

1 SuperDCA for genome-wide epistasis analysis

2 Santeri Puranen^{1,2}, Maiju Pesonen^{1,2}, Johan Pensar², Ying Ying Xu^{1,2}, John A. Lees³, Stephen D. Bentley³,
3 Nicholas J Croucher⁴, Jukka Corander^{*,†,2,3,5}, Erik Aurell^{*,†,1,6,7,8}

4 ¹Department of Computer Science, Aalto University, 00076 Espoo, Finland

5 ²Helsinki Institute of Information Technology (HIIT), Department of Mathematics and Statistics,
6 University of Helsinki, 00014 Helsinki, Finland

7 ³Pathogen Genomics, Wellcome Trust Sanger Institute, Cambridge CB10 1SA, UK

8 ⁴Department of Infectious Disease Epidemiology, Imperial College London, London W2 1PG, UK

9 ⁵Department of Biostatistics, University of Oslo, 0317 Oslo, Norway

10 ⁶Department of Computational Biology, KTH –Royal Institute of Technology, 100 44 Stockholm, Sweden

11 ⁷Department of Applied Physics, Aalto University, 00076 Espoo, Finland

12 ⁸Institute of Theoretical Physics, Chinese Academy of Sciences, Beijing, China

13 [†]These authors contributed equally: Erik Aurell, Jukka Corander

14 ^{*}Corresponding authors: Erik Aurell (eaurell@kth.se), Jukka Corander (jukka.corander@medisin.uio.no);

15 Abstract

16 The potential for genome-wide modeling of epistasis has recently surfaced given the possibility of
17 sequencing densely sampled populations and the emerging families of statistical interaction models.
18 Direct coupling analysis (DCA) has earlier been shown to yield valuable predictions for single protein
19 structures, and has recently been extended to genome-wide analysis of bacteria, identifying novel
20 interactions in the co-evolution between resistance, virulence and core genome elements. However,
21 earlier computational DCA methods have not been scalable to enable model fitting simultaneously to
22 10^4 - 10^5 polymorphisms, representing the amount of core genomic variation observed in analyses of
23 many bacterial species. Here we introduce a novel inference method (SuperDCA) which employs a new
24 scoring principle, efficient parallelization, optimization and filtering on phylogenetic information to
25 achieve scalability for up to 10^5 polymorphisms. Using two large population samples of *Streptococcus*
26 *pneumoniae*, we demonstrate the ability of SuperDCA to make additional significant biological findings
27 about this major human pathogen. We also show that our method can uncover signals of selection that
28 are not detectable by genome-wide association analysis, even though our analysis does not require
29 phenotypic measurements. SuperDCA thus holds considerable potential in building understanding about
30 numerous organisms at a systems biological level.

31

32 **Author Summary**

33 Recent work has demonstrated the emerging potential in statistical genome-wide modeling to uncover
34 co-selection and epistatic interactions between polymorphisms in bacterial chromosomes from densely
35 sampled population data. Here we develop the Potts model based approach further into a fully mature
36 computational method which can be applied to most existing bacterial population genomic data sets in
37 a straightforward manner. Our advances are relying on more efficient parameter scoring, highly
38 optimized and parallelized open source C++ code, which does not rely on the computation-intensive
39 polymorphism subsampling approximations used earlier. By analyzing the two largest available
40 population samples of *Streptococcus pneumoniae* (the pneumococcus), we highlight several biological
41 discoveries related to the survival of the pneumococcus and co-evolution of penicillin-binding loci, which
42 were not uncovered by the earlier analyses. Our method holds considerable potential for building
43 understanding about numerous organisms at a systems biological level.

44 **Introduction**

45 Direct Coupling Analysis (DCA) emerged less than a decade ago and has opened up a new direction of
46 biological research by demonstrating that large population based protein sequence analysis can be
47 leveraged to make accurate predictions about protein structure[1-7]. DCA has been successfully
48 extended to predict secondary and tertiary RNA structure[8], synergistic effects on fitness of mutations
49 in the *E. coli* lactamase TEM-1[9], the fitness landscapes of HIV proteins[10], prediction of mutation
50 effects from sequence co-variation[11], and to genome-wide epistasis analysis for bacterial population
51 genomics[12]. Our focus here is to significantly extend the applicability of DCA methodology by enabling
52 scalable inference for two orders of magnitude larger than previously modeled dimensionality of
53 sequence positions.

54 Maximum likelihood inference for the Potts models employed in DCA is intractable due to the form of
55 the normalizing constant of the model distribution, hence various weaker criteria or approximations
56 have been used to derive estimators of the model parameters. Notably, maximum pseudolikelihood is a
57 statistically consistent inference method which has typically outperformed variational methods[13],
58 such as the mean-field estimator[14]. The different software implementations based on regularized
59 maximum pseudolikelihood for DCA applications (plmDCA)[3,14-17] have been designed for at most
60 1000-2000 sequence positions, after which the computation times tend to become prohibitive for
61 practical purposes.

62 To enable use of plmDCA in a much higher dimensional setting, with the order of 10^5 polymorphisms in a
63 bacterial genome, Skwark et al.[12] stratified a genome into non-overlapping windows and sampled
64 randomly a single SNP from each window to form haplotypes of approximately 1,500 sequence
65 positions, on which the plmDCA implementation by Ekeberg et al.[15] could be directly applied. They
66 then used a large number of repeated random sampling of positions from the stratified genome to
67 aggregate information about interactions between polymorphisms across the genome. While this
68 approach was demonstrated to successfully capture both known and novel interactions, it remains very
69 computationally intensive and may still leave important interactions undiscovered as only a fraction of

70 all possible combinations of interactions will be covered even when using large numbers of repeated
71 samples. It is also a hybrid method which does not fully implement global model learning which is a
72 conceptually central point of DCA. To avoid these problems, here we introduce a method termed
73 SuperDCA, which can perform inference simultaneously for all SNP positions in a much higher
74 dimension. These advances are based on a new computational architecture exploiting efficient
75 parallelization and optimization to achieve scalability for up to 10^5 polymorphisms. In addition to being
76 significantly faster with more modest computational resources, we also show that the global inference
77 with SuperDCA allows the discovery of previously undetected epistatic interactions that inform our
78 understanding of bacterial biology related to survival of the pneumococcus at lower temperatures.
79 SuperDCA is freely available from <https://github.com/santeripuranen/SuperDCA>

80 **Results**

81 **Results of SuperDCA and comparison with genomeDCA**

82 The Potts model for genome-wide epistasis analysis was fitted to two largest existing pneumococcal
83 population data sets using the SuperDCA method; the Maela[12,18] and Massachusetts populations[19].
84 Two variants of the Maela population data were considered: one with only bi-allelic SNPs (81045 loci),
85 filtered as in Skwark et al.[12] in order to maintain compatibility for comparison of the results, and the
86 second with no restriction to bi-allelic SNP sites (94028 loci, **Methods**). For Massachusetts 78731 SNP
87 loci were analyzed (**Methods**). Figure 1 shows the cumulative distributions of the estimated coupling
88 strengths between SNP sites for the Maela and Massachusetts populations. In both cases a vast majority
89 of the couplings were of negligible magnitude and could be discarded from further detailed investigation
90 using the thresholds shown in Figure 1 (**Methods**).

91 Figure 1 Log histograms of the cumulative distributions of estimated between-site couplings for Maela
92 (left) and Massachusetts (right). The thresholds indicate the learned boundary between negligible and
93 moderate to strong couplings.

94 Supplementary Figure 1 shows the overlap between the predicted genomeDCA and SuperDCA links on a
95 gene level for Maela population. SuperDCA replicated the previously identified links between PBP gene
96 pairs, as well as the network containing the *smc* gene. In contrast, SuperDCA did not identify significant
97 links between *pspA*, *divIVA*, and the triplet upstream of *ply*, SPN23F19480-19500. In the simultaneous
98 analysis which is not affected by chromosome stratification and random sampling of positions, the
99 respective couplings no longer clearly deviated from the background dependence distribution, which is
100 considerably wider for SuperDCA than for genomeDCA. This illustrated by a closer examination of the
101 pairwise mutual information (MI) values (for further details see **Methods**) between the SNP loci in *pspA*,
102 *divIVA*, and SPN23F19480-19500. The few stronger pairwise dependencies between the three genes
103 disappear when all SNP loci are considered simultaneously. As a consequence of performing a full DCA
104 analysis, in contrast to only partial DCA, the SuperDCA approach is less susceptible to highlighting
105 weaker dependencies than genomeDCA.

