

1 **Designing Minimal Genomes Using Whole-Cell Models**

2 Joshua Rees^{1,2^} and Oliver Chalkley^{1,3,4^}, Sophie Landon^{1,3}, Oliver Purcell⁵,

3 Lucia Marucci^{1,3,6+} and Claire Grierson^{1,2+*}

4 ¹BrisSynBio, University of Bristol, Bristol BS8 1TQ, UK;

5 ²School of Biological Sciences, University of Bristol, Bristol Life Sciences Building, 24 Tyndall

6 Avenue, Bristol, BS8 1TQ, UK;

7 ³Department of Engineering Mathematics, University of Bristol, Bristol BS8 1UB, UK;

8 ⁴Bristol Centre for Complexity Science, Department of Engineering Mathematics, University of

9 Bristol, Bristol BS8 1UB, UK;

10 ⁵Prospect Bio, 150 N Hill Drive, Ste 14, Brisbane, CA 94005, USA;

11 ⁶School of Cellular and Molecular Medicine, University of Bristol, Bristol BS8 1UB, UK;

12 ^Co-first authors + Co-last authors * Corresponding author

13 Corresponding author: Prof. Claire Grierson (claire.grierson@bristol.ac.uk)

14 **Abstract**

15 In the future, entire genomes tailored to specific functions and environments could be designed using
16 computational tools. However, computational tools for genome design are currently scarce. Here we
17 present algorithms that enable the use of design-simulate-test cycles for genome design, using
18 genome minimisation as a proof-of-concept. Minimal genomes are ideal for this purpose as they
19 have a simple functional assay, the cell either replicates or not. We used the first (and currently only
20 published) whole-cell model, for the bacterium *Mycoplasma genitalium*. Our computational
21 design-simulate-test cycles discovered novel *in-silico* minimal genomes smaller than *JCVI-Syn3.0*, a
22 bacteria with, currently, the smallest genome that can be grown in pure culture. In the process, we
23 identified 10 low essentiality genes, 18 high essentiality genes, and produced evidence for at least
24 two *Mycoplasma genitalium in-silico* minimal genomes. This work brings combined computational
25 and laboratory genome engineering a step closer.

26 **Introduction**

27 For genome-scale engineering and design, minimal genomes are currently the best proof-of-concept
28 ¹. These are reduced genomes containing only genes essential for life, provided there is a rich growth
29 medium and no external stressors ^{1,2}. The largest scale efforts in genome minimisation to date
30 include: *JCVI-Syn3.0*, a 50% gene reduction of *Mycoplasma mycoides* ²; several strains of
31 *Escherichia coli* reduced by 38.9% ³ and 35% ⁴ of their base pairs *in-vivo*; an *E.coli* gene reduction of
32 77.6% in *Saccharomyces cerevisiae* ⁵; and two 36% gene reductions of *Bacillus subtilis* ⁶.
33 Initially, these were either prescriptively designed, with requirements based on current biological
34 knowledge, or based on extensive laboratory testing of individual genes. These were then developed
35 iteratively in the lab, a time consuming and expensive process due to the limitations of current
36 techniques and unexpected cell death, likely caused by unknown genetic interactions. This hinders
37 progress as laboratories can only follow a small number of high-risk research avenues with limited
38 ability to backtrack ¹.
39 Another approach, building novel organisms from the bottom-up, is currently infeasible in the
40 majority of bacteria due to technological and economic constraints. Megabase sized genomes can

41 be constructed within yeast ^{5,7}, but one of the most promising approaches, genome transplantation,
42 has only been demonstrated in a subset of *Mycoplasmas* ⁸⁻¹⁰ and is mutagenic ⁹.

43 A further barrier to genome minimisation is the dynamic nature of gene essentiality. A simple
44 definition of a cell as “living” is if it can reproduce, an “essential” gene being indispensable for cell
45 division. A “non-essential” gene can be removed and leave division intact ^{1,11}. But a cell’s need for
46 specific genes (and their products) is dependent on the external cellular environment and on the
47 genomic context ¹ (the presence or absence of other genes, and resulting gene products, in the
48 genome), which can change each time a gene is removed. Some essential genes can become
49 dispensable with the removal of a particular gene (i.e. a toxic byproduct is no longer produced, so its
50 removal is unnecessary), referred to as “protective essential” genes ^{1,12,13}. Likewise, some
51 non-essential genes become essential when a functionally equivalent gene is removed, leaving a
52 single pathway to a metabolite (a “redundant essential” gene pair). Additionally, gene products can
53 perform together as a complex, with individually non-essential genes involved in producing an
54 essential function ¹⁴; when enough deletions accumulate to disrupt the group, the remaining genes
55 become essential. The cellular death that occurs when redundant essential genes are removed
56 together, or complexes are disrupted, is referred to as synthetic lethality ^{2,15,16}. A recent review ¹
57 updates gene essentiality from a binary categorisation to a gradient with four categories: no
58 essentiality (if dispensable in all contexts), low essentiality (if dispensable in some contexts, i.e.
59 redundant essential and complexes), high essentiality (if indispensable in most contexts, i.e.
60 protective essential), and complete essentiality (if indispensable in all contexts). These broad labels
61 describe an individual gene’s essentiality in different genomic contexts, and are compatible with
62 other labels that explain underlying mechanisms and interactions in greater levels of detail.

63 To overcome the above, large-scale problems we used existing computational models with novel
64 genome design algorithms to investigate 10,000s of gene knockout combinations *in-silico*, with rapid
65 feedback and iteration. Testing potential genome reductions at scale for lethal interactions should
66 produce functional *in-silico* genomes, which can be implemented *in-vivo* with a lower risk of failure.

67 This generation of non-prescriptive designs, with no assumed biological requirements outside those
68 inherent in the model, increases the likelihood of novel findings.

69 We used the *Mycoplasma genitalium* (*M.genitalium*) whole-cell model ¹⁷, which describes the
70 smallest culturable, self-replicating, natural organism ¹⁸ (at the time the model was built). It is the only
71 existing model of a cell's individual molecules that includes the function of every known gene
72 product (401 of the 525 *M.genitalium* genes), making it capable of modelling genes in their genomic
73 context ¹⁷. A single cell is simulated from random initial conditions until the cell divides or reaches a
74 time limit. The model combines 28 cellular submodels, with parameters from >900 publications and
75 >1,900 experimental observations, resulting in 79% accuracy for single-gene knockout essentiality ¹⁷.
76 Outside of single-gene knockout simulations, it has been used to investigate discrepancies between
77 the model and real-world measurements ^{17,19}, design synthetic genetic circuits in the context of the
78 cell ²⁰, and make predictions about the use of existing antibiotics against new targets ²¹.

79 We produced two genome design algorithms (Minesweeper and the Guess/Add/Mate Algorithm
80 (GAMA)) which use the *M.genitalium* whole-cell model to generate minimal genome designs. Using
81 these computational tools we found functional *in-silico* minimal genomes, between 33 and 53 genes
82 smaller than the most recent predictions for a reduced *Mycoplasma* genome of 413 genes ^{2,15,16}.
83 These *in-silico* genomes are ideal candidates for further *in-vivo* testing.

84 **Results**

85 **Genome Design Tools: Minesweeper and GAMA**

86 Minesweeper and GAMA conduct whole-cell model simulations in three step cycles: design
87 (algorithms select possible gene deletions); simulate (the genome minus those deletions); and test
88 (analyse the *in-silico* cell produced). Simulations that produce dividing cells go through to the next
89 cycle of simulations. The number of gene deletions increases in each cycle, producing progressively
90 smaller genomes. Minesweeper and GAMA have generated 2157 and 53,451 of *in-silico* genomes
91 respectively to date, but for brevity only the smallest genomes are presented here.

92 Minesweeper is a four stage algorithm inspired by divide and conquer algorithms ²², initially
93 investigating genes individually to identify complete/high essentiality genes, before breaking the
94 genome into differently sized subsets to broadly test, then accumulating deletions and identifying low
95 essential genes as they appear. It deletes genes in groups that get progressively smaller until it
96 reaches individual gene deletions, and only deletes non-essential genes (as determined by
97 single-gene knockout simulations, see Initial Input below). By not considering essential genes the
98 search area is reduced, which makes it capable of producing minimal genome size reductions
99 quickly (within two days). It uses between 8 and 359 CPUs depending on the stage, with data
100 storage handled by user submitted information and simulation execution conducted manually.

