness had on required sequencing depth, communities with similar and more  
210 realistic lognormal structures (7, 16) at different levels of evenness were compared to one another (Fig  
211 4A). Decreasing evenness (increasing  $a$ ; equation 9) led to both increases in the sequencing depth  
212 required to sequence a given target fraction of  $|K_{MG}|$  (Fig 4B). For communities with more uneven  
213 species distributions, rarer community members required more sequencing. While only semi-  
214 quantitative, this analysis demonstrates that community evenness can have a significant impact on the  
215 sequencing depth necessary to characterize an entire community.

216 In practice, information about a target community structure may not be available for estimating  
217 sequencing depth. The spline model built here illustrates the minimum number of sequences necessary  
218 to sequence a given fraction of a target genome, assuming genome size and proportion that the genomic  
219 content represents in the community metagenome ( $G_{MG}$ ) (Fig 5). This proves useful for constraining  
220 the observed MAG properties from one's bioinformatic pipeline (e.g., Fig 1B-D) in the context of what  
221 proportion of a given microbe's metagenome ( $g_{MG}$ ; equation 4) has been sequenced to exhaustion. For  
222 example, taking the 5 Gb human gut dataset analyzed here (Table 2), if a microbe with a genome size  
223 of ~5 Mbp existed from this environment, then Fig 5C suggests that a 5 Mbp genome  
224 representing >10% of the whole metagenome ( $G_{MG}$ ; equation 5) will be sequenced to a minimum of  
225 50% to exhaustion. More so, one has constrained perspective of how a given genome may be

226 represented in the retrieved MAGs. Although the simple nature of sequencing a genome may not  
227 necessarily translate into the production of more MAGs, one can safely say that additional sequencing  
228 of that 5 Mbp genome which represents >10% of the community will not lead to the addition of more  
229 MAGs. More so, the bioinformatic pipeline would act as the limiting step (opposed to sequencing) in  
230 the production of MAGs.

## 231 **Materials and Methods**

### 232 SEQUENCE DATA SOURCES

233 All sequence data were downloaded from NCBI's Sequence Read Archive (SRA) using the SRA  
234 Toolkit (fastq-dump --split-files) (19). Exact duplicate reads for both forward and reverse reads were  
235 removed using PRINSEQ (-derep 1; v0.20.4) (20). All sequencing datasets were limited to Illumina  
236 shotgun metagenomic paired-end reads. Four datasets were analyzed for this analysis. The first dataset  
237 was from oceanic surface water collected at 5m depth in the Caribbean Sea as a part of the Tara Oceans  
238 expedition (21). The second dataset was from sediment from a depth of 8-10 cm below the surface  
239 (sulfate-rich zone) and collected at the White Oak River Estuary, Station H, North Carolina, USA (4).  
240 The third dataset was collected from maize soil (22). The last dataset was collected from human fecal  
241 samples and represented a human gut microbiome (23). All datasets analyzed in this study are  
242 summarized in Table 1.

### 243 MAG ASSEMBLY PIPELINE

244 The pipeline developed here followed similar pipelines described by other authors (3, 24). All sequence  
245 datasets were analyzed as follows. Trimmomatic (v0.36) (25) removed adapters as well as trimmed  
246 low-quality bases from the ends of individual reads. Read leading and trailing quality scores were  
247 required to be >3. The sliding window was set to 4 base pairs and filtered base pair windows with a

248 mean score <15. Quality controlled reads were assembled into contigs using MEGAHIT (v1.1.2; --  
249 presets meta-large) (26). Due to RAM limitations, assembled contigs <3000 bp in length were excluded  
250 from the analysis. Redundant contigs were removed using CD-HIT (v4.6.8; cd-hi-est -c 0.99 -n 10)  
251 (27). Similarity among the remaining contigs was further evaluated via intra-contig sequence  
252 alignments using Minimus2 (-D OVERLAP=100 MINID=95). The quality-controlled reads (i.e., after  
253 using Trimmomatic) were then mapped to the remaining contigs using Bowtie 2 (v2.3.3) (28) to  
254 generate a coverage score for individual contigs.

255 Resultant contigs were iteratively clustered into MAGs using the unsupervised clustering  
256 algorithm Binsanity (v0.2.6) (24). Similar to Tully et al. (3), six initial clustering iterations were  
257 performed with the parameter, *preference* (-p), set to -10 (iteration 1), -5 (iteration 2), -3 (iteration 3-6).  
258 Between iterations, a refinement step (Binsanity-refine) was performed on the putative MAGs with  
259 constant *preference* (-p) of -25. The refined putative MAGs were evaluated for contamination and  
260 completeness using the software CheckM (v1.0.6) (12), which uses HMMER (v3.1) and Prodigal  
261 (v2.6.3) (29). Contigs associated with putative MAGs meeting one of the following criteria: 1) had a  
262 completeness > 90% and contamination < 10%, 2) had a completeness > 80% and contamination < 5%,  
263 or 3) had a completeness > 50% and contamination < 5% were treated as high-quality. All other MAGs  
264 were considered low-quality MAGs. MAGs defined as high-quality were not modified any further.  
265 Contigs associated with the high-quality MAGs were not used in the subsequent reclustering and  
266 refinement steps. The contigs associated with low-quality MAGs were pooled together and reclustered  
267 during the next iteration of Binsanity clustering. After the sixth iteration, the remaining MAGs which  
268 did not fall into one of the three categories underwent additional refinement using Binsanity-refine.  
269 During this step, MAGs were iteratively refined with *preference* set to -10 (iteration 1), -3 (iteration 2),  
270 and -1 (iteration 3). Between each refinement step, metrics of contamination and completeness were

271 evaluated using CheckM. Again, MAGs which met the criteria of one of the high-quality categories  
272 described above were not further modified. The respective contigs to the putative MAGs were not used  
273 in proceeding refinement steps. After the last iteration of refinement, all MAGs were reevaluated for  
274 completeness and contamination as well as assigned a final taxonomic rank using CheckM.  
275 Completeness and contamination values for MAGs with the resolved taxonomic rank of phylum were  
276 integrated together. The integrated completeness was then divided by 100 to produce effective number  
277 of MAGs.

#### 278 SUBSAMPLING SEQUENCE READ DATASETS

279 The effect of decreased sequencing depth was simulated by subsampling the initial sequence read  
280 datasets described above. Downloaded sequence read datasets were randomly sampled at set fractions  
281 of 1%, 10%, 20%, 40%, 60%, 80%, 90%, 95%, and 100%. To account for variability in the reads  
282 sampled at a given fraction, each fraction was resampled, assembled, and binned in triplicate. All  
283 triplicates were analyzed using the MAG assembly pipeline described above.

#### 284 MODELING MAG RESPONSE TO SEQUENCING DEPTH

285 Effective MAGs as a function of sequencing depth was modeled for environmental sequence datasets  
286 using the Gompertz equation, as reformulated by Zweitering et al. (30) for use with microbial growth  
287 curves:

$$288 \text{ Effective MAGs} = A \times e^{-e^{-\frac{\mu \times e}{A}(\lambda - b) + 1}} \quad (1)$$

289 where  $A$ ,  $\mu$ , and  $\lambda$  are fit coefficients and  $b$  is high-quality bases. To assess the validity of this function,  
290 AIC (6) was calculated for all Gompertz equation fits and compared to AIC values for linear  
291 regressions models for same dataset.

#### 292 DEFINING THE MICROBIAL METAGENOME AND SEQUENCING PROBABILITY

293 Here we draw on set theory to provide a theoretical grounding for our *in silico* simulations described

294 below. The application of probability theory for predicting the expected number sequences to sequence  
295 a metagenome became founded by defining a metagenome as the set of available metagenomic DNA  
296 that can be sequenced in a sequencing experiment. Fig 6A-E provides a cartoon example illustrating the  
297 application of this set theory on a hypothetical microbial population,  $G$ .  $G$  is a community of genomes  
298 ( $g$ ) with finite abundances ( $n$ ). As the definition of microbial species is somewhat contentious (31),  $g$  is  
299 taken as the average genome for all individual genomes defined as a meeting some criteria defining a  
300 taxonomic rank. Thus, the richness ( $s$ ) of  $G$ , or the total number of  $g$ , depends on the definition of  $g$ . In  
301 the example  $G$  (Fig 6A-E),  $s=6$  and the total  $n=13$ . Thus,  $G$  can be represented as (Fig 6A):

$$302 \quad G = \{n_1g_1, n_2g_2 \dots n_s g_s | n \in N\} \quad (2)$$

303 where  $s$  is the total number of unique species within the community (richness). When characterizing  $G$   
304 via shotgun metagenomics, the  $i^{\text{th}}$  genome,  $g_i$ , can be sequenced at  $K$  unique sections given a  
305 characteristic read length,  $k$ , and average genome size,  $l$ , in number of base pairs (Fig 6B). Thus, the  
306 number of unique  $k$ -sized reads,  $K$ , associated with the  $i^{\text{th}}$  genome,  $g_i$ , within  $G$  is equal to:

$$307 \quad K_{g_i} = l(g_i) - k + 1 \quad (3)$$

308 From equation 3, the metagenome,  $g_{MG}$ , for  $g_i$  is defined as the set of all unique possible  $k$ -sized reads  
309 (Fig 6C) or:

$$310 \quad g_{MG,i} = \{g_{i,1,1+k}, g_{i,2,2+k} \dots, g_{i,Kg_i,Kg_i+k}\} \quad (4)$$

