

1                                   **Title: Evolutionary assembly patterns of prokaryotic genomes**

2

3   **Authors:** Maximilian O. Press<sup>1</sup>, Christine Queitsch<sup>1</sup>, Elhanan Borenstein<sup>1,2,3\*</sup>

4

5   **Affiliations**

6   <sup>1</sup>: Department of Genome Sciences, University of Washington, Seattle, WA, USA

7   <sup>2</sup>: Department of Computer Science and Engineering, University of Washington, Seattle, WA,

8   USA

9   <sup>3</sup>: External Faculty, Santa Fe Institute, Santa Fe, NM, USA

10   \*Correspondence to: [elbo@uw.edu](mailto:elbo@uw.edu)

11

12   **Running title:** Evolutionary assembly of prokaryotic genomes

13

14   **Keywords**

15   Genome evolution, prokarya, prokaryote, constraint, epistasis, comparative method, evolutionary  
16   predictability, RuBisCO, parallel evolution.

1 **Abstract:**

2 Evolutionary innovation must occur in the context of some genomic background, which limits  
3 available evolutionary paths. For example, protein evolution by sequence substitution is  
4 constrained by epistasis between residues. In prokaryotes, evolutionary innovation frequently  
5 happens by macrogenomic events such as horizontal gene transfer (HGT). Previous work has  
6 suggested that HGT can be influenced by ancestral genomic content, yet the extent of such gene-  
7 level constraints has not yet been systematically characterized. Here, we evaluated the  
8 evolutionary impact of such constraints in prokaryotes, using probabilistic ancestral  
9 reconstructions from 634 extant prokaryotic genomes and a novel framework for detecting  
10 evolutionary constraints on HGT events. We identified 8,228 directional dependencies between  
11 genes, and demonstrated that many such dependencies reflect known functional relationships,  
12 including, for example, evolutionary dependencies of the photosynthetic enzyme RuBisCO.  
13 Modeling all dependencies as a network, we adapted an approach from graph theory to establish  
14 chronological precedence in the acquisition of different genomic functions. Specifically, we  
15 demonstrated that specific functions tend to be gained sequentially, suggesting that evolution in  
16 prokaryotes is governed by functional assembly patterns. Finally, we showed that these  
17 dependencies are universal rather than clade-specific and are often sufficient for predicting  
18 whether or not a given ancestral genome will acquire specific genes. Combined, our results  
19 indicate that evolutionary innovation via HGT is profoundly constrained by epistasis and  
20 historical contingency, similar to the evolution of proteins and phenotypic characters, and  
21 suggest that the emergence of specific metabolic and pathological phenotypes in prokaryotes can  
22 be predictable from current genomes.

23

## 1 INTRODUCTION:

2 A fundamental question in evolutionary biology is how present circumstances affect future  
3 adaptation and phenotypic change (Gould and Lewontin 1979). Studies of specific proteins, for  
4 example, indicate that epistasis between sequence residues limits accessible evolutionary  
5 trajectories and thereby renders certain adaptive paths more likely than others (Weinreich et al.  
6 2006; Gong et al. 2013; de Visser and Krug 2014; Harms and Thornton 2014). Similarly, both  
7 phenotypic characters (Ord and Summers 2015) and specific genetic adaptations (Christin et al.  
8 2015; Conte et al. 2012) show strong evidence of parallel evolution rather than convergent  
9 evolution. That is, a given adaptation is more likely to repeat in closely related organisms than in  
10 distantly related ones. This inverse relationship between the repeatability of evolution and  
11 taxonomic distance implies a strong effect of lineage-specific contingency on evolution, also  
12 potentially mediated by epistasis (Orr 2005).

13 Such observations suggest that genetic adaptation is often highly constrained and that the  
14 present state of an evolving system can impact future evolution. Yet, the studies above are  
15 limited to small datasets and specific genetic pathways, and a more principled understanding of  
16 the rules by which future evolutionary trajectories are governed by the present state of the system  
17 is still lacking. For example, it is not known whether such adaptive constraints are a feature of  
18 genome-scale evolution or whether they are limited to finer scales. Moreover, the mechanisms  
19 that underlie observed constraints are often completely unknown. Addressing these questions is  
20 clearly valuable for obtaining a more complete theory of evolutionary biology, but more  
21 pressing, is essential for tackling a variety of practical concerns including our ability to combat  
22 evolving infectious diseases or engineer complex biological systems.