101 GAMA is a biased genetic algorithm ²³. It first conducts two stages (Guess and Add) of only
102 non-essential gene deletions, which form a biased initial generation for the next (Mate) stage. The
103 latter follows a standard genetic algorithm process. GAMA produces deletion segments that vary by
104 individual genes, requiring 100s-1000s of CPUs. It takes two months to generate minimal genome
105 size reductions as it uses between 400 and 3000 CPUs depending on the stage. Custom
106 management code is used to coordinate and execute simulations, and store data.

107 **Initial Input**

108 To generate an initial input for Minesweeper and GAMA we simulated single-gene knockouts in an
109 otherwise unmodified *M.genitalium in-silico* genome (as previously reported ^{17,19}, Supplementary

110 Information A). The 359 protein-coding genes were simulated individually (10 replicates each), with
111 152 genes being classified as non-essential and 207 genes classified as essential (i.e. producing a
112 dividing or nondividing *in-silico* cell, respectively). The majority of genes (58%) are essential; this was
113 expected, as *Mycoplasmas* are obligate parasites with reduced genetic redundancy ²⁴.

114 318 genes showed consistent results across knockout replicates, the same phenotype in 10/10
115 cases, with 41 showing inconsistent results. Statistical analysis (binomial proportion confidence
116 interval, Pearson-Klopper, 95% CIs for: one 6/10 replicate [5.74, 6.87], 7/10 replicates [6.66, 7.93],
117 8/10 replicates [7.56, 8.97], 9/10 replicates[8.45, 9.99]) resulted in the genes being classified by the
118 majority phenotype (see Methods and Supplementary Information B & N). Overall, our results agree
119 97% with Karr et.al ¹⁷, see Supplementary Information C.

120 **Minesweeper Method and Results**

121 The first stage of Minesweeper conducts individual gene knockouts *in-silico* to identify complete/high
122 essentiality genes, removing them as gene deletion candidates.

123 The second stage sorts the remaining non-essential genes into deletion segments (from 12.5 to
124 100% of the remaining genes (Figure 1) resulting in 26 segments, broadly sweeping for potential low
125 essential genes. The deletion segments that produce a dividing *in-silico* cell are carried forward to
126 the next stage.

127 The third stage progresses with the largest deletion segment that produced a dividing cell, which is
128 matched with other dividing, non-overlapping segments. A powerset (all possible unique
129 combinations of the matched segments) is generated, and each combination of deletion segments is
130 simulated in an *in-silico* cell.

131 The fourth stage is cyclical. The largest deletion combination that produces a dividing cell is used to
132 generate a remaining gene list, those yet to be deleted, which narrows down potential conditional
133 essential genes. It splits the remaining genes into eight groups (see Methods) and a powerset is

134 generated. Each combination is individually appended to the current largest deletion combination
135 and simulated. Again, the largest deletion combination that produces a dividing cell is used to
136 generate a remaining gene list, which is used to start the next cycle of the stage.

137 If none of the combinations produces a dividing cell, the remaining genes are singly appended to the
138 largest deletion combination and simulated. The individual remaining genes that don't produce a
139 dividing cell are temporarily excluded and a reduced remaining gene list is produced, which is used
140 at the start of the next cycle.

141 The fourth stage continues until there are eight or less remaining genes (where a final appended
142 powerset is run) or all individually appended remaining genes do not produce a dividing cell. Both
143 outcomes result in a list of deleted genes and identified low essential genes.

144 Minesweeper produced results quickly, within two days the third stage removed 123 genes (a 34%
145 reduction) comparable to current lab-based efforts in other species^{3,4,6}. The repeating fourth stage
146 increased the overall number of deletions.

147 In total, Minesweeper deleted 145 genes (Figure 1), creating an *in-silico* *M.genitalium* cell containing
148 256 genes (named Minesweeper_256), which replicates DNA, produces RNA and protein, grows, and
149 divides.

150 **GAMA Method and Results**

151 The first and second stages of GAMA (Guess and Add) are pre-processing stages that provide input
152 for the third stage (Mate), a genetic algorithm. Typically a genetic algorithm would start with random
153 gene knockouts, but to reduce the number of generations required to produce minimal genome size
154 reductions, the Mate stage starts with large gene knockouts produced by Guess and Add (Figure 2).

155 In the first stage, Guess, all the non-essential genes from the initial input are segmented into four
156 sets, to reduce the size and number of combinations to search through. Each set is then used to

157 generate ~400 subsets, by randomly choosing combinations of 50 - 100% of the genes (~40) in the
158 set to delete. The build and test steps are then conducted. If a cell divides, the deletion subset is
159 labelled “viable” and carried forward to the next stage.

160 During the second stage, Add, a number of “viable” subsets are randomly selected from two, three
161 or four of the sets, which are combined into a larger set. Being able to select smaller numbers of
162 subsets reduces the chance of only producing non-dividing cells. ~3000 combined subsets are
163 created, simulated and tested. Those producing a dividing cell are ranked based on the number of
164 genes deleted. The 50 smallest genomes are taken forward to the mate stage.

165 During the third stage, Mate, the 50 smallest genomes are used to speed up the discovery of minimal
166 genomes. The mate stage is cyclical, consisting of generations containing 1000 simulations. Each
167 simulation in a generation combines two of the 50 smallest *in-silico* genomes at random, and
168 introduces random gene knockouts and knock-ins from a pool of all protein-coding genes (including
169 complete and high essentiality genes). The genomes produced are ranked and compared to the
170 smallest 50 genomes, with the new smallest 50 being carried through to the next generation. The
171 mate step automatically stops after 100 generations, but was manually stopped at 46 generations,
172 after 20 generations without producing a smaller genome.

173 In total, the smallest GAMA-reduced *in-silico* genome deleted 165 genes, creating an *in-silico*
174 *M.genitalium* genome of 236 genes (named GAMA_236). GAMA removed more genes than the
175 Minesweeper method, while still producing a simulated cell which replicates DNA, produces RNA
176 and protein, grows, and divides.

177 **GAMA_236 and Minesweeper_256 Genomes**

178 We investigated the characteristics of our two minimal genomes in terms of how consistently they
179 produced a dividing *in-silico* cell, and the range of possible behaviour they displayed. We simulated
180 100 replicates of an unmodified *M.genitalium in-silico* genome, Minesweeper_256, GAMA_236, and a
181 single-gene knockout of a known essential gene (MG_006) to provide a comparison (see

182 Supplementary Information G). The rate of division (or not in the MG_006 knockout simulations) was
183 analysed to assign a phenotype penetrance percentage, quantifying how often an expected
184 phenotype occurred. The unmodified *M.genitalium* and MG_006 knockout *in-silico* genomes
185 demonstrated consistent phenotypes (99% and 0% divided, respectively). Minesweeper_256 was
186 slightly less consistent (89% divided), while GAMA_236 was substantially less consistent, producing
187 a dividing *in-silico* cell 18% of the time. This is not entirely unexpected given the greater number of
188 gene deletions affecting essential gene functions (according to the GO term analysis).

189 The 100 replicates for the unmodified *M.genitalium* genome, Minesweeper_256, and GAMA_236
190 were plotted to assess the range of behaviour (Figure 3). The unmodified *M.genitalium* whole-cell
191 model (Figure 3, top row) shows the range of expected behaviour for a dividing cell (in line with
192 previous results¹⁷). Growth, protein production, and cellular mass increase over time, with most cells
193 dividing at around 10 hours, though division can occur between 6 and 11 hours. RNA production
194 fluctuates but increases over time. DNA replication follows a characteristic shape, with some
195 simulations delaying the initiation of DNA replication past ~9 hours.

196 By comparison, Minesweeper_256 (Figure 3, middle row) displays slower, and in some cases
197 decreasing, growth over time which is capped to a lower maximum. Protein production and cellular
198 mass are generated more slowly and present some erratic behaviour. The range of RNA production is
199 narrower compared to the unmodified *M.genitalium* whole-cell model. DNA replication takes longer
200 and initiation can occur later (at 11 hours). Cell division occurs later, between 8 and 13.889 hours. A
201 number of simulations can be seen failing to replicate DNA and divide.

202 Compared to the other genomes, GAMA_236 (Figure 3, bottom row) shows a much greater range of
203 growth rates. Some grow as fast as the unmodified genome, some are comparable to
204 Minesweeper_256, and some show very low or decreasing growth. Observable protein levels appear
205 between 2 and 5 hours, followed by a slower rate of protein production in some simulations. Cellular
206 mass is either similar to Minesweeper_256 or slower. The range of RNA production is reduced and
207 the rate of RNA production is slower.