311 where the subscripts for  $g_i$  represent a given  $k$ -sized read spanning from an arbitrary starting base pair  
312 to the arbitrary starting base pair plus  $k$ . By substituting  $g_{MG,i}$  into all  $g$  for equation 2 (Fig 6D), the  
313 metagenome for a microbial community,  $G_{MG}$ , is derived to be:

$$314 \quad G_{MG} = \{n_1g_{MG,1}, n_2g_{MG,2} \dots n_s g_{MG,s} | n \in N\} \quad (5)$$

315 while the population of unique  $k$ -sized reads in the metagenome,  $G_{MG}$  (Fig 6E), is represented as:

$$316 \quad K_{MG} = \{g_{MG,1}, g_{MG,2} \dots g_{MG,s}\} \quad (6)$$

317 From equation 4, one can determine the cardinality, or the total number, of unique  $k$ -sized reads in  
318 associated with  $G_{MG}$  (expressed as  $|K_{MG}|$ ). When attempting to fully sequence  $G_{MG}$  using shotgun  
319 metagenomics, we assume that sampling events (sequence reads) are independent and are sampled with  
320 replacement. In fact, Illumina sequencing technology sequences reads in parallel via the individual  
321 DNA fragments binding to individual clusters. Furthermore, the fragmented DNA cannot be sequenced  
322 twice as the sequencing process is destructive (32). Nonetheless, the mass of DNA extracted from a  
323 target environment will represent a negligible fraction of the total DNA which exists in that  
324 environment. As the relative abundance of the  $k$ -sized reads in  $K_{MG}$  does not change when DNA is  
325 extracted from an environment, sampling events can be treated as independent and thus, DNA sampling  
326 reduces to sampling with replacement. If the proportion DNA mass extracted had a significant impact  
327 on the remaining mass of DNA in the environment, then one would be more suited to sequence all the  
328 DNA versus a smaller proportion of the DNA. The sequencer should have no impact on sampling  
329 assuming no sequencing errors due to misreading or spatial sampling issues (i.e., clonal density issues).  
330 Obviously, these issues do exist, but for the sake of a first order, general approximation, these biases  
331 can be ignored.

332 By making the above assumptions, the probability of sequencing all elements in  $G_{MG}$  reduces to  
333 a coupon collectors problem (33). Using the general functional form for calculating expected samples  
334 for sampling all unique elements in a set (equation 13b in 8), one can predict the number of sequences  
335 necessary to sequence all elements in  $K_{MG}$ , such that the expected number of sequences,  $E(G_{MG})$ , is:

$$336 \quad E(G_{MG}) = \int_0^{\infty} (1 - \prod_{j \in K_{MG}} (1 - e^{-p_j t})) dt \quad (7)$$

337 where  $j$  is a given element within  $K_{MG}$ ,  $t$  is the number of sampling events, and  $p_j$  is equal to the  
338 proportion of the  $j^{\text{th}}$   $k$ -sized read within a given population of  $k$ -sized reads.  $p_j$  can be expressed as  
339 follows:



$$340 \quad p_j = \frac{n_i \times j \in K_{MG}}{|G_{MG}|} \quad (8)$$

341 where  $n_i$  is the respective abundance for the species whose MAG contains the  $j^{\text{th}}$   $k$ -sized read within  
342  $K_{MG}$ , and  $|G_{MG}|$  is the cardinality of  $G_{MG}$ , or the total number of  $k$ -sized reads in the metagenome,  $G_{MG}$ .

### 343 MODELING EXPECTED SEQUENCES

344 Equation 7 provides an estimate for the total number of sequences to sequence all  $K_{MG}$ . The  
345 influence of increasing species richness (i.e.,  $s$  in equation 2) on the expected number of sequences was  
346 tested for four hypothetical communities. The first community had an even structure such that all the  
347 metagenomic DNA segments were equally distributed across all  $K_{MG}$ . In the second community, 90%  
348 of the metagenomic DNA segments were equally distributed in 50% of  $K_{MG}$ , and the remaining 10% of  
349 the metagenomic DNA segments were distributed equally across the remaining 50% of  $K_{MG}$ . This  
350 community represented a community with relatively moderate species evenness. In the third  
351 community, 90% of the metagenomic DNA segments were equally distributed across 10% of  $K_{MG}$ , and  
352 the remaining 10% of the metagenomic DNA segments were distributed equally across the remaining  
353 90% of  $K_{MG}$ . This community represented a community with relatively low species evenness. The last  
354 community had 10 equally-sized groups, or octaves (i.e.,  $s$  was the same in all groups). The abundance  
355 of the metagenomic DNA segments in each group followed a lognormal distribution which has been  
356 observed in true microbial populations (e.g., (7, 16)). The functional form for modeling abundances  
357 was based on the functional form of a lognormal community (34):

$$358 \quad S(R) = S_0 e^{-a^2 R^2} \quad (9)$$

359 where  $S_0$  was treated as the maximum relative of abundance ( $S_0 = 1$ ),  $a$  was the inverse width of the  
360 distribution,  $R$  was treated as the positive octave range spanning 0 to 9, and  $S(R)$  represented the  
361 abundance for a given octave. For the lognormal abundance distribution in Fig 2D,  $a$  was set to a value  
362 of 0.2. Each hypothetical community started with a unique number of  $k$ -sized reads  $|K_{MG}| = 1 \times 10^2$ .

363  $|K_{MG}|$  was incrementally increased at 10 equally-spaced, linear steps to a maximum of  $|K_{MG}| = 1 \times$   
364  $10^6$ . As  $|K_{MG}|$  increased, all community structures remained constant. Graphical representation of rank  
365 abundance in Fig 2a was normalized by a given  $|K_{MG}|$  to reflect that populations retained the same  
366 structure even as population size varied. We defined a normalized rank abundance  $r_n$  such that

$$367 \quad r_n = \frac{r}{s} \quad (10)$$

368 where  $r$  and  $s$  are untransformed rank abundance and richness, respectively. Thus, the most abundant  $k$ -  
369 mer in a in a metagenome population has a normalized rank abundance of  $1/s$  and the least abundant  
370 has a normalized rank abundance of 1. For each community, at each step, the expected number of  
371 sequences was calculated using equation 7. The expected number of sequences as a function of  $|K_{MG}|$   
372 were modeled with linear regressions.

373 Equation 7 gives the expected number of sequences required to sequence any sized community  
374 to exhaustion. Numerical sequencing simulations were performed to determine the number of  
375 sequences necessary to sequence a subset of all unique DNA ( $K_{MG}$ ). These numerical sequencing  
376 simulations were applied to four hypothetical community structures described above. Numerical  
377 simulations were performed such that  $|K_{MG}| = 3 \times 10^7, 4 \times 10^7, 5 \times 10^7, 7 \times 10^7, 9 \times 10^7$ , and  
378  $1 \times 10^8$ . During each of these simulations, the parameters read length ( $k$ ) and average genome size ( $l$ )  
379 were set to 100 and  $1 \times 10^6$ , respectively, for all  $g$ . Random elements from  $K_{MG}$  were selected with  
380 replacement to simulate a sequencing event. Numerical simulations were performed until the fraction  
381 of  $|K_{MG}|$  sequenced was 50%, 70%, 90%, 95%, 99%, or 100%. A weight distribution was applied to  
382 elements in a given  $K_{MG}$ . The weight distribution biased sequencing to reflect the relative abundances  
383 of the four hypothetical communities described above. The fraction of  $|K_{MG}|$  sequenced was evaluated  
384 every  $1 \times 10^7$  sequences. Numerical simulations were performed in triplicate for all  $|K_{MG}|$  and all  
385 target fractions of  $|K_{MG}|$ .

386 We explored the influence of community evenness on required sequencing depth by performing  
387 numerical sequencing simulations on 6 different lognormally-distributed communities. The numerical  
388 sequencing simulations were similar to the simulations described above. The 6 lognormal communities  
389 were modeled such that each community had  $S_0=1$ , 10 equally-sized octaves, and  $|K_{MG}| = 1 \times 10^7$ .  
390 The difference between the 6 lognormal distributions was due to variations in  $a$  where  $a=0$ ,  $a=0.005$ ,  
391  $a=0.008$ ,  $a=0.01$ ,  $a=0.015$ , and  $a=0.02$ . Evenness was represented using Pielou evenness index (9),  
392 which is the ratio of the Shannon diversity index (35) for a given community to that of an even  
393 community of the same richness. Shannon diversity was calculated in the context of a metagenomes  
394 such that:

$$395 H_{MG} = \sum_{j \in K_{MG}} -p_j \log(p_j)$$

396 (11)

397 where  $p_j$  is the proportion that the  $j^{\text{th}}$   $k$ -sized read represents among all unique DNA sequences in the  
398 metagenome. Thus, the Pielou evenness index (9) was calculated such that:

$$399 J = \frac{H_{MG}'}{H_{MG,max}} \quad (12)$$

400 where  $J$  was the Pielou evenness index,  $H_{MG}'$  was the metagenome Shannon diversity index, and  
401  $H_{MG,max}$  represented the metagenome Shannon diversity index when all  $p_j$  were equal (i.e.,  $a=0$ ).