1           Here, we address this challenge by analyzing horizontal gene transfer (HGT) in  
2 prokaryotes. HGT is an ideal system to systematically study genome-wide evolutionary  
3 constraints because it involves gene-level innovation, occurs at very high rates relative to  
4 sequence substitution (Nowell et al. 2014; Puigbò et al. 2014a), and is a principal source of  
5 evolutionary novelty in prokaryotes (Gogarten et al. 2002; Jain et al. 2003; Lerat et al. 2005;  
6 Puigbò et al. 2014b). Clearly, many or most acquired genes are rapidly lost due to fitness costs  
7 (van Passel et al. 2008; Baltrus 2013; Soucy et al. 2015), indicating that genes retained in the  
8 long term are likely to provide a selective advantage. Moreover, not all genes are equally  
9 transferrable (Jain et al. 1999; Sorek et al. 2007; Cohen et al. 2011), and not all species are  
10 equally receptive to the same genes (Smillie et al. 2011; Soucy et al. 2015). However,  
11 differences in HGT among species have been attributed not only to ecology (Smillie et al. 2011)  
12 or to phylogenetic constraints (Nowell et al. 2014; Popa et al. 2011), but also to interactions with  
13 the host genome (Jain et al. 1999; Cohen et al. 2011; Popa et al. 2011). Indeed, studies involving  
14 single genes or single species support the influence of genome content on the acquisition and  
15 retention of transferred genes (Pal et al. 2005; Iwasaki and Takagi 2009; Chen et al. 2011; Press  
16 et al. 2013; Sorek et al. 2007; Johnson and Grossman 2014). For example, it has been  
17 demonstrated that the presence of specific genes facilitates integration of others into genetic  
18 networks (Chen et al. 2011), and that genes are more commonly gained in genomes already  
19 containing metabolic genes in the same pathway (Pal et al. 2005; Iwasaki and Takagi 2009).  
20 However, to date, a systematic, large-scale analysis of such dependencies has not been presented.

21           In this paper, we therefore characterize a comprehensive collection of genome-wide  
22 HGT-based dependencies among prokaryotic genes, uncover potential rules of genome evolution  
23 in prokaryotes, and demonstrate that the acquisition of genes is to some extent predictable based

1 on these rules. Overall, our study suggests that genetic innovation and adaptation are  
2 substantially constrained through gene-level epistatic interactions such as those that we describe  
3 influencing HGT.

4

5

## 6 **RESULTS:**

### 7 **PGCE Inference**

8 We first set out to detect pairs of genes for which the presence of one gene in the genome  
9 promotes the gain of the other gene (though not necessarily *vice versa*) (Figure 1). Such “pairs of  
10 genes with conjugated evolution” (PGCEs) represent putative epistatic interactions at the gene  
11 level and may guide genome evolution. To this end, we obtained a collection of 634 prokaryotic  
12 genomes, annotated by KEGG (Kanehisa et al. 2012), and linked through a curated phylogeny  
13 (Dehal et al. 2010). For each of the 5801 genes that varied in presence across these genomes, we  
14 reconstructed the probability of this gene’s presence or absence on each branch of the  
15 phylogenetic tree using a previously introduced method (Cohen and Pupko 2010), as well as the  
16 probability that it was gained along these branches. We further confirmed that genes’  
17 presence/absence was robust to the reconstruction method employed (99.5% agreement between  
18 reconstruction methods used; Methods). From these reconstructions, we estimated the frequency  
19 with which each gene was gained in the presence of each other gene, and followed previous  
20 studies (Maddison 1990; Cohen et al. 2012) in using parametric bootstrapping (Figure S1) to  
21 detect PGCEs – gene pairs for which one gene is gained significantly more often in the presence  
22 of the other (Figure S2, SI Text). In total, we identified 8,415 PGCEs. We finally applied a  
23 transitive reduction procedure to discard potentially spurious PGCEs, resulting in a final network

1 containing 8,228 PGCEs connecting a total of 2,260 genes (Figures S3, S4, SI Text). A detailed  
2 description of the procedures used can be found in Methods, and the final list of PGCEs is  
3 supplied as File S1.

4

#### 5 **PGCEs represent biologically relevant dependencies**

6 Comparing this final set of PGCEs to known biological interactions, we confirmed that the  
7 obtained PGCEs represent plausible biological dependencies. For example, genes sharing the  
8 same KEGG Pathway annotations were more likely to form a PGCE (Figure 2A), as were genes  
9 that are linked in an independently-derived network of bacterial metabolism (Levy and  
10 Borenstein 2013) (Figure 2B). Moreover, PGCEs often linked genes in functionally related  
11 pathways (Figure S5, SI Text). We similarly identified specific examples in which PGCEs  
12 connected pairs of genes with well-described functional relationships. One such example is the  
13 PGCE connecting *rbsL* and *rbsS* (sometimes written *rbcL/rbcS*), two genes that encode the large  
14 and small subunits of the well-described photosynthetic enzyme ribulose-1-5-bisphosphate  
15 carboxylase-oxygenase (RuBisCO), respectively. The *rbsL* subunit alone has carboxylation  
16 activity in some bacteria, but the addition of *rbsS* increases enzymatic efficiency, consistent with  
17 its PGCE dependency on *rbsL* (Figure 3A) (Andersson and Backlund 2008). Moreover, these  
18 genes are known to undergo substantial horizontal transfer (Delwiche and Palmer 1996).