208 Some simulations replicate DNA at a rate comparable to the unmodified genome, others replicate
209 more slowly, and some do not complete DNA replication. Cell division occurs across a greater range
210 (6 - 13.889 hours). A number of simulations showing metabolic defects can be seen. These do not
211 produce any growth, and can also be seen failing to replicate DNA and divide.

212 We investigated what processes were removed in the creation of Minesweeper_256, using gene
213 ontology (GO) biological process terms (see Methods and Supplementary Information I-K). The
214 baseline *M.genitalium* whole-cell model has 259 genes of 401 genes (72% coverage) with GO terms
215 on UniProt²⁵. Minesweeper_256 has 186 (73%) genes with GO terms and 70 (27%) genes without.
216 The 140 gene deletions reduced 22 (14%) GO categories, and removed 41 (27%) GO categories
217 entirely, of which 29 (70%) were associated with a single gene (see Supplementary Information L).

218 The GO categories reduced include: DNA (replication, topological change, transcription regulation
219 and initiation); protein (folding and transport); RNA processing; creation of lipids; cell cycle; and cell
220 division. As the *in-silico* cells continue to function, we can assume that these categories could
221 withstand low-level disruption.

222 Removed GO categories that involved multiple genes include: proton transport; host interaction; DNA
223 recombination and repair; protein secretion and targeting to membrane; and response to oxidative
224 stress.

225 Removed GO categories that contain single genes include: transport (proton, carbohydrate,
226 phosphate and protein import, protein insertion into membrane); protein modification (refolding,
227 repair, targeting); chromosome (segregation, separation); biosynthesis (coenzyme A, dTMP, dTTP,
228 lipoprotein); breakdown (deoxyribonucleotide, deoxyribose, mRNA, protein); regulation (phosphate,
229 carbohydrate, and carboxylic acid metabolic processes, cellular phosphate ion homeostasis);
230 cell-cell adhesion; foreign DNA cleavage; SOS response; sister chromatid cohesion; and uracil
231 salvage.

232 These deletions reduce the ability of *M.genitalium* to interact with the environment and defend
233 against external forces. This results in a reduction in control, from transport to regulation to genome

234 management, and pruned metabolic processes and metabolites. This leaves Minesweeper_256's
235 *in-silico* cell alive, but more vulnerable to external and internal pressures, less capable of responding
236 to change, and more reliant on internal processes occurring by chance.

237 In comparison, GAMA_236 has 163 genes (69% coverage) with GO terms on UniProt ²⁵, with 73
238 genes with no GO terms. The 165 genes deleted reduced 17 (11%) GO categories, and removed 55
239 (35%) GO categories, 38 (69%) of which were associated with a single-gene (see Supplementary
240 Information M).

241 8 unaffected and five reduced GO categories in Minesweeper_256 were removed in GAMA_236, with
242 one unaffected GO category unique to GAMA_236 (phosphate ion transmembrane transport). Four
243 GO categories were reduced further in GAMA_236: DNA (transcription, transcription regulation,
244 transport) and glycerol metabolic process.

245 The 13 additional GO categories removed include: DNA (transcription (termination, regulation of
246 elongation, antitermination, initiation)); RNA (processing (mRNA, tRNA, rRNA), rRNA catabolic
247 process, tRNA modification, pseudouridine synthesis); thiamine (biosynthetic process, diphosphate
248 biosynthetic process); and protein lipoylation.

249 GO analysis of GAMA_236, when compared to Minesweeper_256, suggests a further reduction of
250 both internal control and reactivity to external environment.

251 **Genes with Low and High Essentiality**

252 We analysed Minesweeper_256 and GAMA_236 to determine whether these were different minimal
253 genomes, or GAMA_236 was an extension of Minesweeper_256. We conducted a gene content
254 comparison of an unmodified *M.genitalium*, Minesweeper_256, and GAMA_236 genomes (Figure 4,
255 Supplementary Information F), highlighting gene deletions unique to each minimal genome. We took
256 this a step further and compared Minesweeper_256 to all of the GAMA genomes 256 to 236 genes in
257 size. Figure 5 shows the GAMA algorithm's avenue of gene reductions converging to a minimal
258 genome, but Minesweeper_256 is not on the same path of convergence.

259 Our comparison of the genomes found 18 genes knocked out in GAMA_236 that have high
260 essentiality ¹. They were defined as essential by single knockout in an unmodified *M.genitalium*
261 whole-cell model, but could be removed in the genomic context of GAMA_236 without preventing
262 division (see Supplementary Information A & E). We also found that four of these 18 genes could be
263 removed as a group in the genomic context of Minesweeper_256, but doing so greatly increased the
264 number of non-dividing cells produced (see Supplementary Information E).

265 Our genome comparison also found that Minesweeper_256 removed four genes, and GAMA_236
266 removed five genes (Table 1), which could not be removed either individually or as a group from its
267 counterpart, without causing cellular death or mutations that prevented cellular division. We
268 confirmed that these nine genes were individually non-essential. One additional gene, MG_305,
269 could not be additionally removed in both GAMA_236 and Minesweeper_256. Our results
270 demonstrate that these nine genes have low essentiality ¹. To identify the cause of this synthetic
271 lethality we attempted to match the functions of these low essentiality genes (Table 1), as we
272 anticipated finding redundant essential gene pairs or groups. We found two genes in GAMA_236
273 (MG_289, MG_291) had matching GO terms with the gene MG_411 in Minesweeper_256. These, and
274 three other adjacent genes on the genome, were tested by combinatorial gene knockouts in an
275 unmodified *M.genitalium* whole-cell model genome (see Supplementary Information H). MG_289,
276 MG_290, MG_291 were found to form a functional group, as were MG_410, MG_411, MG_412.
277 These genes could be deleted individually and in functional groups from an otherwise unmodified
278 *M.genitalium* whole-cell genome, and produce a dividing *in-silico* cell. However, any double gene
279 deletion combination that involved one gene from each functional group resulted in a cell that could
280 not produce RNA, produce protein, replicate DNA, grow or divide.

281 *M.genitalium* only has two external sources of phosphate, inorganic phosphate and phosphonate.
282 MG_410, MG_411, and MG_412 transport inorganic phosphate into the cell, with MG_289, MG_290,
283 and MG_291 transporting phosphonate into the cell ^{18,26}. These phosphate sources proved to be a
284 key difference between our minimal genomes. Minesweeper_256 removed the phosphate transport

285 genes, relying on phosphonate as the sole phosphate source. GAMA_236 removed the phosphonate
286 transport genes, relying on inorganic phosphate as the sole phosphate source. This can be seen in
287 the GO term analysis, the phosphate ion transmembrane transport is still present in GAMA_236 but
288 not in Minesweeper_256.

289 It has previously been theorised that individual bacterial species will have multiple minimal genomes
290 ^{27,28}, with different gene content depending on the environment and which evolutionary redundant
291 cellular pathways were selected during reduction. We would argue that one of these selected
292 pathways is phosphate source, with minimal genomes differing by choice of phosphate transport
293 genes and associated processing stages, equivalent to the *phn* gene cluster in *Escherichia coli* ²⁹. We
294 could not however find any annotated phosphonate processing genes that had been subsequently
295 removed in GAMA_236. We suspect that further “pivot points”, the selection of one redundant
296 cellular pathway over another during reduction, will be identified in future *in-vivo* and *in-silico*
297 bacterial reductions increasing the base number of minimal genomes per bacterial species.

298 Discussion

299 We created two genome design algorithms (Minesweeper and GAMA) that used computational
300 design-simulate-test cycles to produce *in-silico* *M.genitalium* minimal genomes (achieving 36% and
301 41% reductions, respectively). Our minimal genomes are smaller than *JCVI-syn3.0* (currently the
302 smallest genome that can be grown in pure culture ²) and 33 - 53 genes smaller than the most recent
303 predictions for a reduced *Mycoplasma* genome ¹⁶.

304 Additionally, we identified 10 low essentiality genes, 18 high essentiality genes ¹, and produced
305 evidence for at least two minima for *Mycoplasma genitalium in-silico*. We plan to test these results
306 experimentally to ascertain the accuracy of the model and the functionality of our minimal genomes.