402 Lastly, numerical simulations were performed to determine the sequencing depth necessary to  
403 achieve a target fraction for an individual metagenome ( $g_{MG}$ ). Target fractions were increased from 0.5  
404 to 1 at 100 linearly-spaced intervals. The fraction of the metagenome community ( $G_{MG}$ ) that  $g_{MG}$   
405 represented varied from 1% to 100% in 30 lognormally-spaced intervals. The target genome sizes ( $l$ )  
406 varied such that  $l=0.5 \times 10^6$ ,  $l=1 \times 10^6$ ,  $l=2 \times 10^6$ ,  $l=3 \times 10^6$ ,  $l=5 \times 10^6$ ,  $l=10 \times 10^6$ ,  $l=15 \times 10^6$ , and  $l=20 \times 10^6$ .  
407 The sequencing depth for a given combination of target fraction, genome size, and fraction of the  
408 metagenome community was modeled using the gam function (mgcv R package; (36)). For modeling

409 purposes, target fraction was raised to the 12<sup>th</sup> power and both genome size and sequences were log-  
410 transformed. The number of smooth dimensions for fraction of community, genome size, and target  
411 fraction were heuristically varied till the resulting fit demonstrated residuals with a normal distribution.  
412 Note that the objective here was not build a predictive model but simply a first order approximation for  
413 simulations performed here.

#### 414 DATA AVAILABILITY

415 All simulations and codes used for modeling sequencing depth are freely available on Github at:  
416 [https://github.com/taylorroyalty/sequence\\_simulation\\_code](https://github.com/taylorroyalty/sequence_simulation_code).

#### 417 **Acknowledgements**

418 This research was supported by the National Science Foundation and a C-DEBI subaward (contribution  
419 number to be determined).

420 **References**

- 421 1. Knight R, Vrbanac A, Taylor BC, Aksenov A, Callewaert C, Debelius J, Gonzalez A, Kosciolk  
422 T, McCall L-I, McDonald D, Melnik A V, Morton JT, Navas J, Quinn RA, Sanders JG,  
423 Swafford AD, Thompson LR, Tripathi A, Xu ZZ, Zaneveld JR, Zhu Q, Caporaso JG, Dorrestein  
424 PC. 2018. Best practices for analysing microbiomes. *Nat Rev Microbiol*.
- 425 2. Parks DH, Rinke C, Chuvochina M, Chaumeil P-A, Woodcroft BJ, Evans PN, Hugenholtz P,  
426 Tyson GW. 2017. Recovery of nearly 8,000 metagenome-assembled genomes substantially  
427 expands the tree of life. *Nat Microbiol* 903:1–10.
- 428 3. Tully BJ, Graham ED, Heidelberg JF. 2018. The reconstruction of 2,631 draft metagenome-  
429 assembled genomes from the global oceans. *Sci Data* 5:1–8.
- 430 4. Baker BJ, Lazar CS, Teske AP, Dick GJ. 2015. Genomic resolution of linkages in carbon,  
431 nitrogen, and sulfur cycling among widespread estuary sediment bacteria. *Microbiome* 3:14.
- 432 5. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. 2017. Shotgun metagenomics, from  
433 sampling to analysis. *Nat Biotechnol* 35:833–844.
- 434 6. Zumel N, Mount J. 2014. *Practical Data Science with R*, 1st ed. Manning Publications Co.,  
435 Greenwich, CT, USA.
- 436 7. Galand PE, Casamayor EO, Kirchman DL, Lovejoy C. 2009. Ecology of the rare microbial  
437 biosphere of the Arctic Ocean. *Proc Natl Acad Sci* 106:22427–22432.
- 438 8. Locey KJ, Lennon JT. 2016. Scaling laws predict global microbial diversity. *Proc Natl Acad Sci*  
439 113:5970–5975.
- 440 9. Pielou EC. 1966. The measurement of diversity in different types of biological collections. *J*  
441 *Theor Biol* 13:131–144.
- 442 10. Land M, Hauser L, Jun S-R, Nookaew I, Leuze MR, Ahn T-H, Karpinets T, Lund O, Kora G,

- 443 Wassenaar T, Poudel S, Ussery DW. 2015. Insights from 20 years of bacterial genome  
444 sequencing. *Funct Integr Genomics* 15:141–161.
- 445 11. Bowers RM, Kyrpides NC, Stepanauskas R, Harmon-Smith M, Doud D, Reddy TBK, Schulz F,  
446 Jarett J, Rivers AR, Eloe-Fadrosh EA, Tringe SG, Ivanova NN, Copeland A, Clum A, Becraft  
447 ED, Malmstrom RR, Birren B, Podar M, Bork P, Weinstock GM, Garrity GM, Dodsworth JA,  
448 Yooseph S, Sutton G, Glöckner FO, Gilbert JA, Nelson WC, Hallam SJ, Jungbluth SP, Etema  
449 TJG, Tighe S, Konstantinidis KT, Liu WT, Baker BJ, Rattei T, Eisen JA, Hedlund B, McMahon  
450 KD, Fierer N, Knight R, Finn R, Cochrane G, Karsch-Mizrachi I, Tyson GW, Rinke C, Lapidus  
451 A, Meyer F, Yilmaz P, Parks DH, Eren AM, Schriml L, Banfield JF, Hugenholtz P, Woyke T.  
452 2017. Minimum information about a single amplified genome (MISAG) and a metagenome-  
453 assembled genome (MIMAG) of bacteria and archaea. *Nat Biotechnol* 35:725–731.
- 454 12. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM: assessing the  
455 quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome*  
456 *Res* 25:1043–55.
- 457 13. Roesch LFW, Fulthorpe RR, Riva A, Casella G, Hadwin AKM, Kent AD, Daroub SH, Camargo  
458 FAO, Farmerie WG, Triplett EW. 2007. Pyrosequencing enumerates and contrasts soil microbial  
459 diversity. *ISME J* 1:283–290.
- 460 14. Feng BW, Li XR, Wang JH, Hu ZY, Meng H, Xiang LY, Quan ZX. 2009. Bacterial diversity of  
461 water and sediment in the Changjiang estuary and coastal area of the East China Sea. *FEMS*  
462 *Microbiol Ecol* 70:236–248.
- 463 15. Rintala A, Pietilä S, Munukka E, Eerola E, Pursiheimo JP, Laiho A, Pekkala S, Huovinen P.  
464 2017. Gut microbiota analysis results are highly dependent on the 16s rRNA gene target region,  
465 whereas the impact of DNA extraction is minor. *J Biomol Tech* 28:19–30.

- 466 16. Kang S, Rodrigues JLM, Ng JP, Gentry TJ. 2016. Hill number as a bacterial diversity measure  
467 framework with high-throughput sequence data. *Sci Rep* 6:1–4.
- 468 17. Zaheer R, Noyes N, Ortega Polo R, Cook SR, Marinier E, Van Domselaar G, Belk KE, Morley  
469 PS, McAllister TA. 2018. Impact of sequencing depth on the characterization of the microbiome  
470 and resistome. *Sci Rep* 8:5890.
- 471 18. Sims D, Sudbery I, Ilott NE, Heger A, Ponting CP. 2014. Sequencing depth and coverage: Key  
472 considerations in genomic analyses. *Nat Rev Genet* 15:121–132.
- 473 19. Leinonen R, Sugawara H, Shumway M. 2010. The Sequence Read Archive. *Nucleic Acids Res*  
474 39:2010–2012.
- 475 20. Schmieder R, Edwards R. 2011. Quality control and preprocessing of metagenomic datasets.  
476 *Bioinformatics* 27:863–864.
- 477 21. Karsenti E, Acinas SG, Bork P, Bowler C, de Vargas C, Raes J, Sullivan M, Arendt D, Benzoni  
478 F, Claverie JM, Follows M, Gorsky G, Hingamp P, Iudicone D, Jaillon O, Kandels-Lewis S,  
479 Krzic U, Not F, Ogata H, Pesant S, Reynaud EG, Sardet C, Sieracki ME, Speich S, Velayoudon  
480 D, Weissenbach J, Wincker P. 2011. A holistic approach to marine Eco-systems biology. *PLoS*  
481 *Biol* 9:7–11.
- 482 22. Howe AC, Jansson JK, Malfatti SA, Tringe SG, Tiedje JM, Brown CT. 2014. Tackling soil  
483 diversity with the assembly of large, complex metagenomes. *Proc Natl Acad Sci* 111:4904–  
484 4909.
- 485 23. Schirmer M, Smeekens SP, Vlamakis H, Jaeger M, Oosting M, Franzosa EA, Jansen T, Jacobs  
486 L, Bonder MJ, Kurilshikov A, Fu J, Joosten LAB, Zhernakova A, Huttenhower C, Wijmenga C,  
487 Netea MG, Xavier RJ. 2016. Linking the Human Gut Microbiome to Inflammatory Cytokine  
488 Production Capacity. *Cell* 167:1125–1136.e8.