19 Multiple additional genes were found to promote *rbsS* gain (88 PGCEs in total, Table  
20 S1), many of which, as expected, are associated with carbon metabolism. Other genes in this set,  
21 however, unexpectedly implicated nitrogen acquisition, as well as other pathways (Table S2), in  
22 promoting *rbsS* gain. For example, all components of the *urt* urea transport complex had a PGCE  
23 link with *rbsS*, as shown by the reconstructed phylogenetic history of *urtA* and *rbsS* (Figure 3B).

1 This strict dependency could reflect nitrogen's role as a rate-limiting resource for primary  
2 production in phytoplankton and other photosynthetic organisms (Eppley and Peterson 1979;  
3 Sohm et al. 2011). In comparing the reconstructions from which *urtA-rbsS* and *rbsL-rbsS*  
4 dependencies were inferred, we further observed that *rbsS* is gained only in lineages where both  
5 dependencies were previously present. This indicates that while both *rbsL* and *urtA* may be  
6 necessary for the acquisition of *rbsS*, neither *rbsL* nor *urtA* are independently sufficient for the  
7 acquisition of *rbsS*. Other PGCEs may interact in similarly complex fashions in controlling the  
8 acquisition of genes, and thus such relationships may be gene-specific and involve a variety of  
9 biological mechanisms that may be difficult to generalize. For further analyses, we therefore  
10 focused on analyzing large-scale patterns of PGCE connectivity and on exploring how the  
11 dependencies between various genes structure the relationships between functional pathways.

## 12 **PGCE network analyses reveal evolutionary assembly patterns**

13 The *rbsS*-associated PGCEs described above show how PGCEs captured an assembly pattern  
14 involving multiple pathways. Therefore, we next set out to infer global evolutionary assembly  
15 patterns based on the complete set of PGCEs identified. Specifically, we used a network-based  
16 topological sorting approach (SI Text) to rank all genes in the PGCE network. According to this  
17 procedure, genes without dependencies occupy the first rank, genes in the second rank have  
18 PGCE dependencies only on first rank genes, genes in the third rank have dependencies only on  
19 first and second rank genes, and so on until all genes are associated with some rank. In other  
20 words, the obtained ranking represents general patterns in the order by which genes are gained  
21 throughout evolution, with the gain of higher-ranked genes succeeding the presence of the lower-  
22 ranked genes on which they depend. Using this approach, we found that genes could be fully  
23 classified into five ranks (Fig 4A). The first rank was by far the largest at 1,593 genes (most

1 genes do not have detectable dependencies), the second rank had 498 genes, and successive  
2 ranks showed declining membership until the last (fifth) rank, with only 5 genes (Table S3).

3 To identify evolutionary assembly patterns from these ranks, we examined the set of  
4 genes in each rank and identified overrepresented functional categories (Table 1). These enriched  
5 functional categories indicate that certain functional groups of genes consistently occupy specific  
6 positions in these evolutionary assembly patterns, whether in controlling other genes' gain or in  
7 being controlled by other genes. For example, we found that the first rank was enriched for  
8 flagellar and pilus genes involved in motility, in addition to Type II secretion genes (many of  
9 which are homologous to or overlap with genes encoding pilus proteins) and certain two-  
10 component genes. The second rank was enriched for various metabolic processes, whereas later  
11 ranks were enriched for Type III and Type IV secretion systems and conjugation genes (Table 1).  
12 This finding suggests that habitat commitments are made early in evolution, mediated by motility  
13 genes that could underlie the choice and establishment of physical environments. This  
14 environmental choice is followed by a metabolic commitment to exploiting the new habitat. Last,  
15 genes for interaction with the biotic complement of these habitats are gained, and replaced  
16 frequently in response to evolving challenges. Considering two distinct but highly homologous  
17 pilus assembly pathways, one (fimbrial) was enriched in a low rank and one (conjugal) was  
18 enriched in a high rank, suggesting that the specific function of the gene rather than other  
19 sequence-level gene properties drove the ranking (Figure S6A). We additionally confirmed that  
20 the observed rank distribution for these functions is not explained by variation in the frequency  
21 of gene gain (Figure S6B). Furthermore, as expected, we observed that the gains of genes  
22 appearing late in the sort were overrepresented in later branches of the tree compared to the gains



1 of lower-ranked genes (Figures 4B, S7), suggesting that the chronology of gene acquisition  
2 reflects the overall assembly patterns in gain order.