307 We believe that single-gene knockout classifications are unreliable for genome minimisation, as they
308 fail to take into account genomic context. Single-gene knockout studies will underestimate minimal
309 genome size as low essentiality genes will be scored as non-essential ^{2,15,16}, but they will also
310 overestimate minimal genome size as high essentiality genes will be scored as essential. We found
311 10 low essential genes within 358 protein-coding genes. As a single synthetic lethality event will
312 prevent a genome from surviving, this gives a 3% chance of error for untested genome designs in
313 even this evolutionarily reduced genome. Additionally, single-gene knockout studies narrow the
314 scope of genome design; the 18 high essentiality genes identified as dispensable within GAMA_236
315 would not have been traditionally targeted by laboratory methods.

316 There are limitations to the approach presented here. Models are not perfect representations of
317 reality: through necessity this model bases some of its parameters on data from other bacteria ¹⁷;
318 multi-generation simulations are only possible by isolating one submodel from the rest of model
319 (which loses genomic context); and *M.genitalium* has genes of unknown function that the model
320 cannot account for.

321 The success of our *in-silico* genomes *in-vivo* is dependent on the accuracy of the model, which is
322 untested at this scale of genetic modification. Minesweeper_256 and GAMA_236 may only function
323 in the first generation of cells and the impact of the unmodelled genes is unknown. These genes may

324 change the genomic context such that our minimal genomes are not successful, or as found with
325 JCVI-Syn3.0 the genes of unknown function will be required for viability ².

326 Our algorithms are currently adaptable to future, under development whole-cell models, as the
327 algorithms interact with the models only via the input of gene deletion lists and analysing the output.
328 This includes the *E.coli* whole-cell model at the Covert Lab, Stanford and the *Mycoplasma*
329 *pneumoniae* whole-cell model at the Karr Lab, Mount Sinai, New York ³⁰.

330 We believe that a hybrid of computational and lab based genome design and construction is now
331 possible. This could produce quicker and cheaper laboratory results than currently possible, opening
332 up this research to broader and interdisciplinary research communities. It also expands our research
333 horizons raising the possibility of building truly designer cells, with increased efficiency and functional
334 understanding.

335 **Methods**

336 **Model Availability**

337 The *M.genitalium* whole-cell model is freely available: <https://github.com/CovertLab/WholeCell>. The
338 model requires a single CPU and can be run with 8GB of RAM. We run the *M.genitalium* whole-cell
339 model on Bristol's supercomputers using MATLAB R2013b, with the model's standard settings.
340 However, we use our own version of the SimulationRunner.m. MGGRunner.m is designed for use
341 with supercomputers that start hundreds of simulations simultaneously, artificially incrementing the
342 time-date value for each simulation, as this value is subsequently used to create the initial conditions
343 of the simulation. This incrementation prevents the running of multiple simulations with identical initial
344 conditions.

345 Our research copy of the whole-cell model was downloaded 2017-01-10.

346 **Code Availability**

347 The code used for this research is openly available on Github (public code provided on publication).
348 This includes the code for Minesweeper and GAMA genome design tools, scripts for statistical
349 analysis, scripts for analysing GO terms, our custom simulation runner, analysis scripts, a template
350 bash script, as well as the bash scripts and text files used to generate the simulations in this paper.

351 **Statistics**

352 We used the R binom package (<https://www.rdocumentation.org/packages/binom>) to conduct
353 one-tailed binomial proportion confidence intervals on our 41 genes showing inconsistent results
354 (success ranging from 6 to 9 replicates, out of a total of 10 replicates). We used binom.confint.exact
355 (Pearson-Klopper) using 95% CIs, producing for: 6/10 replicates [0.26, 0.87], 7/10 replicates [0.34,
356 0.93], 8/10 replicates [0.44, 0.97], 9/10 replicates [0.55, 0.99]. We graphed these results in R and in
357 Python using Seaborn (<https://seaborn.pydata.org/>), the exact values, code, and graphs produced
358 are available in Supplementary Information B & N.

359 Figure 5 was generated by creating a similarity matrix between all of the 2955 genomes, with the
360 gene information represented in a binary format (present or absent). The matrix calculated a distance
361 metric (1 - Adjusted Rand Index), with each genome comparison given a normalised score (0 = the

362 genomes were identical, 1 = as different as would be expected if each genome was generated
363 randomly, 2 = completely different). The resulting 2955 x 2955 matrix was then reduced to two
364 dimensions with a standard PCA.

365 **Minesweeper**

366 Minesweeper is written in Python3 and consists of four scripts (one for each stage). It uses no
367 external libraries, so should be able to be run on any modern operating system (as they come with
368 Python preinstalled) via a terminal. Each stage/script requires a text file(s) as input, with each stage
369 outputting simulation files. These are run on a supercomputer and the automatically produced
370 summary file is used as input for the next stage. Stages one to three are sequential, with stage four
371 repeating until Minesweeper stops. Detailed instructions are provided in the README and progress
372 is recorded in the deletion log in /OUTPUT_final.

373 The first stage of Minesweeper is optional, if you already have single gene knockout simulation
374 results, you can proceed to the second stage. The second stage creates 26 deletion segments:
375 100%, 90%A, 90%B, 80%A, 80%B, 70%A, 70%B, 60%A, 60%B, 50%A, 50%B, 33%A-C, 25%A-D,
376 12.5%A-H. The A segments start from the top of the list of genes, whereas the B segments start
377 from the bottom of the list of genes. The third stage progresses with the three largest deletion
378 segments that produced a dividing cell, these three variants are referred to as red, yellow, blue.
379 These perform as replicates and as a check on if the results are converging. The three variants are
380 matched with smaller, dividing, non-overlapping segments using a list of allowed matches
381 (implementation is detailed in third stage script), and unique combinations generated using a python
382 implementation of powersets. The fourth stage splits the remaining genes into eight groups. The
383 reason for selecting eight groups and three variants, is that a set of eight produces 256 unique
384 combinations. Three variants each with 256 simulations (768 total) is 85% of the capacity of
385 BlueGem. A set of nine groups with three variants (1536 simulations total) is 170% the capacity of
386 BlueGem. Queueing systems mean that you don't require this number of CPUs in total, but the
387 execution time is multiplied as you wait for the simulations to process. The number of variants and
388 groups can be lowered or increased depending on the number of CPUs you have available.

389 **GAMA**

390 GAMA is written in Python3 and relies on a variety of different packages. These dependencies can be
391 easily taken care of by installing it from PyPI using either 'pip install genome_design_suite' or 'conda
392 install genome_design_suite' (it is recommended that you do this from within a virtual environment
393 since this is pre-alpha and has not been extensively tested with different versions of all the libraries).
394 A dependencies list is available in the main directory of the github repository if you would like to do
395 this manually. The main dependency is the 'genome_design_suite' which is a suite of tools created
396 by Oliver Chalkley at the University of Bristol which enables it to be easily run on different (or even
397 multiple) clusters and well as enabling automatic data processing and database management. Due to
398 the large amount of data produced by the Whole-Cell model, the simulation output data was reduced
399 to essential data, converted into Pandas DataFrames (<https://pandas.pydata.org/>) and saved in
400 Pickle files. GAMA would have produced 100s of TBs of data in the model's native output format
401 (compressed matlab files) which we are not able to store so this was an essential step. In order to run
402 this code you must have a computer dedicated to remotely managing the simulations. A PC with a
403 quad-core Intel(R) Xeon(R) CPU E5410 (2.33GHz) and 1GB of RAM running CentOS-6.6 was used as
404 our computer manager, which is referred to as OC2. GAMA was run on OC2 using the scripts
405 contained in gama_manangement.zip Each stage of GAMA was run individually and manually
406 updated as it was in proof-of-concept stage when GAMA_236 was found. ko.db is an SQLite3
407 database used to stored key information about simulations like average growth rate and division
408 time.

409 The guess stage splits the singularly non-essential genes in roughly equally sized partitions. The four
410 files, focus_on_NE_split_[1-4].py, run the exploration of each of the four partitions of the guess stage
411 from OC2 - after unzipping gama_management.zip these can be found in gama/guess. The
412 submission scripts and other files automatically created to run the simulations on the cluster can be
413 found in gama_run_files.zip -> gama_run_files/guess. The simulation output is saved in Pickle files
414 and can be found in gama_data/guess. Due to a technical problem the growth rate and division time
415 of the genomes simulated in this stage are not in ko.db. viability_of_ne_focus_sets_pickles.zip
416 contains the viability data of these simulations and the Python script used to collect it.