- 489 24. Graham ED, Heidelberg JF, Tully BJ. 2017. BinSanity: unsupervised clustering of  
490 environmental microbial assemblies using coverage and affinity propagation. *PeerJ* 5:e3035.
- 491 25. Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence  
492 data. *Bioinformatics* 30:2114–20.
- 493 26. Li D, Luo R, Liu C-M, Ting H-F, Sadakane K, Yamashita H, Lam T-W. 2016. MEGAHIT v1.0:  
494 A fast and scalable metagenome assembler driven by advanced methodologies and community  
495 practices. *Methods* 102:3–11.
- 496 27. Fu L, Niu B, Zhu Z, Wu S, Li W. 2012. CD-HIT: accelerated for clustering the next-generation  
497 sequencing data. *Bioinformatics* 28:3150–3152.
- 498 28. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
499 9:357–9.
- 500 29. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. 2010. Prodigal:  
501 prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics*  
502 11:119.
- 503 30. Zwietering MH, Jongenburger I, Rombouts FM, Van't Riet K. 1990. Modeling of the bacterial  
504 growth curve. *Appl Environ Microbiol* 56:1875–1881.
- 505 31. Rosselló-Móra R, Amann R. 2015. Past and future species definitions for Bacteria and Archaea.  
506 *Syst Appl Microbiol* 38:209–216.
- 507 32. Bentley DR, Balasubramanian S, Swerdlow HP, Smith GP, Milton J, Brown CG, Hall KP, Evers  
508 DJ, Barnes CL, Bignell HR, Boutell JM, Bryant J, Carter RJ, Keira Cheetham R, Cox AJ, Ellis  
509 DJ, Flatbush MR, Gormley NA, Humphray SJ, Irving LJ, Karbelashvili MS, Kirk SM, Li H, Liu  
510 X, Maisinger KS, Murray LJ, Obradovic B, Ost T, Parkinson ML, Pratt MR, Rasolonjatovo IMJ,  
511 Reed MT, Rigatti R, Rodighiero C, Ross MT, Sabot A, Sankar S V., Scally A, Schroth GP,



512 Smith ME, Smith VP, Spiridou A, Torrance PE, Tzonev SS, Vermaas EH, Walter K, Wu X,  
513 Zhang L, Alam MD, Anastasi C, Aniebo IC, Bailey DMD, Bancarz IR, Banerjee S, Barbour SG,  
514 Baybayan PA, Benoit VA, Benson KF, Bevis C, Black PJ, Boodhun A, Brennan JS, Bridgham  
515 JA, Brown RC, Brown AA, Buermann DH, Bundu AA, Burrows JC, Carter NP, Castillo N,  
516 Catenazzi MCE, Chang S, Neil Cooley R, Crake NR, Dada OO, Diakoumakos KD, Dominguez-  
517 Fernandez B, Earnshaw DJ, Egbujor UC, Elmore DW, Etchin SS, Ewan MR, Fedurco M, Fraser  
518 LJ, Fuentes Fajardo K V., Scott Furey W, George D, Gietzen KJ, Goddard CP, Golda GS,  
519 Granieri PA, Green DE, Gustafson DL, Hansen NF, Harnish K, Haudenschild CD, Heyer NI,  
520 Hims MM, Ho JT, Horgan AM, Hoschler K, Hurwitz S, Ivanov D V., Johnson MQ, James T,  
521 Huw Jones TA, Kang GD, Kerelska TH, Kersey AD, Khrebtukova I, Kindwall AP, Kingsbury  
522 Z, Kokko-Gonzales PI, Kumar A, Laurent MA, Lawley CT, Lee SE, Lee X, Liao AK, Loch JA,  
523 Lok M, Luo S, Mammen RM, Martin JW, McCauley PG, McNitt P, Mehta P, Moon KW,  
524 Mullens JW, Newington T, Ning Z, Ling Ng B, Novo SM, O'Neill MJ, Osborne MA, Osnowski  
525 A, Ostadan O, Paraschos LL, Pickering L, Pike AC, Pike AC, Chris Pinkard D, Pliskin DP,  
526 Podhasky J, Quijano VJ, Raczy C, Rae VH, Rawlings SR, Chiva Rodriguez A, Roe PM, Rogers  
527 J, Rogert Bacigalupo MC, Romanov N, Romieu A, Roth RK, Rourke NJ, Ruediger ST, Rusman  
528 E, Sanches-Kuiper RM, Schenker MR, Seoane JM, Shaw RJ, Shiver MK, Short SW, Sizto NL,  
529 Sluis JP, Smith MA, Ernest Sohna Sohna J, Spence EJ, Stevens K, Sutton N, Szajkowski L,  
530 Tregidgo CL, Turcatti G, Vandevondele S, Verhovsky Y, Virk SM, Wakelin S, Walcott GC,  
531 Wang J, Worsley GJ, Yan J, Yau L, Zuerlein M, Rogers J, Mullikin JC, Hurles ME, McCooke  
532 NJ, West JS, Oaks FL, Lundberg PL, Klenerman D, Durbin R, Smith AJ. 2008. Accurate whole  
533 human genome sequencing using reversible terminator chemistry. *Nature* 456:53–59.

534 33. Flajolet P, Gardy D, Thimonier L. 1992. Birthday paradox, coupon collectors, caching

- 535 algorithms and self-organizing search. *Discret Appl Math* 39:207–229.
- 536 34. Magurran AE. 1988. *Ecological Diversity and Its Measurement*, 1st ed. Croom Helm Ltd.
- 537 35. Shannon CE. 1948. A mathematical theory of communication. *Bell Syst Tech J* 27:379–423.
- 538 36. Wood S. 2017. mgcv: Mixed GAM Computation Vehicle with GCV/AIC/REML Smoothness
- 539 Estimation.
- 540 CRAN <https://cran.r-project.org/package=mgcv> Retrieved 14 May 2018

541 **Tables**

542 **Table 1.** Estimates of fit coefficients for the Gompertz equation (equation 1) for the effective MAGs as  
543 a function of sequencing depth in published datasets from ocean surface water, estuarine sediment,  
544 maize soil, and the human gut.  $p$  values for all coefficients were  $\ll 0.05$ .

Environment	$A (\pm SE)$	$\mu (\pm SE)$	$\lambda (\pm SE)$	MAG Yield*
Ocean Surface Water	97.67 (4.15)	5.84 (0.31)	1.16 (0.33)	0.88
Estuary Sediment	26.25 (2.11)	1.63 (0.06)	3.70 (0.24)	0.86
Maize Soil	43.65 (1.98)	1.43 (0.07)	6.13 (0.60)	0.70
Human Gut	17.49 (1.02)	5.01 (0.46)	0.67 (0.14)	0.90

545 \*Calculated as the ratio of the maximum effective MAGs experimentally observed to maximum  
546 effective MAGs ( $A$ )

547 Table 2. Summary of sequence datasets analyzed with the MAG pipeline.

Environment	NCBI SRA Accession	Sequencing Platform	Total Reads*	High Quality Bases*	General Notes	Citation
Ocean Surface Water	ERR599029	Illumina HiSeq 2000	337,228,196	33,396,930,215	Caribbean Sea (5 mbsl)	(21)
Estuary Sediment	SRR5248164	Illumina HiSeq 2000	113,025,112	15,887,161,501	Sulfate Zone (8-10 cmbsf)	(4)
Maize Soil	SRR351473	Illumina HiSeq 2000	472,686,494	38,246,948,858	Surface Soil	(22)
Human Gut	SRR5127631	Illumina HiSeq 2000	50,951,710	4,846,948,241	--	(23)

548 \*Combination of forward and backwards pair-end reads.

549 **Figures**

550 **Fig 1.** The influence that the parameters  $A$ ,  $\mu$ , and  $\lambda$  had on the Gompertz equation (A). The property of  
551 the Gompertz equation that each parameter influences is colored red. Mean MAG completeness (B),  
552 and mean MAG contamination as a function of simulated sequencing depth (Gb) for sequence datasets  
553 of the human gut, maize soil, estuarian sediment, and surface ocean microbiomes, using the pipeline  
554 described in the methods section. Translucent lines in (A) correspond to nonlinear least squares fits of  
555 the Gompertz equation to the respective environmental dataset.

556 **Fig 2.** Average expected sequences required to fully sequence four different community structures, one  
557 with relatively high community evenness (A), relatively moderate community evenness (B), relatively  
558 low community evenness (C), and one with a lognormal community structure (D), were predicted using  
559 linear regressions (E) and the log of  $|K_{MG}|$  from equation 6 as a predictor.