### 3 **Evolution by HGT is predictable**

4 The chronological ordering of ranks was relatively consistent across the tree (Figure 4B),  
5 indicating that PGCE dependencies are universal across prokaryotes. Notably, this universality  
6 also implies that gene acquisition is predictable from genome content. Put differently, if PGCEs  
7 are universal, then PGCEs inferred in one clade of the tree are informative in making predictions  
8 about gene acquisition in a different clade. Indeed, studies of epistasis-mediated protein  
9 evolution indicate that the constriction of possible mutational paths should lead to predictability  
10 in evolution, if epistasis is sufficiently strong (Weinreich et al. 2006). To explore this hypothesis  
11 explicitly, we partitioned the tree into training and test sets (Figure 5A). As test sets, we selected  
12 the Firmicutes phylum, and the Alphaproteobacteria/Betaproteobacteria subphyla. Choosing  
13 whole clades as test sets (rather than randomly sampling species from throughout the tree)  
14 guarantees that true predictions are based on universal PGCEs, rather than clade-specific PGCEs.  
15 For each test set, we used a model phylogeny that excluded the test subtree as a training set, and  
16 inferred PGCEs based on this pruned tree (Table S4, Figure S8A). We then used these inferred  
17 PGCEs to score the likelihood of the gain of dependent genes on each branch in the test set,  
18 based on the genome content of the branch's ancestor (Figure 5A, Table S4, SI Text). We used a  
19 naïve and simplistic score: the proportion of genes upon which the gained gene depends that are  
20 present in the reconstructed ancestor of each branch. In both test sets, we found that prediction  
21 quality was surprisingly high (Figure 5B, Figure S8B), suggesting that PGCEs are taxonomically  
22 universal and statistically robust in describing relationships between genes. This predictability is  
23 consistent with the hypothesis that gene-gene dependencies constrain the evolution of genomes

1 by HGT. More broadly, this analysis and our finding that PGCEs predictably determines future  
2 evolutionary gains provide substantial evidence that the preponderance of parallel evolution over  
3 convergent evolution (Ord and Summers 2015; Conte et al. 2012) may be the result of specific,  
4 identifiable genetic dependencies that similarly impact the evolutionary trajectory taken by  
5 similar genomes.

## 6 **DISCUSSION:**

7 Combined, our findings provide substantial evidence to suggest that gene acquisitions in bacteria  
8 are governed by genome content through numerous gene-level dependencies. Our ability to  
9 detect these underlying dependencies is clearly imperfect, owing to various data and  
10 methodological limitations (SI Text, Figure S2). In reality the complete dependency network is  
11 therefore likely much denser than that described above and includes numerous dependencies and  
12 constraints that our approach may not be able to detect. Consequently, our estimates should be  
13 considered as a lower bound on the extent of gene-gene interactions and accordingly the  
14 predictability of HGT.

15 Notably, even considering such caveats, our observations dramatically expand our  
16 knowledge of the constraints on HGT. Previous studies of such constraints demonstrated that  
17 genes frequently acquired by HGT tend to occupy peripheral positions in biological networks  
18 (Jain et al. 1999; Cohen et al. 2011), are often associated with specific cellular functions, and are  
19 phylogenetically clustered. These observations suggested that properties of transferred genes are  
20 also important determinants of HGT regardless of recipient genome content (Jain et al. 1999;  
21 Cohen et al. 2011; Gophna and Ofran 2011) and that the acquisition of certain genes is clade-  
22 specific (Popa et al. 2011; Andam and Gogarten 2011). In contrast, our analysis demonstrates the  
23 importance of recipient genome content in strongly influencing the propensity of a new gene to

1 be acquired. In fact, to some extent, properties previously reported as determining the general  
2 “acquirability” of genes across all species may reflect some average constraint across genomes.  
3 By considering also variation in genomes acquiring genes, our analysis focused on specific  
4 biological effects, whose strengths may vary from genome to genome.

5       Importantly, our model that gene acquisition is affected by recipient genome content is  
6 consistent with the observed enrichment of HGT among close relatives, which presumably have  
7 similar genome content (Gogarten et al. 2002; Andam and Gogarten 2011; Popa et al. 2011;  
8 Popa and Dagan 2011). This taxonomic clustering of innovation by HGT is also in agreement  
9 with previous studies that demonstrated that phenotypic (Gould and Lewontin 1979; Ord and  
10 Summers 2015) and genetic (Conte et al. 2012; Christin et al. 2015) parallel evolution is more  
11 common than convergent evolution, potentially due to the effects of contingency (Gould and  
12 Lewontin 1979; Conte et al. 2012; Christin et al. 2015; Ord and Summers 2015). However, in  
13 contrast to other studies, we present direct evidence that the mechanism by which contingency  
14 controls evolution is epistasis. Furthermore the universality of the PGCEs shows that the  
15 constraints underlying the effect of contingency operate outside the context of parallel evolution.

16       It should also be noted that while our analysis revealed several intriguing patterns, the  
17 precise interpretation of some of these patterns remains unclear. For instance, the observed  
18 correspondence of topological ranks of genes to chronology suggests that evolutionary age is a  
19 potential contributor to such ranking, especially considering that our reconstructions likely lack  
20 many genes that have not been retained in any extant genomes. However, the biological  
21 plausibility and statistical robustness of PGCEs demonstrated above strongly argue that the  
22 observed evolutionary patterns are the result of constraint-inducing dependencies. Future work

1 may therefore aim to quantify the trade-off between functional and chronological determinants in  
2 apparent evolutionary constraints.

3 Finally, we demonstrate the predictability of genomic evolution by horizontal transfer  
4 from current genomic content. As stated above, this finding also suggests that such dependencies  
5 are fairly universal across the prokaryotic tree. It should be noted that our approach was designed  
6 specifically to understand the PGCE network's significance and universality, rather than predict  
7 gene acquisition. It is likely that an approach specifically engineered for gene acquisition  
8 prediction would substantially outperform our approach. The estimates of predictability of  
9 genomic evolution presented here are accordingly quite conservative.