417 The add stage was executed on OC2 by running the files in gama_management.zip -> gama/add.

418 The submission scripts and other files automatically created to run the simulations on the cluster can
419 be found in gama_run_files.zip -> gama_run_files/add. The simulation output can be found in
420 gama_data/add and an overview of the simulation results can be found in ko.db where the
421 batchDescription.name is some derivative of 'mix_ne_focus_split'.

422 The mate stage was executed on OC2 by running the file in gama_management.zip -> gama/mate.

423 The submission scripts and other files automatically created to run the simulations on the cluster can
424 be found in gama_run_files.zip -> gama_run_files/mate. The simulation output can be found in
425 gama_data/mate and an overview of the simulation results can be found in ko.db where
426 batchDescription.name is some derivative of 'big_mix_of_split_mixes'.

427 **Equipment**

428 We used the University of Bristol Advanced Computing Research Centres's BlueGem, a 900-core
429 supercomputer, which uses the Slurm queuing system, to run whole-cell model simulations. GAMA
430 also used BlueCrystal, a 3568-core supercomputer, which uses the PBS queuing system.

431 We used a standard office desktop computer, with 8GB of ram, to write new code, interact with the
432 supercomputer, and run single whole-cell model simulations. We used the following GUI software on
433 Windows/Linux Cent OS: Notepad++ for code editing, Putty (ssh software)/the terminal to access the
434 supercomputer, and FileZilla (ftp software) to move files in bulk to and from the supercomputer. The
435 command line software we used included: VIM for code editing, and SSH, Rsync, and Bash for
436 communication and file transfer with the supercomputers.

437 **Data Format**

438 The majority of output files are state-NNN.mat files, which are logs of the simulation split into
439 100-second segments. The data within a state-NNN.mat file is organised into 16 cell variables, each
440 containing a number of sub-variables. These are typically arranged as 3-dimensional matrices or time

441 series, which are flattened to conduct analysis. The other file types contain summaries of data
442 spanning the simulation.

443 **Data Analysis Process**

444 The raw data is automatically processed as the simulation ends. runGraphs.m carries out the initial
445 analysis, while compareGraphs.m overlays the output on collated graphs of 200 unmodified
446 *M.genitalium* simulations. Both outputs are saved as MATLAB .fig and .pdfs, though the .fig files
447 were the sole files analysed. The raw .mat files were stored in case further investigation was required.
448 To classify our data we chose to use the phenotype classification previously outlined by Karr (Figure
449 6B¹⁷), which graphed five variables to determine the simulated cells' phenotype. However, the script
450 responsible for producing Figure 6B, SingleGeneDeletions.m, was not easily modified. This led us to
451 develop our own analysis script recreating the classification: runGraphs.m graphs growth, protein
452 weight, RNA weight, DNA replication, cell division, and records several experimental details. There
453 are seven possible phenotypes caused by knocking out genes in the simulation: non-essential if
454 producing a dividing cell; and essential if producing a non-dividing cell because of a DNA replication
455 mutation, RNA production mutation, protein production mutation, metabolic mutation, division
456 mutation, or slow growing.

457 For the single gene knockout simulations produced in initial input, the non-essential simulations were
458 automatically classified and the essential simulations flagged. Each simulation was investigated
459 manually and given a phenotype manually using the decision tree (see Supplementary Information D).

460 For simulations conducted by Minesweeper and GAMA, simulations were automatically classified
461 solely by division, which can be analysed from cell width or the endtime of the simulation.

462 Further analysis, including: cross-comparison of single-gene knockout simulations, comparison to
463 Karr et al's¹⁷ results, analysis of Minesweeper and GAMA genomes (genetic content and similarity,
464 behavioural analysis, phenotypic penetrance, gene ontology), and identification and investigation of
465 high and low essentiality genes and groupings, were completed manually. The GO term analysis of
466 gene deletion impacts was processed by a created script (see Github for code), then organised into
467 tables of GO terms that were unaffected, reduced, or removed entirely.

468 **Modelling: Scripts, Process and Simulations**

469 Generally, there are six scripts we used to run the whole-cell model. Three are the experimental files
470 created with each new experiment (the bash script, gene list, experiment list), and three are stored
471 within the whole-cell model and are updated only upon improvement (MGRunner.m, runGraphs.m,
472 and compareGraphs.m). The bash script is a list of commands for the supercomputer(s) to carry out.
473 Each new bash script is created from the GenericScript.sh template, which determines how many
474 simulations to run, where to store the output, which analysis to run, and where to store the results of
475 the analysis. The gene list is a text file containing rows of gene codes (in the format 'MG_XXX'). Each
476 row corresponds to a single simulation and determines which genes that simulation should knockout.
477 The experiment list is a text file containing rows of simulation names. Each row corresponds to a
478 single simulation and determines where the simulation output and results of the analysis are stored.
479 In brief, to manually run the whole-cell model: a new bash script, gene list, and experiment list are
480 created on the desktop computer to answer an experimental question. The supercomputer is
481 accessed on the desktop via ftp software, where the new experimental files are uploaded, the
482 planned output folders are created, and MGRunner.m, runGraphs.m, compareGraphs.m files are
483 confirmed to be present. The supercomputer is then accessed on the desktop via ssh software,
484 where the new bash script is made executable and added to the supercomputer's queuing system to
485 be executed. Once the experiment is complete, the supercomputer is accessed on the desktop via
486 ssh software, where the results of the analysis are moved to /pdf and /fig folders. These folders are
487 accessed on the desktop via ftp software, where the results of the analysis are downloaded. More
488 detailed instructions are contained within the template bash script.
489 Each wild-type simulation consists of 300 files requiring 0.3GB. Each gene manipulated simulation
490 can consist of up to 500 files requiring between 0.4GB and 0.9GB. Each simulation takes 5 to 12
491 hours to complete in real time, 7 - 13.89 hours in simulated time.

492 **Data Availability**

493 The databases used to design our *in-silico* experiments, and compare our results to, includes Karr et
494 al¹⁷ and Glass et al²⁴ Supplementary Information, and Fraser et al *M.genitalium* G37 genome¹⁸
495 interpreted by KEGG²⁶ and UniProt²⁵ as strain ATCC 33530/NCTC 10195.
496 Minesweeper simulations raw and transformed output (.mat files) are available upon request, as the
497 they require 4.2 TB of storage. The output .fig files (10 GB) are available for download from the our
498 group's Research Data Repository at the University of Bristol. GAMA simulations transformed output
499 is available in ko.db.

500 **References**

- 501 1. Rancati, G., Moffat, J., Typas, A. & Pavelka, N. Emerging and evolving concepts in gene
502 essentiality. *Nat. Rev. Genet.* **19**, 34–49 (2018).
- 503 2. Hutchison, C. A. *et al.* Design and synthesis of a minimal bacterial genome. *Science* **351**,
504 1414–U73 (2016).
- 505 3. Iwadate, Y., Honda, H., Sato, H., Hashimoto, M. & Kato, J.-I. Oxidative stress sensitivity of
506 engineered *Escherichia coli* cells with a reduced genome. *FEMS Microbiol. Lett.* **322**, 25–33
507 (2011).
- 508 4. Hirokawa, Y. *et al.* Genetic manipulations restored the growth fitness of reduced-genome
509 *Escherichia coli*. *J. Biosci. Bioeng.* **116**, 52–58 (2013).
- 510 5. Zhou, J., Wu, R., Xue, X. & Qin, Z. CasHRA (Cas9-facilitated Homologous Recombination
511 Assembly) method of constructing megabase-sized DNA. *Nucleic Acids Res.* **44**, e124 (2016).
- 512 6. Reuß, D. R. *et al.* Large-scale reduction of the *Bacillus subtilis* genome: consequences for the
513 transcriptional network, resource allocation, and metabolism. *Genome Res.* **27**, 289–299 (2017).
- 514 7. Gibson, D. G. *et al.* Complete chemical synthesis, assembly, and cloning of a *Mycoplasma*
515 *genitalium* genome. *Science* **319**, 1215–1220 (2008).
- 516 8. Lartigue, C. *et al.* Genome transplantation in bacteria: changing one species to another. *Science*
517 **317**, 632–638 (2007).
- 518 9. Baby, V. *et al.* Cloning and Transplantation of the *Mesoplasma florum* Genome. *ACS Synth. Biol.*
519 **7**, 209–217 (2018).
- 520 10. Lartigue, C. *et al.* Creating bacterial strains from genomes that have been cloned and engineered
521 in yeast. *Science* **325**, 1693–1696 (2009).
- 522 11. Szostak, J. W., Bartel, D. P. & Luisi, P. L. Synthesizing life. *Nature* **409**, 387–390 (2001).
- 523 12. Qian, W., Ma, D., Xiao, C., Wang, Z. & Zhang, J. The genomic landscape and evolutionary
524 resolution of antagonistic pleiotropy in yeast. *Cell Rep.* **2**, 1399–1410 (2012).
- 525 13. Commichau, F. M., Pietack, N. & Stülke, J. Essential genes in *Bacillus subtilis*: a re-evaluation
526 after ten years. *Mol. Biosyst.* **9**, 1068–1075 (2013).
- 527 14. Ryan, C. J., Krogan, N. J., Cunningham, P. & Cagney, G. All or nothing: protein complexes flip