106 **Epistasis in the penicillin-binding proteins**

107 Since the bulk of the biological signal of between-site variation dependence presented in Figure 1 is due
108 to linkage disequilibrium (LD) between sites in close proximity, we used a refined version of the
109 phylogenetic ranking of the couplings (Supplementary Tables 1-3, **Methods**), to focus on the strongest
110 candidates of co-selected loci. Figure 2 shows two sets of SNP loci which are involved in the top ranking
111 couplings in Maela, alongside with the phylogenetic distribution of the alleles. The very top ranking
112 couplings are between sites in the three penicillin-binding proteins (PBPs), as discovered in the earlier
113 epistasis analysis which stratified the genome into non-overlapping windows and used the Potts model
114 for sampled subsets of loci to reduce dimensionality[12].

115 Figure 2. Maela population distribution of alleles at top ranked coupled SNP sites. The estimated
116 genome-wide maximum likelihood phylogeny is shown on the left. Each column is labeled by the
117 genome position, gene name and a corresponding functional categorization. Columns marked by red
118 rectangles indicate coupled sites in *pbp2x*, *pbp2b* that have a reversed minor/major allele distribution
119 compared with the remaining displayed SNPs in the same genes.

120 Figure 2 reveals a particular pattern of dependence between PBP mutations that adds significant
121 biological information to the earlier findings[12]. The SNP positions marked by red rectangles in Figure 2
122 have an approximately reversed distribution of minor/major alleles in the population, which may reflect
123 fitness differences regarding co-evolution of emerging mutations. In *pbp2x* the first marked position
124 (codon position 359) corresponds to a synonymous mutation coding for amino acid phenylalanine, part
125 of a conserved cluster of hydrophobic residues (Figures 3A and 3C) consisting of F353, P354, F393, L402,
126 L403, and the E357 to K406 charge interaction located at the upper part of the transpeptidase domain
127 near the active site. This cluster of residues likely has a role in maintaining structural integrity in this
128 region (marked with cyan), as it is positioned next to the more mobile loop (marked with red) at residue
129 positions 362-383 that partially covers the active site. Selection pressure seems to act in favor of the
130 phenylalanine phenotype, since the genotype space clearly is explored here and switching the
131 phenotype to the similarly sized and hydrophobic (but in contrast to phenylalanine non-aromatic)
132 residues leucine or isoleucine, would only require a single non-synonymous mutation.

133 The second and third mutations (codon position 576, N/S/H amino acid changes; codon position 598, I/V
134 amino acid changes) are conservative changes (Figure 3D) that may remotely affect the active site
135 geometry or substrate association/dissociation kinetics, possibly as a compensatory mechanism for
136 changes elsewhere. Active-site reshaping is an established cause of beta-lactam resistance in *S.*
137 *pneumoniae*, where the involved polymorphisms can appear quite subtle at first sight. Our LD adjusted
138 coupling scores indicate a very strong coupling between genome positions 294028/293661 in *pbp2x* and
139 1613045/1613098 in *pbp2b*. The fourth and fifth mutations (codon position 714, conserved L amino
140 acid; codon position 721, E/Q amino acid change) are located in the PASTA-2 domain (Figure 3B; marked
141 with green). The Q721 variant is prevalent in beta-lactam susceptible- and E721 in non-susceptible
142 isolates. PASTA (PBP and Serine/Threonine kinase Associated) domains typically bind beta-lactams,
143 however, a direct mechanistic role for 721 in beta-lactam resistance seems unlikely due to the structural
144 position facing away from the protein core region. Rather, 721 is more likely to be involved in divisome
145 complex formation and functions in a way that supports bacterial resilience in the presence of
146 antibiotics; *pbp2x* and the PASTA domains therein are essential for bacterial division[20,21]. The

147 characteristics and placement of L714 and the fact that all polymorphisms at this site are synonymous,
148 point to a role in assuring structural integrity rather than in direct beta-lactam interaction.

149 In *pbp2b* the second marked position (codon position 458, D/N amino acid change) is located such that
150 it may affect the active site in a mechanistic way via two distinct routes, either by indirectly modifying
151 stability of the loop region (marked with red) proximal to the active site, or by slightly affecting the
152 geometry of active site residues through the helix from 445 to 456 (marked with orange) directly
153 connected to active site residues N445 and S443. The first marked position in *pbp2b* (codon position
154 476, G/E amino acid change) is spatially separated from 458. Although glycine at this site is more
155 prevalent in beta-lactam non-susceptible- and glutamic acid in susceptible isolates, the potential role of
156 the residue at this position in resistance remains unclear and would be a target for further experimental
157 work.

158

159 Figure 3. Structural mapping of the *pbp2x* (panels A-C) and *pbp2b* (panel D) positions marked in Figure 2.
160 The panels show the transpeptidase domains of each PBP with active site residues shown in cyan and
161 positions marked in Figure 2 as sticks in orange or green. Panel A depicts a structure-stabilizing cluster of
162 conserved hydrophobic residues (light gray sticks) and charge interaction (dark gray) in a region
163 proximal to (cyan cartoon) the *pbp2x* active site (with bound inhibitory antibiotic as pink space-filling
164 volume) and a mobile loop (red cartoon) covering the active site. Panel B depicts the PASTA-2 domain
165 essential for divisome complex function (green cartoon) with the bulk of the protein to the right (gray
166 cartoon). Panel C shows an overview of the *pbp2x* transpeptidase domain colored as in the detail views
167 in panels A and B. Panel D depicts the *pbp2b* transpeptidase domain region proximal to the active site
168 with a helix (orange cartoon) mechanically connecting the active site to the 'top' of the protein. An
169 adjacent mobile loop covering the active site is shown in red.

170 Figure 4 shows a clear overlap between the Maela and Massachusetts populations in terms of identified
171 links between genes involved in antibiotic resistance. For the two PBP gene pairs *pbp2x-pbp2b* and
172 *pbp2x-pbp1a* the numbers of strong links between SNPs are large in both populations. For the pair
173 *pbp1a-pbp2b* there is a pronounced asymmetry in this respect, such that the Massachusetts population
174 harbors a large number of links whereas there are only very few in Maela. The latter observation is in
175 line with the findings by Skwark et al.[12] which indicated that most interactions found between the PBP
176 genes were between *pbp2x-pbp2b* and *pbp2x-pbp1a*. The fact that the Massachusetts population clearly
177 deviates from this suggests that the co-evolution of PBPs may follow a non-congruent route in different
178 populations. In the case of Massachusetts versus Maela, this may be a consequence of markedly
179 different serotype distribution in the two populations, or other ecological constraints such as the varying
180 selection pressure from different beta-lactam antibiotic usage. In Maela, beta-lactam prescriptions were
181 almost exclusively amoxicillin, whereas in Massachusetts the pediatric prescription practice is likely to
182 have been considerably more varied. Similar to the asymmetry of the extent of *pbp1a-pbp2b* couplings,
183 the reverse allele distribution pattern discussed previously for Maela was not observed in the
184 Massachusetts population. Given these differences our results suggest that the co-selective pressure on
185 PBP gene polymorphisms acts differently depending on the type of the beta-lactams used in the

186 population, warranting further experimental work to elucidate the mechanistic role of the coupled
187 variations.

188 Figure 4. Overlap of estimated SNP interactions between the Maela and Massachusetts populations.
189 Each dot represents an estimated link (interaction) between two coding sequences (CDSs), the blue
190 CDSs are involved in antibiotic resistance, and the red CDSs are in close proximity to antibiotic resistance
191 loci. Grey dots represent other functional categories not displayed here explicitly for visual clarity. Both
192 axes are on log-scale and the values represent numbers of links in each CDS pair.