10 The determinism and predictability of evolutionary patterns therefore appear to be an  
11 outcome not only of intramolecular epistasis in proteins or phylogenetic constraints, but also of  
12 genome-wide interactions between genes. This suggests that the evolution of medically,  
13 economically, and ecologically important traits in prokaryotes depends on ancestral genome  
14 content and is hence at least partly predictable, potentially informing research in the  
15 epidemiology of infectious diseases, bioengineering, and biotechnology.

16

## 17 **METHODS**

18 All mathematical operations and statistical analyses were performed in R 2.15.3 (2012).  
19 Probabilistic ancestral reconstructions were obtained using the *gainLoss* program (Cohen and  
20 Pupko 2010). Phylogenetic simulations and plots were performed with the APE library (Paradis  
21 et al. 2004). Network analyses and algorithms were implemented using either the *igraph* (Csardi  
22 and Nepusz 2006) or *NetworkX* (Hagberg et al. 2013) libraries, and visualized using Cytoscape  
23 v3.1.1 (Shannon et al. 2003).

1

## 2 **Phylogenies**

3 We used a pre-computed phylogenetic tree (Dehal et al. 2010) as a model of bacterial evolution.

4 We mapped all extant organisms in this tree to organisms in the KEGG database by their NCBI

5 genome identifiers, and pruned all tips that did not directly and uniquely map to KEGG. This

6 yielded a phylogenetic tree connecting 634 prokaryotic species. For analyses involving subtrees

7 of this phylogenetic tree, we used iTOL (Letunic and Bork 2011) to extract subtrees.

8

## 9 **Inferring phylogenetic histories for genes**

10 We used the *gainLoss* v1.266 software (Cohen and Pupko 2010), a set of presence/absence

11 patterns of orthologous genes from KEGG (Kanehisa et al. 2012), and the phylogenetic tree

12 described above to infer 1) the probabilities of presence and absence of genes at internal nodes of

13 the tree, 2) gain and loss rates of each gene, and 3) tree branch lengths within a single model. We

14 obtained a probabilistic ancestral reconstruction based on stochastic mapping for each of 5801

15 genes that were present in at least one species and absent in at least one species, and filtered out

16 genes that were found to be gained less than twice throughout the tree, yielding 5031 genes

17 which we further analyzed. We used the probabilities of presence and absence of each of these

18 5031 genes at each node and tip on the tree to compute the probability of each branch

19 experiencing 1) gain (absent in ancestor and present in descendant) and 2) presence (present in

20 both ancestor and descendant; Supporting Text). For a gene X on a branch with ancestor A and

21 descendant B, we assume:

22 1.  $\Pr(X \text{ present on branch}) = \Pr(X \text{ present in } A \cap X \text{ present in } B) =$

23  $\Pr(X \text{ present in } A) * \Pr(X \text{ present in } B)$

1                    2.  $\Pr(X \text{ gained on branch}) = \Pr(X \text{ absent in A} \cap X \text{ present in B}) =$   
2                     $\Pr(X \text{ absent in A}) * \Pr(X \text{ present in B})$

3  
4 Note that these probability estimates are distinct from those obtained by using the *gainLoss*  
5 continuous-time Markov chain on the same ancestral reconstruction, which consider also  
6 hypothetical gains that are not retained and are thus not relevant to our analysis (Supporting  
7 Text).

### 8 **Robustness analysis of reconstruction method**

9 We used a maximum-parsimony reconstruction as inferred by *gainLoss* to benchmark the  
10 accuracy of the *gainLoss* reconstruction by stochastic mapping. In this analysis, only internal  
11 node reconstructions were considered, as tip reconstructions (for which the states are known) are  
12 not informative about algorithm performance. Since the maximum-parsimony reconstruction is  
13 binary (presence/absence) and the stochastic mapping reconstruction is probabilistic, for  
14 purposes of comparison we rounded the probabilities of the stochastic mapping reconstruction to  
15 obtain a presence/absence reconstruction (*i.e.*, a probability  $>0.5$  denotes presence and  $\leq 0.5$   
16 denotes absence). We computed the agreement between the two reconstructions as the  
17 percentage of internal node reconstructions that agree on the state of the gene.