- 528 essentiality between distantly related eukaryotes. *Genome Biol. Evol.* **5**, 1049–1059 (2013).
- 529 15. Dobzhansky, T. Genetics of natural populations; recombination and variability in populations of
530 *Drosophila pseudoobscura*. *Genetics* **31**, 269–290 (1946).
- 531 16. Glass, J. I., Merryman, C., Wise, K. S., Hutchison, C. A., 3rd & Smith, H. O. Minimal Cells-Real
532 and Imagined. *Cold Spring Harb. Perspect. Biol.* (2017). doi:10.1101/cshperspect.a023861
- 533 17. Karr, J. R. *et al.* A whole-cell computational model predicts phenotype from genotype. *Cell* **150**,
534 389–401 (2012).
- 535 18. Fraser, C. M. *et al.* THE MINIMAL GENE COMPLEMENT OF MYCOPLASMA-GENITALIUM.
536 *Science* **270**, 397–403 (1995).
- 537 19. Sanghvi, J. C. *et al.* Accelerated discovery via a whole-cell model. *Nat. Methods* **10**, 1192–1195
538 (2013).
- 539 20. Purcell, O., Jain, B., Karr, J. R., Covert, M. W. & Lu, T. K. Towards a whole-cell modeling
540 approach for synthetic biology. *Chaos* **23**, 025112 (2013).
- 541 21. Kazakiewicz, D., Karr, J. R., Langner, K. M. & Plewczynski, D. A combined systems and
542 structural modeling approach repositions antibiotics for *Mycoplasma genitalium*. *Comput. Biol.*
543 *Chem.* **59 Pt B**, 91–97 (2015).
- 544 22. Knuth, D. E. *The Art of Computer Programming, Volume 3: (2Nd Ed.) Sorting and Searching*.
545 (Addison Wesley Longman Publishing Co., Inc., 1998).
- 546 23. Back, T. & Bäck, T. *Evolutionary Algorithms in Theory and Practice: Evolution Strategies,*
547 *Evolutionary Programming, Genetic Algorithms*. (Oxford University Press, 1996).
- 548 24. Glass, J. I. *et al.* Essential genes of a minimal bacterium. *Proc. Natl. Acad. Sci. U. S. A.* **103**,
549 425–430 (2006).
- 550 25. Apweiler, R. *et al.* UniProt: the Universal Protein knowledgebase. *Nucleic Acids Res.* **32**, D115–9
551 (2004).
- 552 26. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*
553 **28**, 27–30 (2000).
- 554 27. Burgard, A. P., Vaidyaraman, S. & Maranas, C. D. Minimal reaction sets for *Escherichia coli*
555 metabolism under different growth requirements and uptake environments. *Biotechnol. Prog.* **17**,
556 791–797 (2001).

- 557 28. Huynen, M. Constructing a minimal genome. *Trends Genet.* **16**, 116 (2000).
- 558 29. Kamat, S. S., Williams, H. J. & Raushel, F. M. Intermediates in the transformation of
559 phosphonates to phosphate by bacteria. *Nature* **480**, 570–573 (2011).
- 560 30. Karr, J. Comprehensive whole-cell computational models of individual cells. *Whole-Cell*
561 *Modeling* (2018). Available at: <https://www.wholecell.org/models/>. (Accessed: 6th November
562 2018)

563 **Acknowledgments**

564 We would like to thank the Advanced Computing Research Centre (ACRC) and BrisSynBio, a
565 BBSRC/EPSRC Synthetic Biology Research Centre, at the University of Bristol for access to the
566 BlueCrystal and Bluegem supercomputers. Special thanks to the HPC and RDSF teams of the
567 ACRC, particularly Dr. Christopher Woods, Simon Burbidge, Matt Williams, and Damian Steer for
568 their help with BlueCrystal, BlueGem, data storage and publication.

569 We would like to thank Jonathan Karr for his advice on running gene knockout simulations using the
570 *M. genitalium* whole-cell model, and for his constructive and enlightening feedback.

571 We would like to thank Anthony Vecchiarelli (Assistant Professor, University of Michigan) and class
572 MCDB 401 (“Building the Synthetic Cell”) for conducting a class review of our preprint paper,
573 providing us with constructive and encouraging feedback.

574 We would like to thank John Glass (Professor, JCVI Synthetic Biology and Bioenergy Group) for his
575 constructive and informative feedback.

576 We would like to thank Cameron Matthews and Julia Needham, University of Bristol undergraduates,
577 who conducted simulations to test the inclusion of MG_290 within the phosphonate group as part of
578 summer studentships.

579 LM is supported by the Medical Research Council grant MR/N021444/1 to LM, and by the
580 Engineering and Physical Sciences Research Council grant EP/R041695/1 to LM.

581 OC, LM and CG are supported by a BrisSynBio, a BBSRC/EPSRC Synthetic Biology Research
582 Centre (BB/L01386X/1), flexi-fund grant.

583 OC is supported by the Bristol Centre for Complexity Sciences (BCCS) Centre for Doctoral Training
584 (CDT) EP/I013717/1; JR and SL are supported by EPSRC Future Opportunity Scholarships.

585 **Author Contributions**

586 C.G, L.M, O.C for attaining initial funding.

587 C.G, L.M, O.P, O.C, J.R, S.L were involved in ideation.

588 S.L was involved in analysis and development of Figure 4.

589 O.C was responsible for the development and implementation of the *Mycoplasma genitalium*
590 whole-cell model outside of the Covert Lab, Stanford (on Bristol’s BlueGem and BlueCrystal), initial

591 ideation about uses of whole-cell models, GAMA (method, results, section), Figure 2, Figure 5, and
592 collaborative theorising on essentiality and minimal genomes.