193 **Epistasis in cold tolerance and transmission potential**

194 The current analysis additionally highlights several important between-site dependencies not identified
195 by genomeDCA, showing greater sensitivity for identifying putative epistatic interactions. Firstly, the
196 highest ranked SuperDCA couplings included twenty links between cold resistance-related genes
197 exoribonuclease R (*rnr*), glyceroporin (*glpF1*), and lytic amidase C (*lytC*) (Figure 2), the strongest of which
198 was ranked 668. In total, among the 5000 highest ranked couplings, there were two links between *glpF1*
199 and *rnr*, and 18 links between *glpF1* and *lytC*. *GlpF1* is a transporter that imports glycerol, is involved in
200 maintaining membrane fluidity with temperature changes[22]. The *glpF1* gene is at the 3' end of its
201 operon, with a tightly-folding BOX repeat at its distal end[23]. This would make the corresponding mRNA
202 a potential target for *rnr*, a cold shock response 3'->5' exonuclease that degrades tightly-folded RNAs
203 that might be misfolded at lowered temperatures. Hence these interactions may be involved in tuning
204 the expression of *glpF1* at lowered temperatures. Like *glpF1*, *lytC* is involved in maintaining the cell
205 surface at lower temperatures, as it is the cellular amidase specialized at degrading peptidoglycan at
206 lower temperatures (30 degrees Celsius, rather than 35-37 degrees Celsius)[24].

207 Previous work has demonstrated a significant seasonality in the transmission dynamics for the Maela
208 population while carefully controlling for viral epidemics; the probability of the transmission being
209 higher during the cold and dry winter months in comparison to warmer and more humid spring and
210 summer months[25]. To examine whether the observed epistatic links related to survival at lower
211 temperatures are connected with the seasonal transmission phenomenon, we examined the major and
212 minor allele frequencies at the strongly linked cold resistance loci according to months, averaged over
213 the three years 2007-2010 during which the data were sampled. Figure 5 shows clear temporal signals in
214 terms of when the isolates carrying the linked minor/major alleles were sampled. The temporal changes
215 in allele frequencies for the strongest cold resistance related link between *glpF1* (position 2162687) and
216 *rnr* (position 871912), and also for the most strongly coupled sites between *lytC* (position 1533938) and
217 *glpF1* (position 2162676), display a repetitive pattern of synchrony across years. In the first case, the
218 proportion of major alleles in *glpF1* increases towards the end of the year, while in *rnr* the proportion of
219 the minor alleles varies, being the dominant allele in January, April, and December. In the second case,
220 the pattern in *glpF1* remains the same, but the proportion of minor alleles in *lytC* increases towards the
221 winter months.

222 Figure 5. Seasonal variation of the allele frequencies for the two top cold resistance couplings between
223 *glpF1-rnr* and *glpF1-lytC* averaged over three years, 2007-2010. The shaded areas indicate 95%
224 confidence intervals.

225 These findings combined with the earlier results on Maela hosts being more susceptible for transmission
226 during the cold and dry winter months[25], suggest that the recurrent selective advantage related to
227 increased cold tolerance to facilitate survival outside hosts has been sufficient to shape the variation in
228 population allele frequencies. To investigate whether the selection pressure on cold resistance genes
229 could be discovered using a genome-wide association study (GWAS) approach, we coded the phenotype
230 of each sample as winter or summer depending on the sampling date (**Methods**). We then applied the
231 SEER GWAS method to identify polymorphisms that explain the variation in the phenotype[26].
232 Supplementary Figure 2 shows the Manhattan plot of the SEER analysis based on the annotated
233 reference genome. No clear association signal can be seen and the SNP loci within the cold resistance
234 genes are not associated with any markedly smaller p-values than the level of background variation of
235 the association signal.

236 No cold resistance related couplings were found among the top 5000 couplings in the Massachusetts
237 population, which may represent the less variable environmental conditions to which children are
238 exposed, and the sampling of isolates only during winter, rather than year-round. In contrast, the Maela
239 refugee camp conditions are such that the changes in selection exposure are more directly influential.

240

241 **Filtering on phylogenetic information**

242 Inferred couplings from DCA typically have to be filtered to remove those that refer to trivial or non-
243 informative dependencies. In the protein-structure applications very strong couplings are inferred
244 among close neighbors along the peptide backbone, and are usually removed after model fitting by a
245 simple distance based cut-off. A related issue is sampling bias, which for protein-structure applications
246 has been handled by a reweighting applied to each sequence[1]. In bacterial sequence data produced
247 from a sample taken from a small area over a limited period of time, a further issue is clonal inheritance;
248 the meta-population is in a state of flux, and for a short window of time may not fully relax to the
249 postulated Potts model of DCA. To compensate for this problem we used a refined version (**Methods**) of
250 the phylogenetic re-ranking of the coupling estimates introduced in Skwark et al.[12] To visualize its
251 effect, we consider mutual information to characterize the strength of pairwise dependencies between
252 SNP loci. MI is a widely used information theoretical measure of dependence between discrete-valued
253 variables, and it has been a popular tool as part of bioinformatics methods for DNA sequence
254 analysis[27-29]. Here we use MI to characterize the strength of pairwise dependence between SNP loci
255 as a function of their ranked estimated couplings alone and a ranking based jointly on couplings and
256 phylogenetic criteria. Figure 6 shows the distribution of inferred MI values (**Methods**) for the two
257 rankings in both Maela and Massachusetts population. The PBP-related couplings are nearly universally
258 associated with higher MI values, indicating their tighter co-evolution despite of the negligible level of
259 background LD between the three PBP segments. The distributions of large MI values have a clear shift

260 towards a higher rank for both Maela and Massachusetts populations, which succinctly demonstrates
261 the usefulness of using a phylogenetic ranking of coupling estimates to highlight co-selected sites above
262 the background LD. A comparison of MI distributions for PBP-related SNPs for the two populations
263 revealed that Maela displays stronger dependencies between the PBP mutations than Massachusetts
264 (Supplementary Figure 3).

265 Figure 6. Estimated mutual information for 60749 pairs of SNPs (Maela) and 125469 pairs of SNPs
266 (Massachusetts).

267 Scalability improvements in SuperDCA

268 Overall, SuperDCA achieved an 18-fold effective performance increase over the earlier reference
269 plmDCA implementation[15] on a single 20-core dual-socket compute node, enabling inference of
270 $1.4 \cdot 10^{11}$ parameters for a 94028 SNP genome dataset in less than 8 days, instead of an estimated 170
271 days. This was achieved through multiple alterations to the central algorithm explained below. Let
272 (s_1, s_2, \dots, s_N) be a haplotype over N SNP loci, where each s_i can take values from an alphabet with
273 cardinality q . Typically this cardinality varies between three (allelic states: minor/major/gap) and five
274 (allelic states: A,C,G,T, gap). A Potts model assigns a probability distribution on such haplotypes defined
275 by the following formula

$$P(s_1, s_2, \dots, s_N) = \frac{1}{Z} e^{E(s_1, s_2, \dots, s_N)}$$

276 where the normalizing constant Z is known as the partition function and the expression in the exponent
277 is

$$E(s_1, s_2, \dots, s_N) = \sum_{i=1}^N \sum_{a=1}^q h_i(a) \delta_{s_i, a} + \sum_{i,j=1}^N \sum_{a,b=1}^q J_{ij}(a, b) \delta_{s_i, a} \delta_{s_j, b}$$

278 In above $\delta_{x,y}$ represents the Kronecker delta function which takes the value one if the arguments x and
279 y are equal, and is otherwise zero. The linear terms are $h_i(a) \delta_{s_i, a}$ for different SNP loci and their alleles.
280 The coefficients $h_i(a)$ parametrize a deviation from the uniform allele distribution for each SNP,
281 independently of the values of all the other variables. The quadratic terms are the matrix elements
282 $J_{ij}(a, b) \delta_{s_i, a} \delta_{s_j, b}$ for different combinations of values of i and j , and a and b . The coefficients $J_{ij}(a, b)$,
283 which are the *couplings* or *interactions* of pairs of SNPs, are defined as zero when the two indices i and j
284 are equal. A coupling matrix with all elements equal to zero for non-identical locus index pairs implies
285 that the alleles at these two loci are distributed independently in the population. Small positive values of
286 the coupling matrix elements correspond to weak dependence between the SNP loci. In this paper we
287 have addressed the issues of gauge invariance and gauge fixing in the Potts model[1] as described
288 previously[12,15].