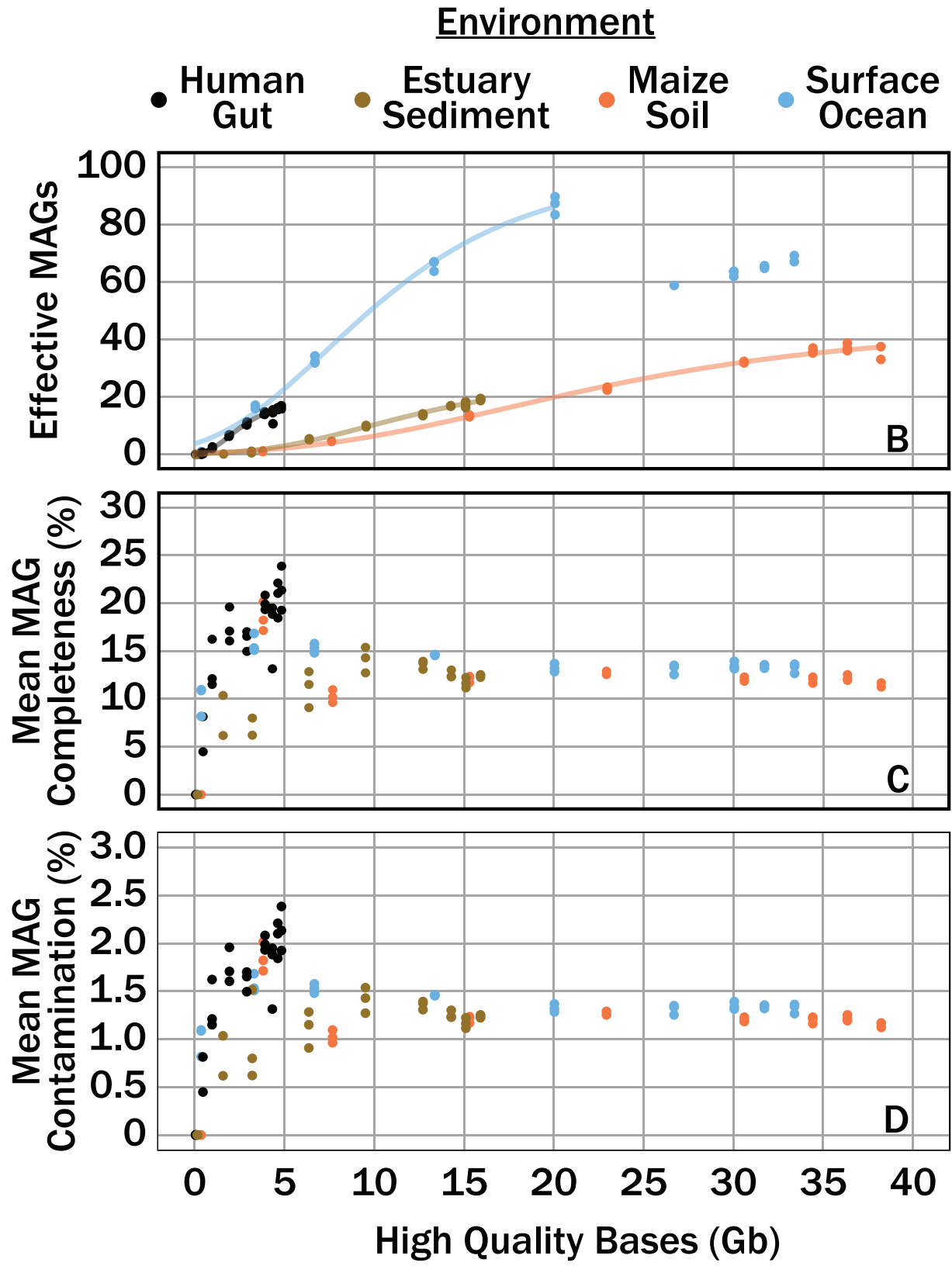
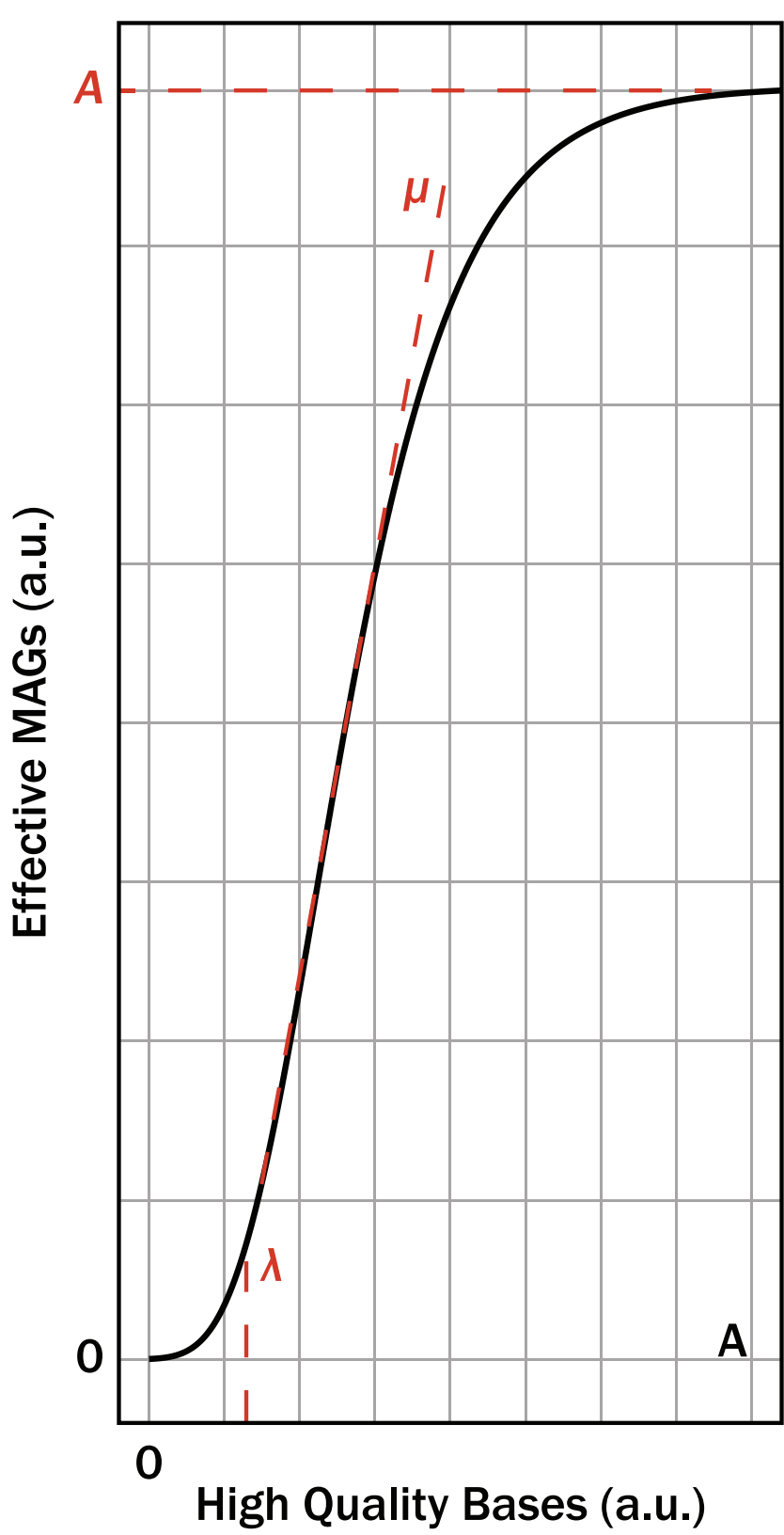
560 **Fig 3.** Sequences necessary to reach variable target sequencing depths (colors) for four different  
561 community structures, one with relatively high community evenness (A), relatively moderate  
562 community evenness (B), relatively low community evenness (C), and one with a lognormal  
563 community structure (D). Red translucent lines correspond with linear regression curves for the  
564 respective community in Fig 2E.