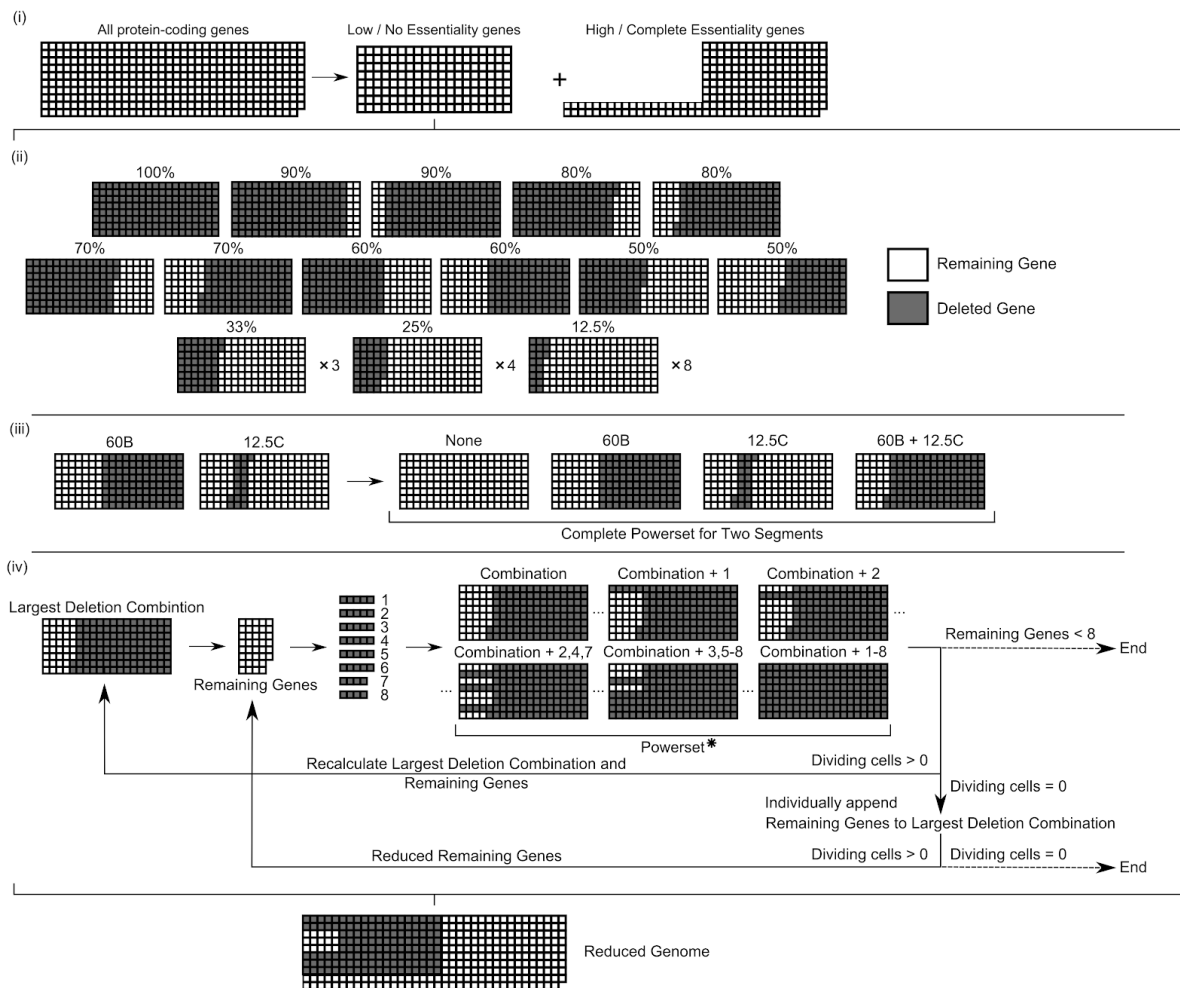
593 J.R was responsible for the development of automated graphing, Minesweeper (method, results),
594 spreadsheet analysis of *in-silico* results, Table 1, Figure 1, Figure 3, collaborative theorising on
595 essentiality and minima, and writing of the paper and Supplementary Information.

596 C.G, L.M, O.C, O.P, S.L were involved in editing and feedback on paper.

597 **Competing Interests**

598 The authors declare no competing interests.

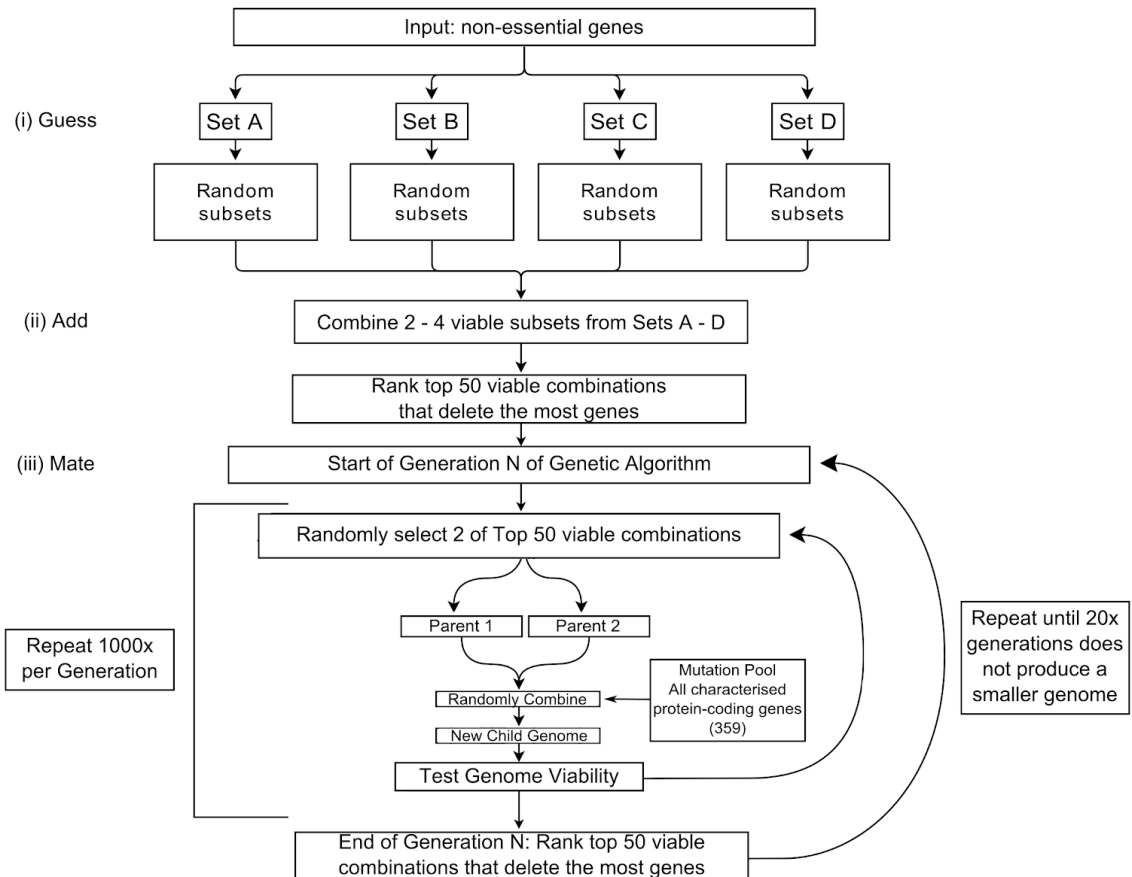
599 **Figures**



600 **Figure 1. Minesweeper Algorithm for Genome Design**

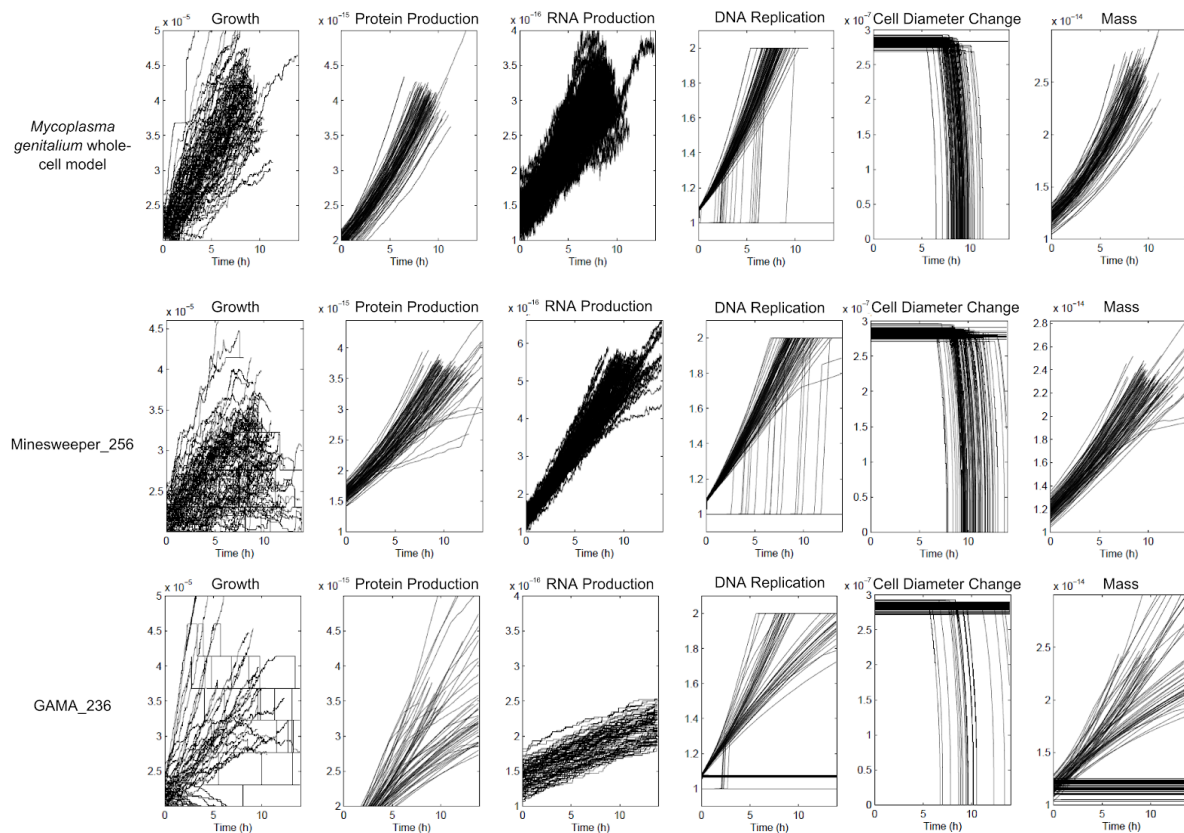
601 (i) *in-silico* single gene knockouts are conducted to identify low / no essentiality genes (whose
 602 knockout does not prevent cell division). (ii) 26 deletion segments, ranging in size from 100% to
 603 12.5% of the low / no essentiality genes, are simulated. Grey indicates a gene deletion, white
 604 indicates a remaining gene. Deletion segments that do not prevent division go to the next stage. (iii)
 605 The largest deletion segment is matched with all dividing, non-overlapping segments. A powerset (all
 606 possible unique combinations of this set of matched deletion segments) is generated and each
 607 combination simulated. Deletion segments that do not prevent division go to the next stage. (iv) The
 608 largest deletion segment determines the remaining low / no essentiality genes that have been
 609 deleted. These remaining genes are divided into eight groups (see Methods), a powerset generated