18

### 19 **Quantifying PGCEs**

20 We defined a “pair of genes with conjugated evolution” (PGCE) as a gene pair ( $i, j$ ) for which  
21 the presence of one gene  $i$  encourages the gain of the other,  $j$ . Considering these genes as  
22 phylogenetic characters, we therefore aim to detect pairs for which “gain” state transitions for  
23 character  $j$  are enriched on branches where character  $i$  remains in the “present” state. This

1 problem is related to previous methods for detecting coevolution or correlation between  
2 phylogenetic characters (Maddison 1990; Huelsenbeck et al. 2003; Cohen et al. 2012). Given  $N$   
3 branches and  $k$  genes, there are  $2 N \times k$  matrices,  $P$  and  $G$ , describing the probabilities,  
4 respectively, of presence and gain of each gene along each branch. The test statistic for each  
5 gene pair  $(i, j)$  is the probabilistic count of branches where the gain of gene  $j$  occurs, while  
6 conditioning on the presence of gene  $i$  (cell  $C_{ij}$  in a  $k \times k$  matrix  $C$ ). To compute  $C$  across  $N$   
7 branches, we sum the conditional probabilities of the gain of gene  $j$  in the presence of gene  $i$   
8 across the tree, *i.e.* the products of the two  $N \times k$  matrices,  $P$  (presence) and  $G$  (gain), for each  
9 gene pair:

$$C_{ij} = \sum_{n=1}^N G_{ni} P_{nj}$$

10  
11 Entries in  $C$  which are significantly larger than a null expectation of gains represent PGCEs  
12 between the row and column genes of  $C$ .

13

#### 14 **Null distribution for PGCEs**

15 For two independently evolving genes  $i$  and  $j$ , the counted gains of  $j$  in the presence of  $i$ ,  $C_{ij}$ , will  
16 be distributed under the null hypothesis (independent evolution) as some function of the  
17 prevalence of  $i$  (the sum of  $P_i$ , the vector of probabilities of presence of  $i$  across branches of the  
18 tree), the probabilistic count of gains experienced by  $j$  (the sum of  $G_j$ , the vector of probabilities  
19 of gains of  $j$  across nodes of the tree), and the topology and branch lengths of the tree ( $\tau$ ):

$$C_{ij} \sim f(P_i, G_j, \tau)$$

20 As this distribution may be difficult to formalize for a specific dataset, we followed previous  
21 studies (Cohen et al. 2012; Huelsenbeck et al. 2003; Maddison 1990) and approximated this null

1 distribution via parametric bootstrapping. Specifically, we simulated the evolution of  $10^5$  genes  
2 along the tree using the APE library function *rTraitDisc()* (Paradis et al. 2004). For the gain and  
3 loss rates used in these simulations, we used *gainLoss* gain and loss rates estimated for the 5801  
4 empirical genes. We fit gamma distributions to these values by maximum likelihood using the  
5 function *fitdistr()* from the MASS library (Venables and Ripley 2002). For both gains and losses,  
6 we increased the shape parameter of the gamma distribution (by a factor of 3 for gains, 1.5 for  
7 losses), to ensure that simulated genes showed sufficiently large numbers of gains. This was  
8 necessary because parametric bootstrapping with the rates inferred by *gainLoss* resulted in left  
9 skewed distributions of gene gains (compare Figures S1A, S1C, and S1E), which were likely to  
10 confound null models, whereas for our null models to be applicable for this analysis, the  
11 distribution of simulated gene gains should be roughly similar to the distribution of gains of  
12 empirical genes (see Figure S1, Supplementary Text).

13         These simulated genes should evolve independently and thus represent a null model for  
14 PGCEs. As above, we constructed matrices representing the probabilities of presence and gain of  
15 these  $10^5$  genes across all of the branches of the phylogeny ( $P_{null}$  and  $G_{null}$ ). We then multiplied  
16 these matrices of simulated genes to compute a  $10^5 \times 10^5$  matrix  $C_{null}$  of probabilistic counts. As  
17 a null distribution for each pair of genes  $i$  and  $j$  with  $C_{ij} > 1$  (those with  $C_{ij} \leq 1$  are not  
18 informative), we used the 1000 simulated genes with prevalence closest to gene  $i$  (rows of  $C_{null}$ ),  
19 and the 1000 simulated genes with a number of gains closest to gene  $j$  (columns of  $C_{null}$ ). We  
20 used the  $10^6$  simulated observations in the resulting submatrix of  $C_{null}$  as a null distribution for  
21  $C_{ij}$ . Notably,  $C_{ij}$  represents probabilistic counts, whereas  $C_{null}$  represents integer counts (the true  
22 reconstruction is known). Consequently, we floored values in  $C_{ij}$ , such that all counts were  
23 truncated at the decimal point. The comparison of  $C_{ij}$  to this null distribution yields an empirical



1 p-value; we rejected the null hypothesis of independence between genes  $i$  and  $j$  for the  $C_{ij}$   
2 observation at a 1% false discovery rate (Benjamini and Hochberg 1995) ( $P < 7 \times 10^{-6}$ ).

3

#### 4 **Constructing a PGCE network.**

5 For each entry in  $C_{ij}$  for which we observed a significant association, we recorded an edge from  
6 gene  $i$  to gene  $j$  in a network of PGCEs. To focus purely on direct interactions, we subjected this  
7 network to a transitive reduction (Hsu 1975). This reduction requires a directed acyclic graph  
8 (DAG). To identify the largest possible DAG in our PGCE network, we identified and removed  
9 the minimal set of edges inducing cycles (Supplementary Text). We performed a transitive  
10 reduction of the resulting DAG using Hsu's algorithm (Hsu 1975) (Supplementary Text).