289 One of the major obstacles for using earlier plmDCA algorithms simultaneously on large numbers of
290 SNPs without locus subset sampling is their large runtime memory requirements. plmDCA memory use
291 is dominated by the storage of q^2 -dimensional parameter matrices J_{ij} , where q is the cardinality of the

292 SNP state space (the maximum value being $q = 5$ when a gap/indel is included). J_{ij} and J_{ji} are needed
293 simultaneously for calculating the pairwise coupling value, and since the elements are inferred row- or
294 column-wise for all i (or j) at a time, a straightforward implementation of the algorithm necessitates
295 simultaneous storage of all couplings in an N -by- N matrix J , therefore storage is needed for $q^2(N^2-N)$
296 scalar elements. The scoring of the estimated coupling matrices would then be calculated according to

$$297 \quad J_{ij} = J_{ji} = \|(J_{ij} + J_{ji})/2\|_F$$

298 where F indicates the Frobenius norm. As an example, if a 10^5 SNP genome alignment was characterized
299 by 5-state alphabet and parameters stored in 64-bit floating point format, then the full interaction
300 matrix J would require approximately 1.8 TB of memory, which is typically beyond the RAM available in
301 state-of-the-art HPC cluster nodes. However, if the scoring of coupling values is instead calculated as

$$J_{ij} = J_{ji} = (\|J_{ij}\|_F + \|J_{ji}\|_F)/2$$

302 then runtime storage requirements are reduced by a substantial factor and the intermediate storage
303 requirements for our example would shrink to 74 GB, which is well in the feasible range for current HPC
304 nodes. Supplementary Figure 4 illustrates numerically that the above two scoring approaches lead to
305 insignificant numerical differences in practice. SuperDCA uses this finding as one of its key
306 improvements of plmDCA.

307 Performance profiling analysis identified high memory requirements and poor cache utilization as a
308 major bottleneck for the performance in earlier plmDCA implementations when applied to higher-
309 dimensional data. Parallel execution scaling also suffered due to memory bandwidth starvation. The
310 maximization step was performed by Ekeberg et al.[15] using L-BFGS gradient-based optimization.
311 However, the objective function required repeated traversal through all input data and the full
312 parameter vector, emphasizing the need for an efficient data structure. To remedy these issues, a space-
313 efficient, block-wise ordered data structure with simple state-pattern dictionary and run-length encoded
314 indexing strategy for the genome data and a cache-friendly blocked memory layout for parameters were
315 developed for SuperDCA and implemented in C++ (Supplementary Figure 5). A particular design choice
316 was made to restrict the maximum value of q to 4. The resulting data structure reduced runtime
317 memory use for nucleotide alignments by more than 4-fold compared with a typical dense data matrix
318 representation. It also helped to reduce computing effort, improved processor cache utilization and
319 enabled efficient utilization of SIMD vector instructions. The aggregate effect of these changes was an 8-
320 fold improvement in single-threaded performance. The reduced main memory bandwidth use also
321 helped improve node-level scaling as we measure a strong scaling factor of >0.7 up to 20 cores.
322 Supplementary Figures 6-8 illustrate the computational scalability aspects for SuperDCA compared with
323 genomeDCA.

324

325 **Discussion**

326 Production of natural population genomic sequence data is currently still exponentially accelerating,
327 highlighting the need for statistical methods that can generate detailed hypotheses for further
328 experimental work regarding loci likely to be important in shaping bacterial evolution. Genome-wide
329 association analysis has for a decade been the major general tool for such purposes, and more recently,
330 its applicability to bacteria has been also demonstrated[26,30-32]. Skwark et al.[12] showed for the first
331 time that statistical genome-wide modeling of joint SNP variation using DCA can uncover valuable
332 information about co-evolutionary pressures on a large scale. This was done without relying on any
333 phenotypic measurements, and using a hybrid scheme that does not fully employ the global model
334 learning aspect of DCA. Here we built upon this initial observation to develop DCA into a powerful tool
335 that is applicable to a majority of the existing bacterial population genome data sets in a
336 computationally scalable manner. The biological insights on the differential evolution of PBPs, and the
337 cold tolerance mechanisms, derived from the results of applying SuperDCA to two of the largest
338 available pneumococcal genome data sets illustrate succinctly how such an approach could provide vital
339 clues to the evolutionary processes under different ecological conditions in natural populations.

340 As the size of genome sequence data sets keeps growing, even our optimized parallel inference
341 algorithm will eventually become too inefficient for practical purposes. Currently the chosen data and
342 algorithmic architecture work extremely effectively for up to around 10^5 polymorphisms. As bacterial
343 whole genome alignments are typically of the order of 10^6 sites, this should be sufficient for most
344 population genomic studies. After this, the runtime will start to increase so rapidly that different
345 computational strategies will be required for data sets including significantly more SNPs. Thus, an
346 important topic for future research is to investigate how the Potts model inference can be performed in
347 a reliable manner without resorting to a quadratic increase in the computational complexity as a
348 function of the number of polymorphisms.

349 **Materials and Methods**

350 **Data pre-processing**

351 Bi- or tri-allelic loci with a minor allele frequency (MAF) greater than 1% were included in the analysis,
352 provided that gap frequency was less than 15%. Gaps were not counted as alleles in the frequency
353 calculations. To facilitate direct comparison with previous results[12] a separate dataset was prepared
354 from the Maela input alignment using otherwise the same filtering rules, but for bi-allelic loci only.
355 Filtering of 305245 SNPs in total resulted in two Maela input datasets for SuperDCA containing 94028
356 SNPs and 3042 samples using the former rules, and 81045 SNPs and 3145 samples using the latter rules.
357 A subset of 103 samples containing mostly low quality reads were included in the data in the previous
358 study, but were here removed from the source alignment prior to locus pre-selection for our 94028 SNP
359 set. For the Massachusetts population the first set of filtering criteria resulted in 78733 SNPs and 670
360 samples.

361 **Hardware and inference details**

362 Parameter inference was performed using a single 20-core HP SL230s G8 compute node with dual Xeon
363 E5 2680 v2 CPUs and 256GB of DDR3-1667 RAM. Total wall clock run times were 186h (Maela with

364 94028 SNPs), 167h (Maela with 81045 SNPs) and 39h (Massachusetts with 78733 SNPs), including file
365 I/O, pre-filtering and parameter inference. Weights correcting for the population structure,
366 regularization and choice of hyper-parameters were calculated exactly as in the genomeDCA
367 method[12]. Coupling estimates for the three data sets which exceeded the cut-off described below are
368 provided as Supplementary Tables 1-3 at <https://github.com/santeripuranen/SuperDCA>.

369 **Prediction cut-off**

370 The Potts models inferred in DCA are heavily over-parametrized. In protein contact applications the
371 benchmark number of parameters is typically in the millions (number of residue pairs times q^2 , where
372 $q=21$), while the number of samples varies typically from thousands to hundreds of thousands.
373 Therefore only a small fraction of largest predictions are retained, commonly in the order of hundreds.
374 Accordingly DCA manifests a varying degree of success when applied to protein families of the same size
375 which has led to sustained efforts in algorithm optimization[6].

376 For the present and future applications to whole-genome data it is of more relevance to deliver a set of
377 predictions at a pre-determined level of deviance from zero. An earlier approach using deviations from
378 an extreme value theory distribution (Gumbel distributions)[12] was not applicable in the present set-
379 up, since we are not only sampling the tail of the coupling coefficients but estimate couplings for all
380 possible pairs of SNPs. As shown in Figure 1, a semi-logarithmic cumulative distribution plot provides a
381 computationally straightforward way to assess whether a particular coupling represents only random
382 fluctuation near zero. The null distribution theory developed in Xu et al. for DCA inference procedures
383 provides a strong motivation for using the linear part of the distribution near the origin as
384 representation of the noise level signals[33]. To obtain a threshold we first perform a systematic scan
385 over the histogram bins to fit a two-component linear spline function to the cumulative distribution. The
386 standard deviation of the null couplings was then estimated using the part of the distribution between
387 zero and the breakpoint. Similar to the Gumbel fit deviance level used by Skwark et al. [12] we then
388 exclude all couplings that are less than 6 standard deviations away from the linear trend from further
389 analysis. Figure 1 illustrates that this procedure effectively filters out the vast majority of all possible
390 couplings as noise, and allows the downstream analysis to focus on the relevant signals.

391 **Phylogenetic ranking of estimated couplings**

392 By default, SuperDCA includes gaps as a state in the Potts model if they are found in the alignment at
393 sites fulfilling the SNP pre-filtering criteria. Some gaps can be considered informative, representing
394 indels, while some simply relate to sites that are difficult to sequence. Hence some strong gap-induced
395 couplings can represent lower quality sequence data instead of true between-site interactions, and they
396 should be automatically de-emphasized to better enable assessment of the biological meaning of the
397 inferred couplings. The superDCA coupling estimates are by default re-ranked using a combination of
398 the three criteria described below, in addition to the actual value of the coupling.