610 for these eight groups, and each member of the powerset individually appended to the current
611 largest deletion combination and simulated. If none of these simulations produces a dividing cell, the
612 remaining genes are appended as single knockouts to the current largest deletion combination and
613 simulated. The individual remaining genes that don't produce a dividing cell are temporarily excluded
614 and a reduced remaining gene list produced. Details of simulations settings are available in the
615 Methods. Powerset* = the complete powerset is not displayed here.



616 **Figure 2. GAMA Algorithm for Genome Design**

617 (i) Only non-essential genes whose knockout does not prevent cell division
 618 and are equally divided into Sets A - D. 400 random subsets are produced and simulated per set,
 619 each containing 50-100% of the genes within the set. Deletion segments that do not prevent division
 620 (“viable”) go to the next stage. (ii) 3000 combinations are generated and simulated. (iii) Is a cyclical
 621 step. The mutation pool targets a random number of genes for alteration (both knockins and
 622 knockouts), including essential genes. Details of simulations settings are available in the Methods.

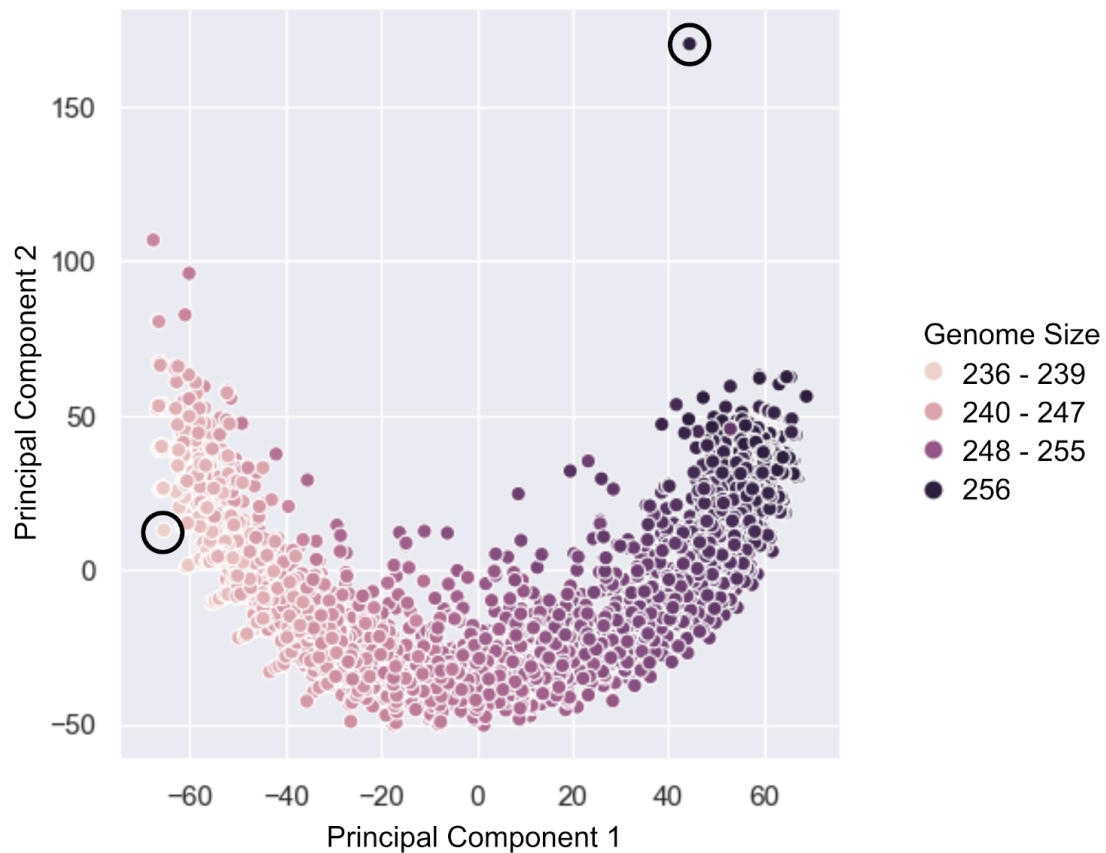


623 **Figure 3. Comparison of unmodified *Mycoplasma genitalium* whole-cell model,**
624 **Minesweeper_256, and GAMA_236 outputs**
625 100 *in-silico* replicates, with second-by-second values plotted for 6 cellular variables over 13.89
626 hours (the default endtime of the simulations). Top row is unmodified genome, showing the expected
627 cellular behaviour (previously show by Karr et al ¹⁷) and is used for comparison. Minesweeper_256
628 and GAMA_236 show deviations in phenotype caused by gene deletions. Non aggregated data for
629 each *in-silico* simulation is available (see Methods).



630 **Figure 4. Comparing the genomes of the *Mycoplasma genitalium* whole-cell model,**
631 **Minesweeper_256, and GAMA_236**

632 The outer ring displays the *M.genitalium* genome (525 genes in total), with modelled genes (401) in
633 navy and unmodelled genes (124, with unknown function) in grey. The middle ring displays the
634 reduced Minesweeper_256 (256 genes) genome in light blue, with genes present in
635 Minesweeper_265 but not in GAMA_236 in dark blue. The inner ring displays the reduced GAMA_236
636 (236 genes) genome in light yellow, with genes present in GAMA_236 but not in Minesweeper_265 in
637 dark yellow. Figure produced from published *M.genitalium* genetic data ^{17,18}, with genetic data for
638 Minesweeper_256 and GAMA_236 available in the Supplementary Information.



639 **Figure 5. Comparing the genomes of Minesweeper_256 and 2954 GAMA genomes**

640 The genome of Minesweeper_256 and all the genomes found by GAMA (that were the same size or
641 smaller) were collated. Each point represents a single genome and is plotted based on a similarity
642 metric (see Methods). The circled genome in the top right is Minesweeper_256 and the circled
643 genome in the bottom left is GAMA_236. The key difference between the genomes is phosphate
644 sources, with Minesweeper_256 using phosphonate and the GAMA genomes using inorganic
645 phosphate.

Gene	Annotation	GO Term (Biological Processes)	Non Essential In	Essential In
MG_039	N/A	N/A	GAMA_236	Minesweeper_256
MG_289	p37	transport	GAMA_236	Minesweeper_256
MG_290	p29	N/A	GAMA_236	Minesweeper_256
MG_291	p69	transport	GAMA_236	Minesweeper_256
MG_427	N/A	OsmC-like protein	GAMA_236	Minesweeper_256
MG_033	glpF	glycerol metabolic process	Minesweeper_256	GAMA_236
MG_410	pstB	N/A	Minesweeper_256	GAMA_236
MG_411	pstA	phosphate ion transmembrane transport process	Minesweeper_256	GAMA_236
MG_412	N/A	N/A	Minesweeper_256	GAMA_236
MG_305	dnaK	protein folding	<i>M.g</i> * whole-cell model	GAMA_236 and Minesweeper_256

646 **Table 1. Low Essentiality Genes from Minesweeper_256 and GAMA_236 genomic contexts**

647 Protein annotation and GO term obtained from KEGG ²⁶ and UniProt ²⁵, based on Fraser et al's

648 *Mycoplasma genitalium** G37 genome ¹⁸.