11

#### 12 **Mapping biological information to the network.**

13 We used network rewiring (as implemented in the *rewire()* function of the *igraph* library (Csardi  
14 and Nepusz 2006)) to generate null distributions of the PGCE network by randomly exchanging  
15 edges between pairs of connected nodes, while excluding self-edges. In each permutation, we  
16 performed  $5N$  rewiring operations, where there are  $N$  edges in the network, to ensure sufficient  
17 randomization. To estimate the relationship between the PGCE network and biological  
18 information we calculated the number of edges shared between the PGCE network and a  
19 metabolic network of all bacterial metabolism obtained from KEGG (Kanehisa et al. 2012; Levy  
20 and Borenstein 2013), and the number of edges shared between members of the same functional  
21 pathway as defined by KEGG, in both the original and randomized networks.

22 To determine whether genes with certain functional annotations were more likely to associate  
23 with one another in the PGCE network, we examined the KEGG Pathway annotations of each

1 pair of genes in the network. We counted the number of edges leading from each pathway to  
2 each other pathway, and obtained an empirical p-value for this count by comparing it to a null  
3 distribution of the expected counts obtained by random rewiring as above.

4

#### 5 **Topological sorting of PGCE networks**

6 To identify global patterns in our PGCE network, we performed topological sorting (Kahn 1962)  
7 with grouping. Topological sorting finds an absolute ordering of nodes in a directed acyclic  
8 graph (DAG), such that no node later in the ordering has an edge directed towards a node earlier  
9 in the ordering. Grouping the sort allows nodes to have the same rank in the ordering if  
10 precedence cannot be established between them, giving a unique solution. For a description of  
11 the algorithm used, see Supplementary Text.

12

#### 13 **Prediction of HGT events on branches.**

14 We used the PGCE network to predict the occurrence of specific HGT events (gene acquisitions)  
15 on the tree in the following fashion. We used two test/training set partitions, with the clades of  
16 Firmicutes and the Alpha/Betaproteobacteria as independent test sets, and the training sets as the  
17 rest of the tree without these clades. To “train” PGCE networks, we performed ancestral  
18 reconstruction of gene presence, PGCE inference, and network processing just as for the entire  
19 tree. We only attempted to predict genes with at least one PGCE dependency (“predictable”  
20 genes). We then considered each branch in the test set independently, attempting to predict  
21 whether each predictable gene was gained on that branch based on the reconstructed genome at  
22 the ancestor node. For each predictable gene-branch combination, our prediction score was the  
23 proportion of the predictable gene’s PGCE dependencies that are present in the ancestor. This is

1 the dot product of the gene presence/absence pattern of the ancestor node ( $A_i$  across  $i$  potentially  
2 present genes) and a binary vector denoting which genes in the PGCE network the predictable  
3 gene depends on ( $P_i$  across  $i$  genes in potential PGCEs), scaled by  $P_i$ :

$$score = \frac{\sum A_i P_i}{\sum P_i}$$

4 Note that this value ranges between 0 and 1 for each predicted gene. As true gains, we used our  
5 reconstructed gene acquisition events for each branch in the test set. We arbitrarily called any  
6 predictable gene-branch pair with a  $\text{Pr}(\text{gain}) > 0.5$  as a gain, and any predictable gene-branch  
7 pair with  $\text{Pr}(\text{gain}) \leq 0.5$  as no gain. We filtered out any gene-branch pair where the gene was  
8 known to be present with  $\text{Pr} > 0.4$ , as in these cases the gene is probably already present. We  
9 analyzed the accuracy of our prediction scores using receiver operating characteristic (ROC)  
10 analysis and by comparing scores of the gain branches to those of the no-gain branches.

11

## 12 **Data Access**

13 Data are available at <http://figshare.com/s/1f341994624c11e5b23706ec4bbcf141>, along with  
14 code for performing analyses.

15

## 16 **Acknowledgements**

17 We are obliged to members of the Borenstein and Queitsch laboratories, and to Evgeny  
18 Sokurenko, Joe Felsenstein, and Willie Swanson for helpful discussions. We thank Ofir Cohen  
19 for help with the *gainLoss* program. MOP was supported in part by National Human Genome  
20 Research Institute Interdisciplinary Training in Genome Sciences Grant 2T32HG35-16. CQ is  
21 supported by National Institute of Health New Innovator Award DP2OD008371. EB is  
22 supported by National Institute of Health New Innovator Award DP2AT00780201. We thank

- 1 UW Genome Sciences Information Technology Services for high-performance computing
- 2 resources.
- 3
- 4

## 1 **FIGURE LEGENDS**

2 **Figure 1. Workflow for deriving the PGCE network.** (A): a model phylogeny and a set of  
3 gene presence/absence patterns at the tips are used to generate an ancestral reconstruction, from  
4 which gains are inferred. Filled circles represent the presence of a gene (distinguished by color),  
5 empty circles represent absence of that gene. Inverted triangles represent points on the phylogeny  
6 where the gene of the indicated color is inferred to be gained. (B): Based on inferred gain and  
7 loss rates, many evolutionary scenarios are independently simulated and used as a null  
8 expectation for evolutionary independence. Filled circles indicate presence of the simulated gene  
9 and empty circles indicate absence, inverted triangles represent gains of the simulated gene on  
10 the phylogeny. (C): A null distribution derived from simulated gene evolution is used to identify  
11 dependencies between real genes. (D): These dependencies are modeled as a network. Filled  
12 circles indicate genes (nodes), arrows indicate dependencies (edges).