399 Let C be a set of estimated couplings and $c_i = [c_{i1}, c_{i2}] \in C$ a pair of SNP loci represented by their
400 genome position indices. Let $y_b = [s_1^{(b)}, \dots, s_N^{(b)}]$ be a haplotype over the N SNP loci. Further, $\mathcal{S}_{i,1}$ is a

401 set of haplotypes carrying a minor allele at locus $c_{i,1}$ and $\mathcal{S}_{i,2}$ a set of haplotypes with a minor allele at
 402 locus $c_{i,2}$. The first phylogenetic ranking criterion is the minimum of the average genome-wide Hamming
 403 distances of all pairs of isolates $y_k, y_l \in \mathcal{S}_{i,o}, o = \{1,2\}, k < l$, i.e. $d_i = \min_o(\bar{d}_{\mathcal{S}_{i,o}}(y_k, y_l)), i =$
 404 $1, \dots, |C|$, where $\bar{d}_{\mathcal{S}_{i,o}}(y_k, y_l) = \frac{1}{|\mathcal{S}_{i,o}|} \sum_{n=1}^N (s_n^{(k)} \neq s_n^{(l)})$.

405 Our second criterion is the normalized number of hierBAPS[34] clusters including isolates carrying the
 406 minor alleles at the two coupled loci, i.e. $a_i = |\{\beta_b | b \in (\mathcal{S}_{i,1} \cap \mathcal{S}_{i,2})\}|$, where β_b is the designated
 407 hierBAPS cluster for haplotype b . Finally, the third criterion is the percentage of isolates where both SNP
 408 loci involved in a coupling had the minor allele, i.e. $m_i = \frac{1}{2} \left[\frac{|\mathcal{S}_{i,1} \cap \mathcal{S}_{i,2}|}{|\mathcal{S}_{i,1}|} + \frac{|\mathcal{S}_{i,1} \cap \mathcal{S}_{i,2}|}{|\mathcal{S}_{i,2}|} \right]$.

409 The above three criteria are normalized by: $d_i^{norm} = \frac{d_i}{\max_i(d_i)}, i = 1, \dots, |C|$, $a_i^{norm} = a_i / \max_i(a_i)$,
 410 and $m_i^{norm} = m_i / \max_i(m_i)$, after which they are combined to a single ranking criterion: $r_i = d_i^{norm} +$
 411 $a_i^{norm} + m_i^{norm}$ having a maximum value of three and a minimum equal to zero. Large values
 412 emphasize cases where both minor alleles at coupled loci are simultaneously widely distributed across
 413 the population. In cases where gaps at any two loci are phylogenetically spread in the population and
 414 would have led to a large estimated coupling values, they are still de-emphasized since they are not
 415 counted as minor alleles. The above criteria are derived by normalizing the individual coupling re-
 416 ranking measures developed by Skwark et al. [12]. The hierBAPS clusterings were obtained from the
 417 original publications introducing genome sequences for the Maela and Massachusetts
 418 populations[18,19].

419

420 Mutual information calculations

421 Mutual information is an information theoretic measure of the mutual dependence between two
 422 variables. Let X_1 and X_2 be two discrete variables with outcome spaces indexed by $i = 1, \dots, r_1$ and
 423 $j = 1, \dots, r_2$, respectively (the outcome indexing differs here from the earlier description of Potts model
 424 for notational simplicity). Let $p = (p_{ij})$ represent the joint distribution over the variables such that p_{ij}
 425 corresponds to the probability of $(X_1 = i, X_2 = j)$. The mutual information between X_1 and X_2 is then
 426 calculated by

$$I(p) = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} p_{ij} \log \frac{p_{ij}}{p_i \cdot p_j},$$

427 where $p_i = \sum_{j=1}^{r_2} p_{ij}$ and $p_j = \sum_{i=1}^{r_1} p_{ij}$ are the marginal probabilities for the corresponding variables
 428 (here SNPs). In practical applications the joint distribution is usually not known and must be estimated
 429 from data. The standard approach of estimating the probabilities is to use the maximum likelihood
 430 estimates given by the relative frequencies $\hat{p}_{ij} = n_{ij}/n$, where n_{ij} denotes the count of the
 431 corresponding configuration and n is the sample size. A drawback of the standard frequentist approach
 432 is that it does not account for the uncertainty of the estimates.

433

434 In this work we adopted a Bayesian approach[35], where we put a prior density function $f(p)$ on the
435 $r_1 \cdot r_2$ unknown joint probabilities. Taking the data into account, the posterior density function
436 $f(p|\text{data})$ can be calculated from the prior using Bayes' theorem. The posterior density over the mutual
437 information given some dataset is then obtained through

$$f(I|\text{data}) = \int \delta(I(p) - I) f(p|\text{data}) dp,$$

438 where $\delta(\cdot)$ is the Dirac delta function.

439

440 The above density function can be approximated using a Monte Carlo simulation by sampling from
441 $f(p|\text{data})$. In particular, assuming a Dirichlet prior on p with hyperparameters α_{ij} enables a
442 straightforward sampling scheme since the posterior $p|\text{data}$ then follows a Dirichlet distribution with
443 updated hyperparameters $\alpha'_{ij} = n_{ij} + \alpha_{ij}$. Our main interest is to calculate the Bayesian point estimate
444 given by the posterior mean

$$E_{p|\text{data}}(I) = \int_0^{\infty} I f(I|\text{data}) dI.$$

445 For this particular purpose, there exists an exact closed-form expression[35]:

$$E_{p|\text{data}}(I) = \frac{1}{\alpha'} \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \alpha'_{ij} [\psi(\alpha'_{ij} + 1) - \psi(\alpha'_{i \cdot} + 1) - \psi(\alpha'_{\cdot j} + 1) - \psi(\alpha' + 1)]$$

446 where $\psi(\cdot)$ is the digamma function and $\alpha' = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \alpha'_{ij}$. When using the above estimator we
447 define the hyperparameters using the reference prior $\alpha_{ij} = \frac{1}{r_1 r_2}$. To adjust for the population structure
448 in the sample, we use the same re-weighting scheme as was applied in our SuperDCA inference with a
449 similarity threshold of 0.90. Finally, to remove the influence of gap-gap interactions, we did not include
450 sequences for which either of the two considered loci had a gap value.

451 **GWAS for the seasonality phenotype**

452 We coded season as a binary variable based on whether isolates were acquired during the winter or the
453 summer. We then tested 123791 SNPs passing simple frequency filtering (>1% MAF) for association with
454 this variable using SEER[26], which performs a logistic regression at every SNP. We used the first three
455 multi-dimensional scaling components of the pair-wise distance matrix as fixed effects to control for
456 population structure[26].

457 **Structural analyses**

458 Crystal structures of *S. pneumoniae* PBPs with the following IDs: 2WAF
459 (pbp2b), 1QMF and 1RP5 (pbp2x) were retrieved from the Protein Data Bank[36] (www.rcsb.org;
460 accession date January 8, 2016) and visualized in The PyMOL Molecular Graphics System, Version 1.8.4.0
461 (Schrödinger, LLC). A chimera of 1QMF (chain A residues 257-618) and 1RP5 (chain A residues 64-256
462 and 619-750; missing sidechain atoms of E721 were reconstructed) was used for visualizing pbp2x.