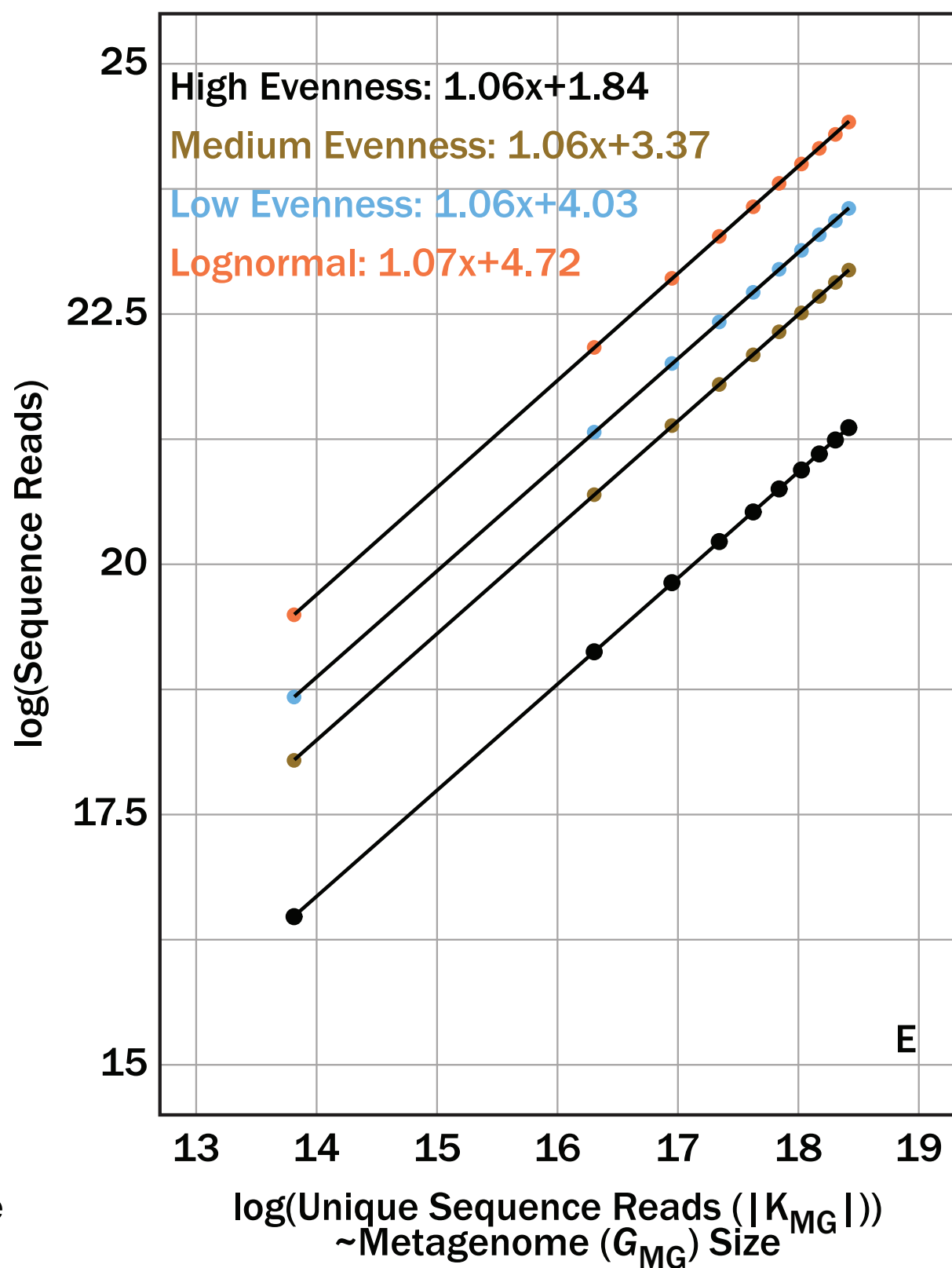
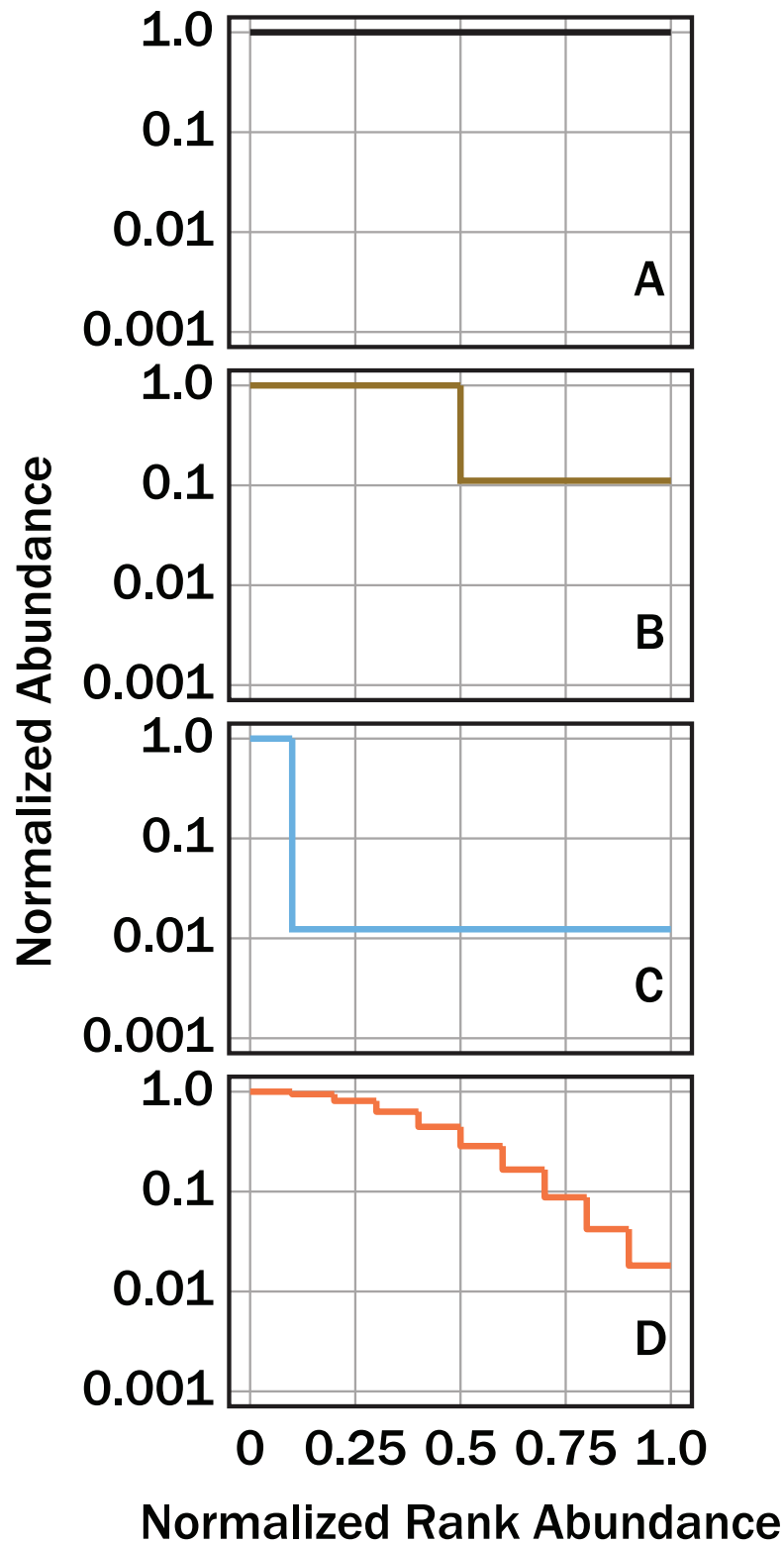
565 **Fig 4.** Numerical sequencing simulations applied to 6 hypothetical communities with different  
566 lognormal distributions that were defined by the parameter,  $\alpha$ , from equation 9 (A). The number of  
567 sequences necessary to sequence a target fraction of  $|K_{MG}|$  (dashed contours) as a function of the  
568 Pielou evenness index,  $J$ , for a given lognormal community structure (B).



569 **Fig 5.** Numerical sequencing simulations show the number of bases (color bar) required to sequence a  
570 target fraction of a genome which represents a given fraction of a community metagenome. Genomes  
571 evaluated were  $0.5 \times 10^6$  (A),  $2 \times 10^6$  (B),  $5 \times 10^6$  (C),  $10 \times 10^6$  (D), and  $20 \times 10^6$  (E) base pairs long.

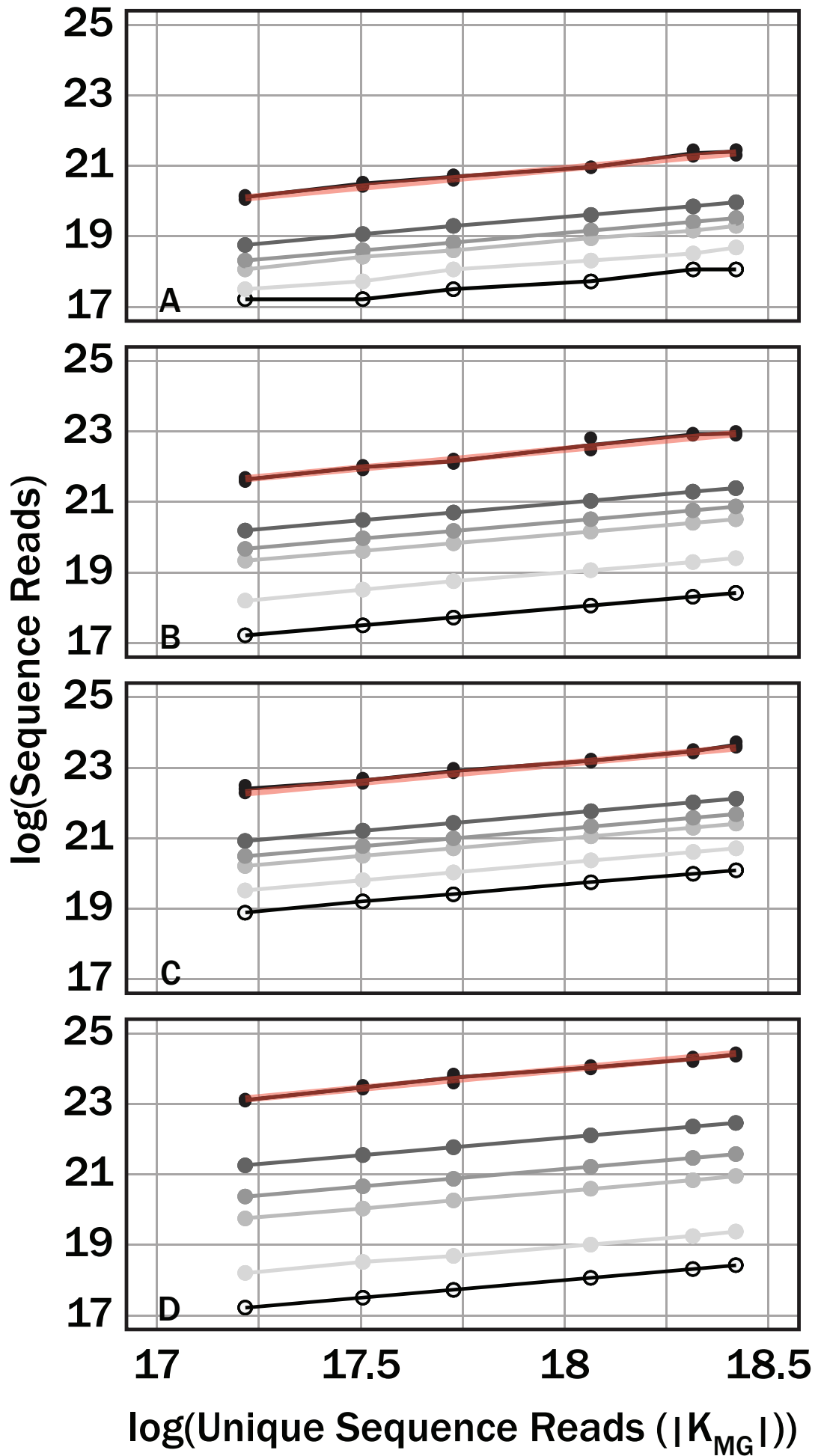
572 **Fig 6.** A cartoon illustrating an example microbial community ( $G$ ), metagenomes for genomes ( $g_{MG,i}$ )  
573 within  $G$ , and the overall metagenome for the given microbial community ( $G_{MG}$ ). In this example, there  
574 are 6 MAGs ( $s=6$ ) and a total of 13 microbes. (A) Black circles represent individual microbes whose  
575 genomes are averaged together,  $g$ . The average genome,  $g$ , are indicated by different color inner-  
576 circles. (B) Individual average genomes can be sequenced at  $K$  unique positions depending on the  
577 characteristic read length,  $k$ , of a sequencer. (C) All unique positions that can be sequenced for a given  
578 genome,  $g$ , defines the metagenome,  $g_{MG}$ , for the  $i^{\text{th}}$  genome,  $g_i$ . (D) Replacing all individual genomes  
579 in (A) with metagenomes,  $g_{MG}$ , gives the metagenome of the microbial community,  $G_{MG}$ .  
580  
581