463

464

465

466 **References**

467

- 468 1. Weigt M, White RA, Szurmant H, Hoch JA, Hwa T (2009) Identification of direct residue contacts in
469 protein-protein interaction by message passing. *Proc Natl Acad Sci U S A* 106: 67-72.
- 470 2. Morcos F, Pagnani A, Lunt B, Bertolino A, Marks DS, et al. (2011) Direct-coupling analysis of residue
471 coevolution captures native contacts across many protein families. *Proc Natl Acad Sci U S A* 108:
472 E1293-1301.
- 473 3. Feinauer C, Skwark MJ, Pagnani A, Aurell E (2014) Improving contact prediction along three
474 dimensions. *PLoS Comput Biol* 10: e1003847.
- 475 4. Morcos F, Hwa T, Onuchic JN, Weigt M (2014) Direct coupling analysis for protein contact prediction.
476 *Methods Mol Biol* 1137: 55-70.
- 477 5. Ovchinnikov S, Kamisetty H, Baker D (2014) Robust and accurate prediction of residue-residue
478 interactions across protein interfaces using evolutionary information. *Elife* 3: e02030.
- 479 6. Ovchinnikov S, Park H, Varghese N, Huang PS, Pavlopoulos GA, et al. (2017) Protein structure
480 determination using metagenome sequence data. *Science* 355: 294-298.
- 481 7. Soding J (2017) Big-data approaches to protein structure prediction. *Science* 355: 248-249.
- 482 8. De Leonardis E, Lutz B, Ratz S, Cocco S, Monasson R, et al. (2015) Direct-Coupling Analysis of
483 nucleotide coevolution facilitates RNA secondary and tertiary structure prediction. *Nucleic Acids*
484 *Res* 43: 10444-10455.
- 485 9. Figliuzzi M, Jacquier H, Schug A, Tenaillon O, Weigt M (2016) Coevolutionary Landscape Inference and
486 the Context-Dependence of Mutations in Beta-Lactamase TEM-1. *Mol Biol Evol* 33: 268-280.
- 487 10. Barton JP, Goonetilleke N, Butler TC, Walker BD, McMichael AJ, et al. (2016) Relative rate and
488 location of intra-host HIV evolution to evade cellular immunity are predictable. *Nat Commun* 7:
489 11660.
- 490 11. Hopf TA, Ingraham JB, Poelwijk FJ, Scharfe CP, Springer M, et al. (2017) Mutation effects predicted
491 from sequence co-variation. *Nat Biotechnol* 35: 128-135.
- 492 12. Skwark MJ, Croucher NJ, Puranen S, Chewapreecha C, Pesonen M, et al. (2017) Interacting networks
493 of resistance, virulence and core machinery genes identified by genome-wide epistasis analysis.
494 *PLoS Genet* 13: e1006508.
- 495 13. Wainwright M, Jordan MI (2008) Graphical models, exponential families, and variational inference.
496 Boston: Now Publishers. 310 p. p.
- 497 14. Ekeberg M, Lovkvist C, Lan Y, Weigt M, Aurell E (2013) Improved contact prediction in proteins: using
498 pseudolikelihoods to infer Potts models. *Phys Rev E Stat Nonlin Soft Matter Phys* 87: 012707.
- 499 15. Ekeberg M, Hartonen T, Aurell E (2014) Fast pseudolikelihood maximization for direct-coupling
500 analysis of protein structure from many homologous amino-acid sequences. *Journal of*
501 *Computational Physics* 276: 341-356.
- 502 16. Kamisetty H, Ovchinnikov S, Baker D (2013) Assessing the utility of coevolution-based residue-
503 residue contact predictions in a sequence- and structure-rich era. *Proc Natl Acad Sci U S A* 110:
504 15674-15679.
- 505 17. Seemayer S, Gruber M, Soding J (2014) CCMpred--fast and precise prediction of protein residue-
506 residue contacts from correlated mutations. *Bioinformatics* 30: 3128-3130.
- 507 18. Chewapreecha C, Harris SR, Croucher NJ, Turner C, Marttinen P, et al. (2014) Dense genomic
508 sampling identifies highways of pneumococcal recombination. *Nat Genet* 46: 305-309.

- 509 19. Croucher NJ, Finkelstein JA, Pelton SI, Mitchell PK, Lee GM, et al. (2013) Population genomics of
510 post-vaccine changes in pneumococcal epidemiology. *Nat Genet* 45: 656-663.
- 511 20. Peters K, Schweizer I, Beilharz K, Stahlmann C, Veening JW, et al. (2014) *Streptococcus pneumoniae*
512 PBP2x mid-cell localization requires the C-terminal PASTA domains and is essential for cell shape
513 maintenance. *Mol Microbiol* 92: 733-755.
- 514 21. Tsui HC, Boersma MJ, Vella SA, Kocaoglu O, Kuru E, et al. (2014) Pbp2x localizes separately from
515 Pbp2b and other peptidoglycan synthesis proteins during later stages of cell division of
516 *Streptococcus pneumoniae* D39. *Mol Microbiol* 94: 21-40.
- 517 22. Blaby IK, Lyons BJ, Wroclawska-Hughes E, Phillips GC, Pyle TP, et al. (2012) Experimental evolution of
518 a facultative thermophile from a mesophilic ancestor. *Appl Environ Microbiol* 78: 144-155.
- 519 23. Croucher NJ, Vernikos GS, Parkhill J, Bentley SD (2011) Identification, variation and transcription of
520 pneumococcal repeat sequences. *BMC Genomics* 12: 120.
- 521 24. Moscoso M, Lopez E, Garcia E, Lopez R (2005) Implications of physiological studies based on genomic
522 sequences: *Streptococcus pneumoniae* TIGR4 synthesizes a functional LytC lysozyme. *J Bacteriol*
523 187: 6238-6241.
- 524 25. Numminen E, Chewapreecha C, Turner C, Goldblatt D, Nosten F, et al. (2015) Climate induces
525 seasonality in pneumococcal transmission. *Sci Rep* 5: 11344.
- 526 26. Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, et al. (2016) Sequence element
527 enrichment analysis to determine the genetic basis of bacterial phenotypes. *Nature*
528 *Communications* doi:10.1038/ncomms12797.
- 529 27. Delcher AL, Harmon D, Kasif S, White O, Salzberg SL (1999) Improved microbial gene identification
530 with GLIMMER. *Nucleic Acids Res* 27: 4636-4641.
- 531 28. Li M, Badger JH, Chen X, Kwong S, Kearney P, et al. (2001) An information-based sequence distance
532 and its application to whole mitochondrial genome phylogeny. *Bioinformatics* 17: 149-154.
- 533 29. Mahony S, Auron PE, Benos PV (2007) Inferring protein-DNA dependencies using motif alignments
534 and mutual information. *Bioinformatics* 23: 1297-1304.
- 535 30. Chen PE, Shapiro BJ (2015) The advent of genome-wide association studies for bacteria. *Curr Opin*
536 *Microbiol* 25: 17-24.
- 537 31. Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, et al. (2014) Comprehensive
538 identification of single nucleotide polymorphisms associated with beta-lactam resistance within
539 pneumococcal mosaic genes. *PLoS Genet* 10: e1004547.
- 540 32. Weinert LA, Chaudhuri RR, Wang J, Peters SE, Corander J, et al. (2015) Genomic signatures of human
541 and animal disease in the zoonotic pathogen *Streptococcus suis*. *Nat Commun* 6: 6740.
- 542 33. Xu Y, Aurell E, Corander J, Kabashima Y (2017) Statistical properties of interaction parameter
543 estimates in direct coupling analysis. arXiv: 170401459
- 544 34. Cheng L, Connor TR, Siren J, Aanensen DM, Corander J (2013) Hierarchical and spatially explicit
545 clustering of DNA sequences with BAPS software. *Mol Biol Evol* 30: 1224-1228.
- 546 35. Hutter M (2002) Distribution of mutual information. *Advances in Neural Information Processing*
547 *Systems* MIT Press, Cambridge 14: 399-406.
- 548 36. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, et al. (2000) The Protein Data Bank. *Nucleic*
549 *Acids Res* 28: 235-242.

550 **Supporting information legends**

551 **S1 Fig. Overlap of the predicted genomeDCA and SuperDCA couplings in Maela population.** Lines are
552 plotted between genes having at least three SNPs linked from the same genes, both in genomeDCA and
553 SuperDCA. This results in 274 overlapping interactions. The thickness of lines is proportional to the

554 number of linked positions within the corresponding genes. Gene annotations shown outside the circle
555 are centered at the positions of the corresponding genes. Red labels are given for genes linked tightly
556 both with genomeDCA and SuperDCA, black labels for genes linked only with genomeDCA and blue
557 labels for genes linked only with SuperDCA.

558 **S2 Fig. SEER GWAS Manhattan plot of p-values for the winter/summer phenotype.**

559 **S3 Fig. Boxplots of MI value distributions for pairs of SNPs in two different PBP genes.**

560 **S4 Fig. Comparison of the norm-of-mean (vertical axis) versus the mean-of-norms (horizontal axis)**
561 **summary score strategies.** Differences between the two are negligible for stronger, statistically
562 significant coupling values and show more pronounced deviations only towards the sub-significant
563 domain. The plot was calculated using a 25% uniformly random sample of loci from the 94028 SNP
564 Maela dataset using the full data as background. Coloring marks coupling value count in log-scale.