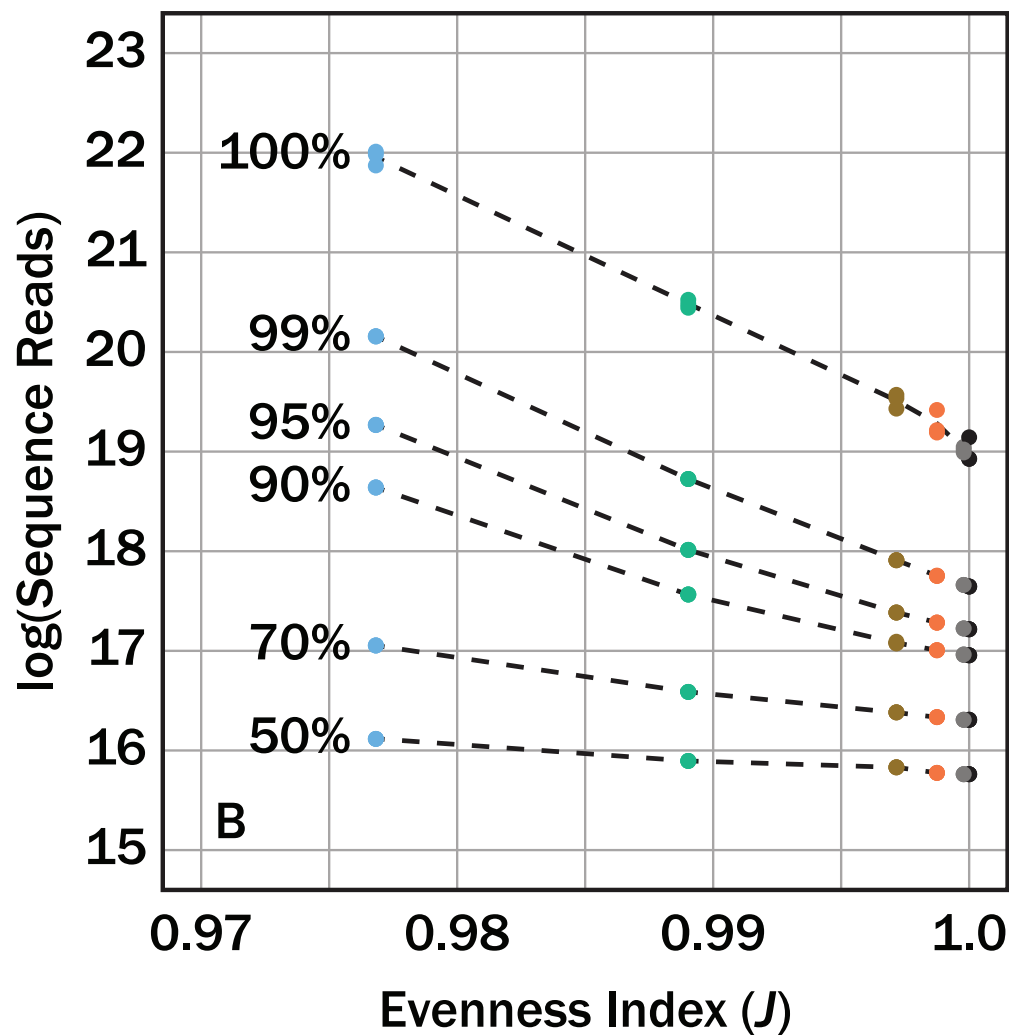
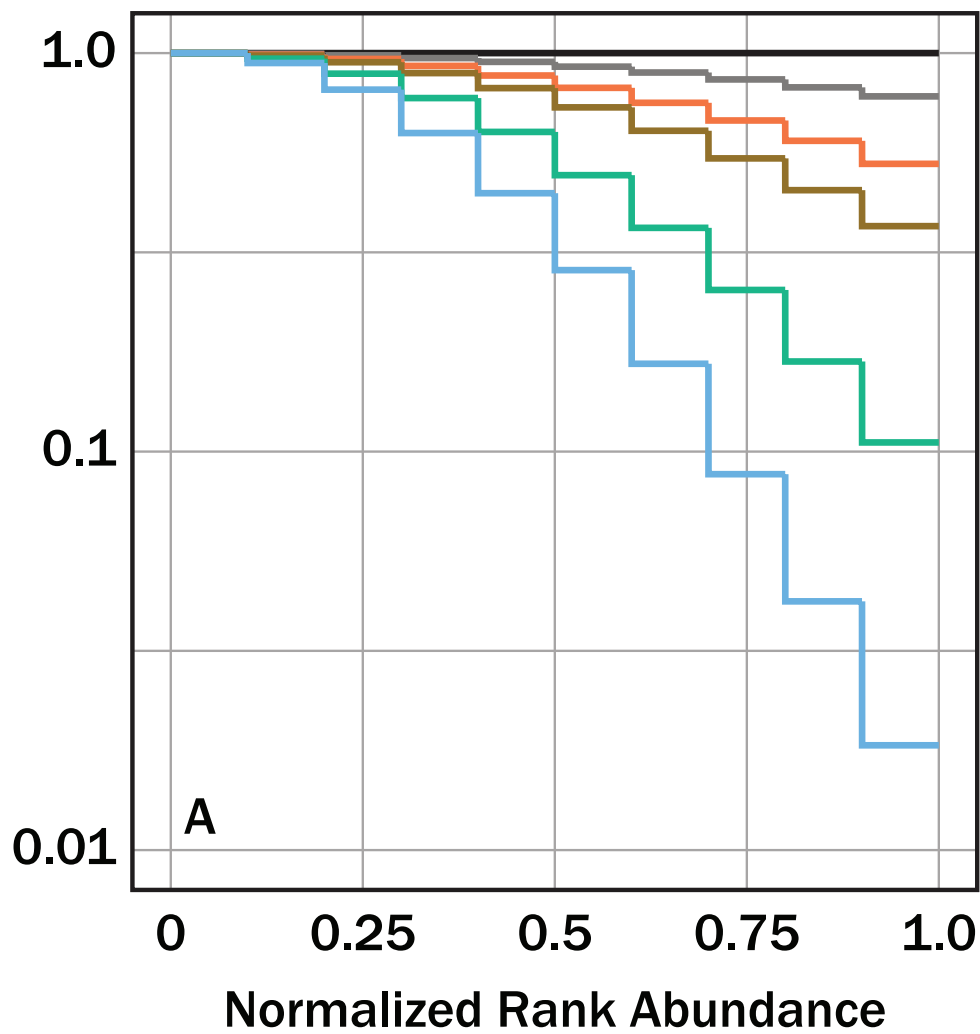
# Fraction of $|K_{MG}|$

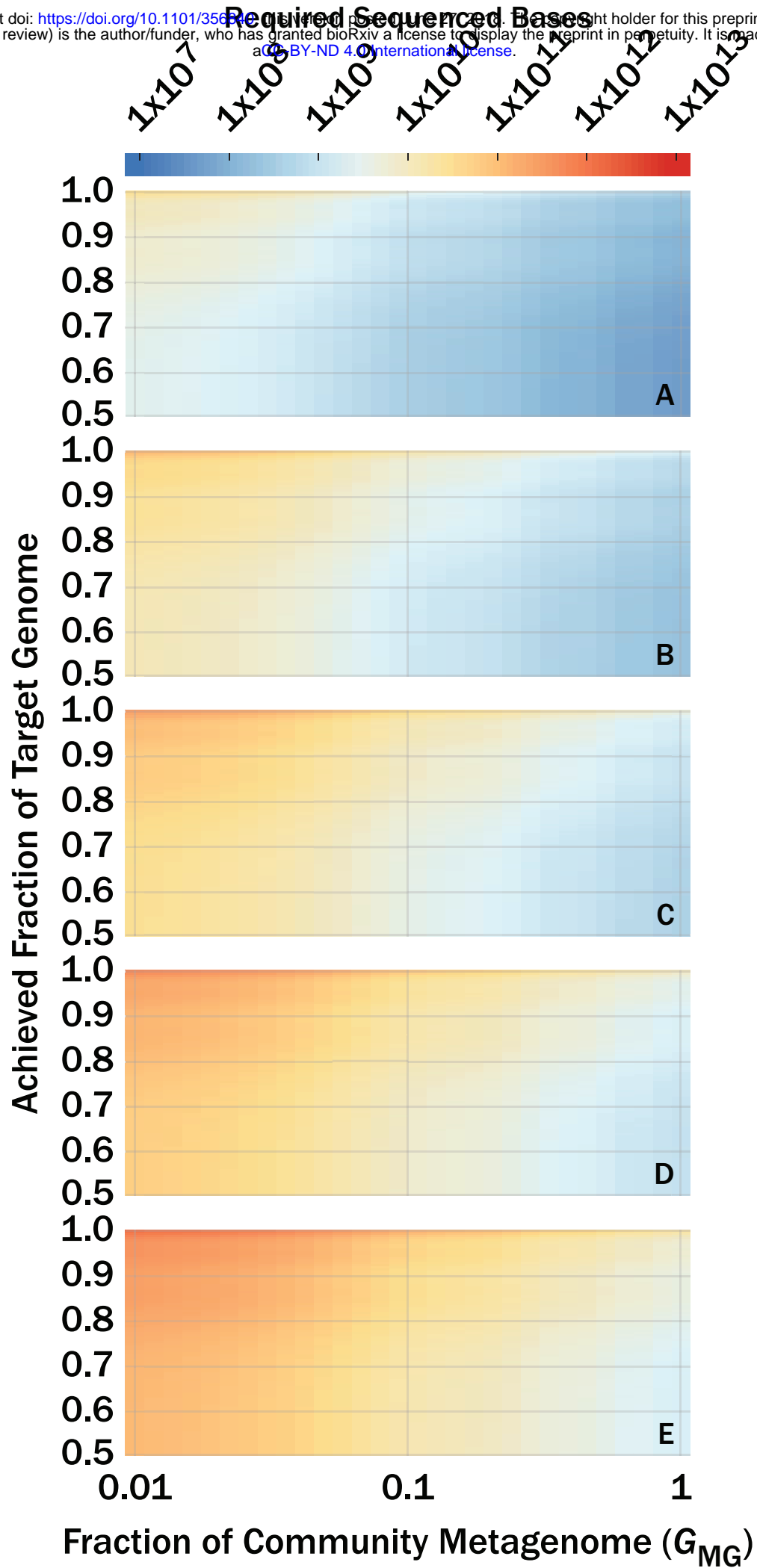
● 100% ● 99% ● 95% ● 90% ● 70% ○ 50%

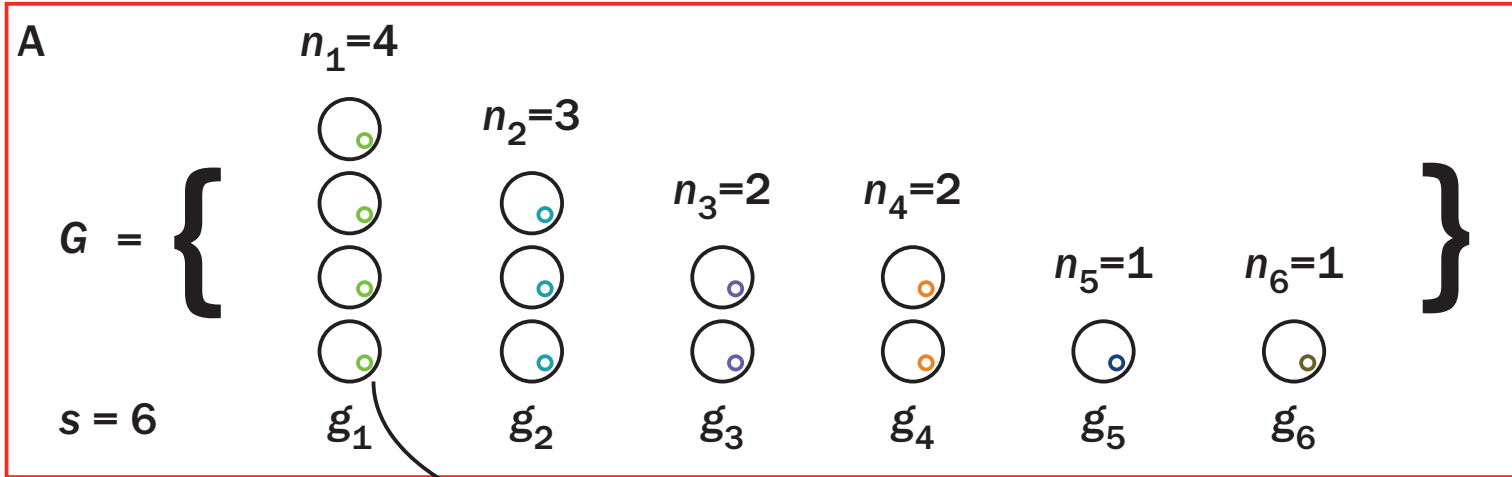


Lognormal Community Structure (a)

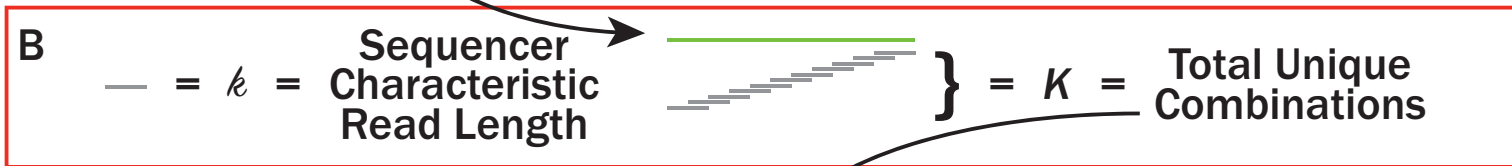
— 0 • — 0.005 • — 0.008 • — 0.01 • — 0.015 • — 0.02 •



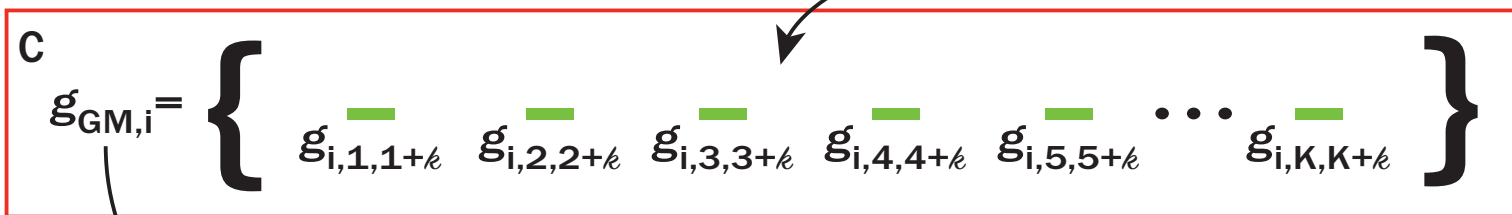




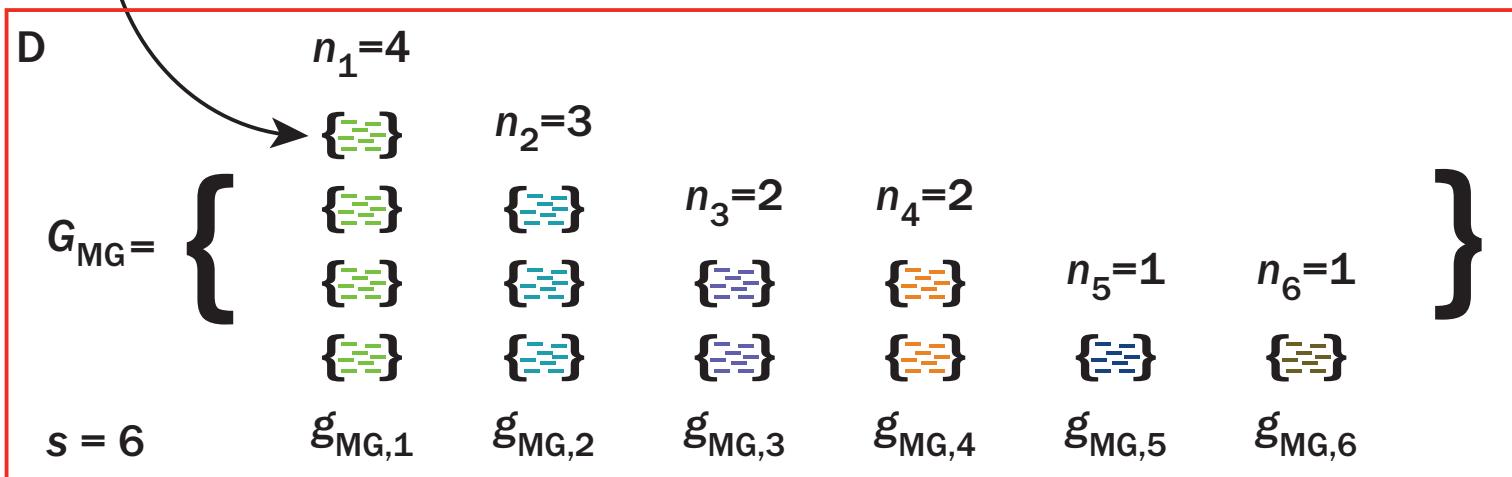
Equation 2



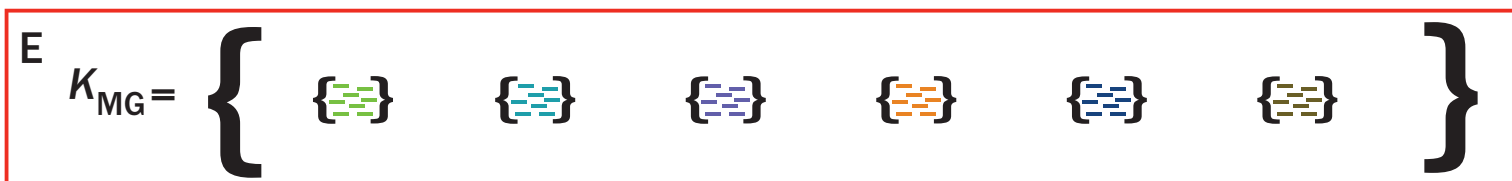
Equation 3



Equation 4



Equation 5



Equation 6