565 **S5 Fig. Schematic drawings of the central data structures used in SuperDCA for storing input state data**
566 **(nucleotide alignments) and the inferred parameters.** The *input data matrix* is stored such that samples
567 (i.e. isolate genomes in our case) are ordered row-wise, with each sample divided into blocks of size b .
568 Only column-wise unique blocks are stored. Block indices are stored for sample-oriented access to the
569 data and sample indices for column-oriented access. Block index lists can optionally be run-length
570 encoded, which leads to very significant space savings in particular when storing full-genome alignments
571 with large regions of low column-wise variation. Index-lists for column-oriented access can similarly be
572 collapsed for saving storage space when indices form contiguous (ascending) sequences. The *inferred*
573 *parameters* are stored in blocked format such that all parameters relating to a particular column block in
574 the input data are grouped together.

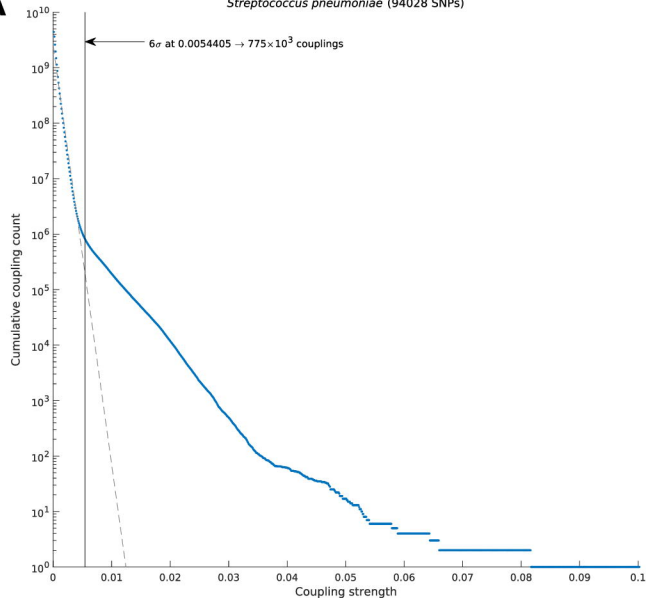
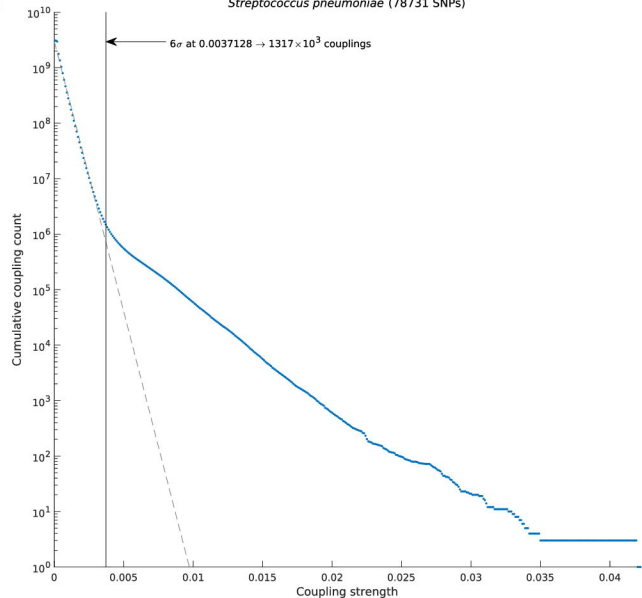
575 **S6 Fig. Comparison of SuperDCA versus plmDCA parallel scaling efficiency.** SuperDCA (blue curve)
576 shows markedly stronger scaling than plmDCA (red curve). The scaling numbers were obtained as a
577 mean of three runs of three 2-permil uniformly random samples of loci (188 loci) from the 94028 SNP
578 Maela dataset and using the full data as background. Inferred parameter storage was disabled in
579 plmDCA for the purpose of benchmarking. All benchmarks were run on a single 20-core HP SL230s G8
580 compute node with dual Xeon E5 2680 v2 CPUs and 256GB of DDR3-1667 RAM.

581 **S7 Fig. Comparison of SuperDCA versus plmDCA sample size scaling.** The sample-compressing
582 datastructure used in SuperDCA enables markedly stronger scaling (blue bars) with increasing sample
583 size than plmDCA (red curve). The scaling numbers were obtained as a mean of 9 runs: three-by-three
584 sets of runs using a uniformly random sample of sequences and run for three 2-permil uniformly
585 random samples of loci (188 loci) from the 94028 SNP Maela dataset and using the full data as
586 background. See caption of Supplementary Figure 6 for details of benchmark hardware.

587

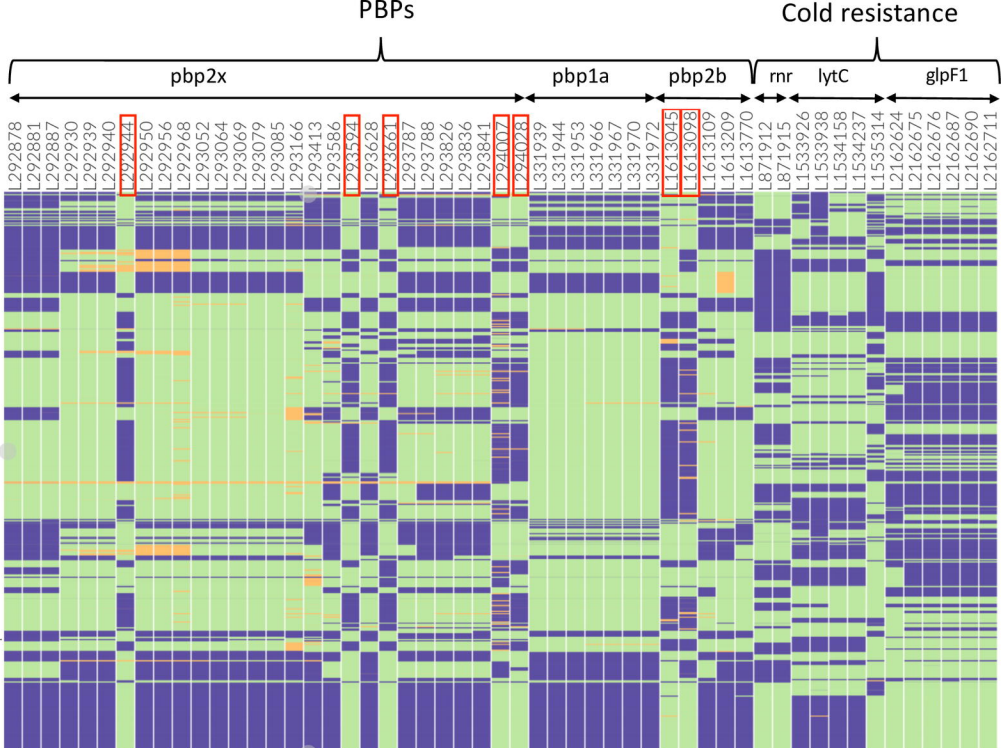
588 **S8 Fig. SuperDCA runtime improvement over plmDCA.** The single-threaded performance of SuperDCA
589 is more than 8-fold that of plmDCA. Due to the greater parallel scalability of SuperDCA the performance

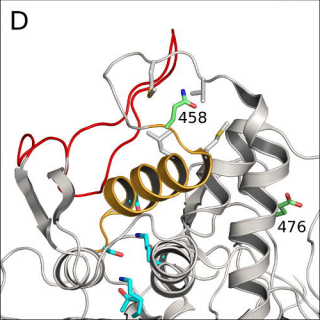
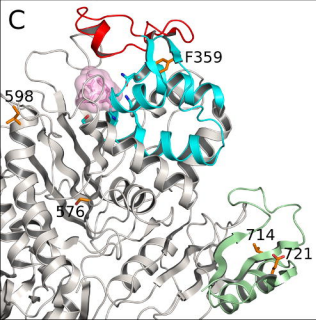
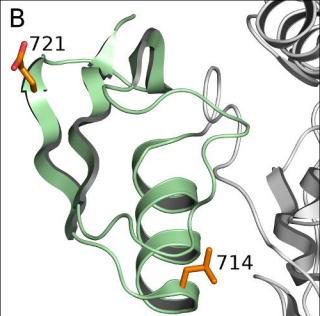
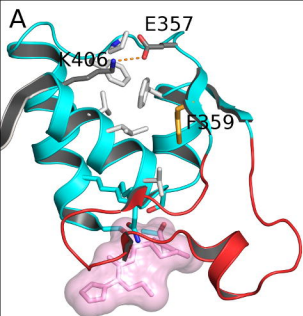
590 delta grows as more compute threads are used, reaching more than 17-fold when run on 20 cores. See
591 caption of Supplementary Figure 6 for details of benchmark settings and hardware.

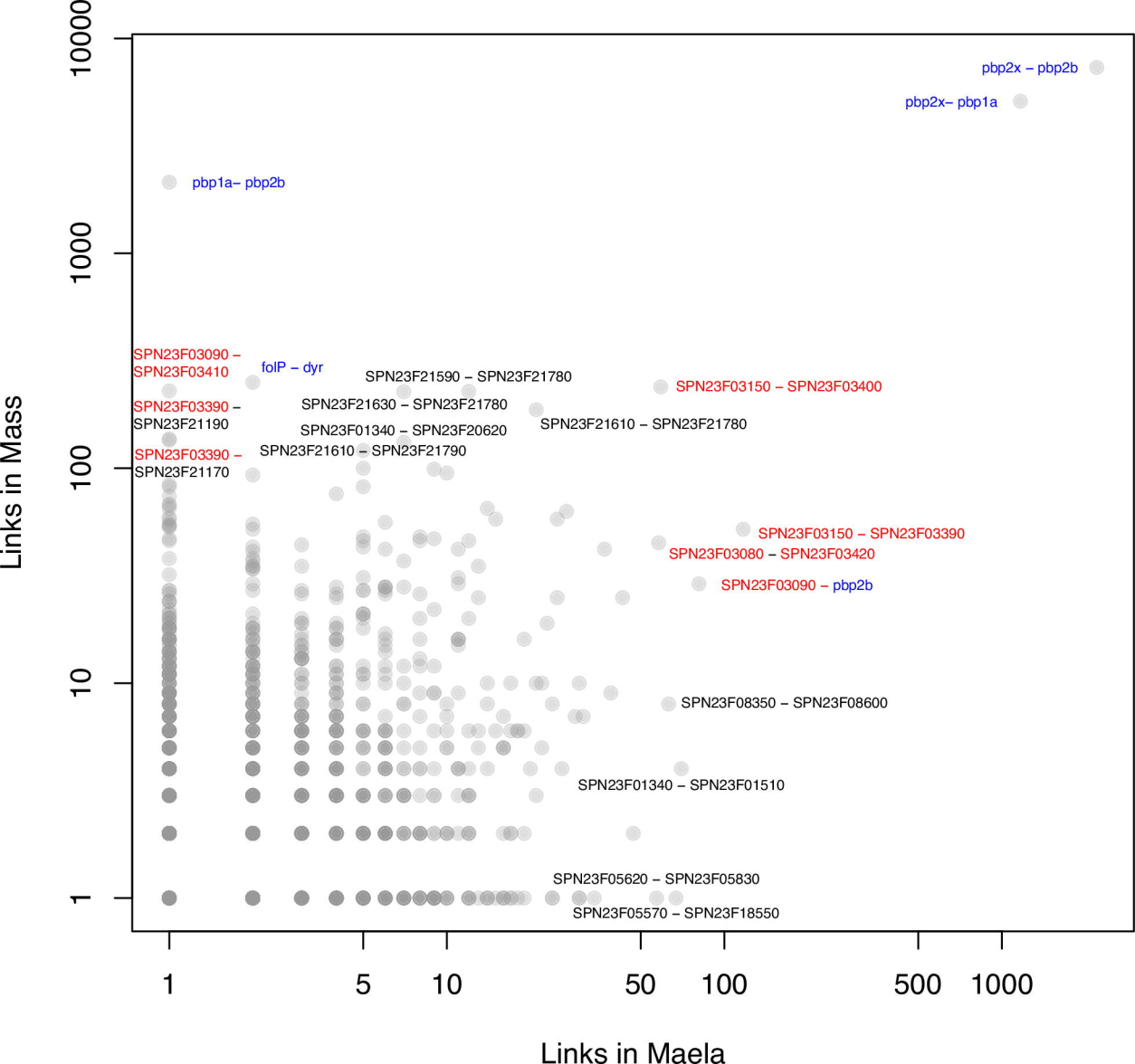
A**Cumulative coupling value distribution**
Streptococcus pneumoniae (94028 SNPs)**B****Cumulative coupling value distribution**
Streptococcus pneumoniae (78731 SNPs)



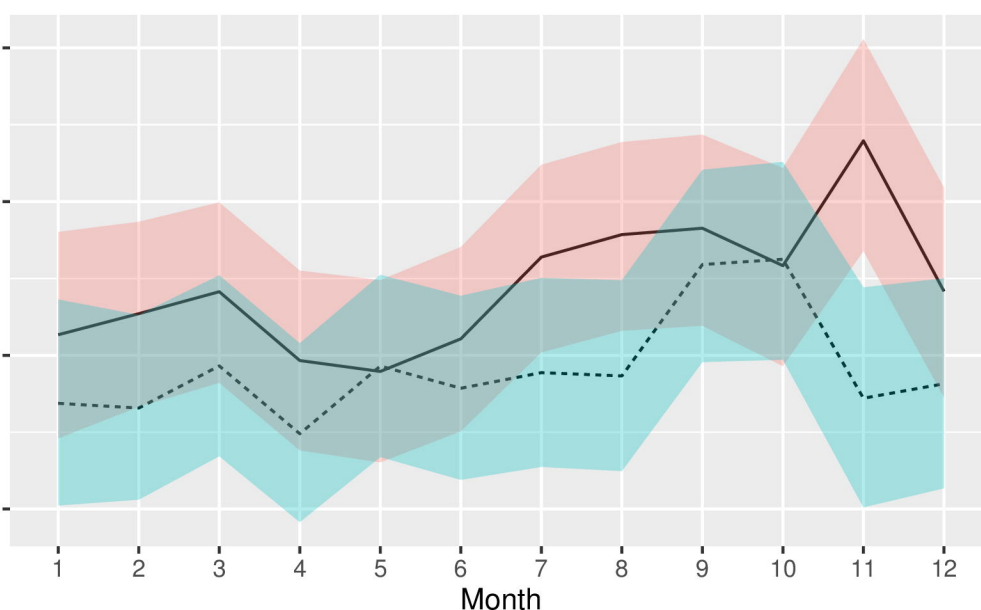
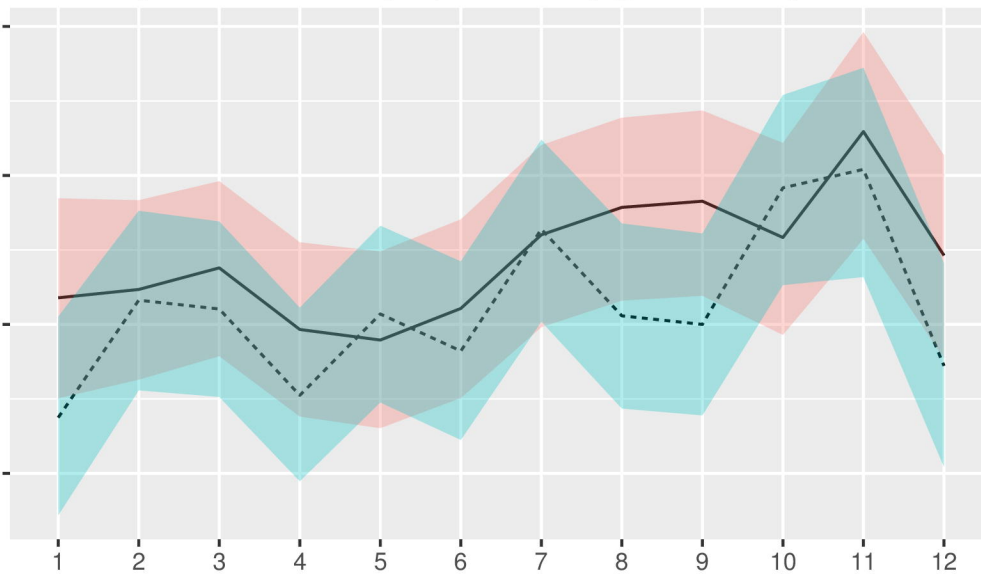
■ = minor
■ = gap
■ = major







Major/ minor allele proportions by gene and by month



Mutual information of significant couplings