13  
14 **Figure 2. PGCEs are enriched for biologically meaningful interactions.** (A): The observed  
15 number of PGCE edges connecting genes in the same pathway (dotted line), compared to the  
16 expected distribution obtained from 1000 rewired networks with identical degree distributions.  
17 (B): The observed number of PGCE edges that also appear in a bacteria-wide metabolic network,  
18 compared to the expected distribution.

19  
20 **Figure 3. The phylogenetic history of *rbsL*, *urtA* and *rbsS*.** The presence of each gene in each  
21 branch in the phylogenetic tree is illustrated with a colored circle, with the circle's diameter  
22 scaled to denote the probability of presence. (A): *rbsL* and *rbsS* evolutionary histories; (B): *urtA*

1 and *rbsS* evolutionary histories. The long branch leading to Archaea (bottom-most clade) was  
2 reduced in size for graphical purposes.

3

4 **Figure 4. Topological sorting of the PGCE dependency network reveals assembly patterns**  
5 **that govern the evolutionary process.** (A): Binned dependencies among the six ranks of genes  
6 in the topological sort (left to right). Node size represents the number of genes in each rank  
7 (using natural logarithm-scale). Edge width represents the number of PGCEs between genes in  
8 different rank (natural logarithm-scale), all edges are directed to the right. (B) The gain of genes  
9 from each rank in each branch of the phylogenetic tree is illustrated (circles). The different colors  
10 represent different ranks. Circle sizes correspond to the proportion of gains on a branch  
11 attributed to genes of that rank (e.g. a large red circle indicates that most gains on a branch  
12 correspond to rank 1). The branch to Archaea (lower clade) has been reduced in size for  
13 graphical purposes. See also Figure S7.

14

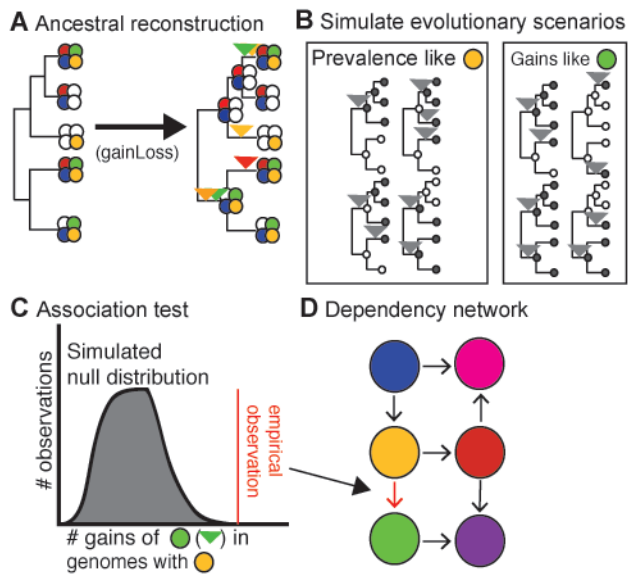
15 **Figure 5. PGCE dependencies lead to taxonomically robust predictability of gene**  
16 **acquisition.** (A): Workflow for predicting gene acquisition between clades of the tree. A training  
17 set is used to build a PGCE dependency model, which is then used to predict on which specific  
18 branches genes are likely to be gained (green circles), based on dependencies inferred from the  
19 training set (red and blue circles). (B): performance of PGCEs in predicting gene acquisitions in  
20 two test sets (indicated clades of the prokaryotic tree). Areas under each curve: Firmicutes, 0.73;  
21 Alpha/Beta-proteobacteria, 0.68. The diagonal dotted line represents the performance of a purely  
22 random prediction. See also Figure S8.

23

24

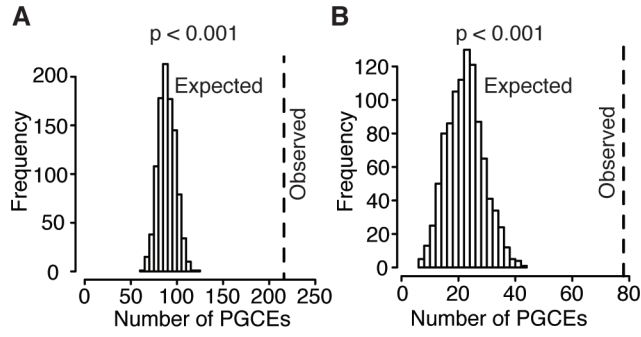
## 1 FIGURES

Figure 1



2  
3

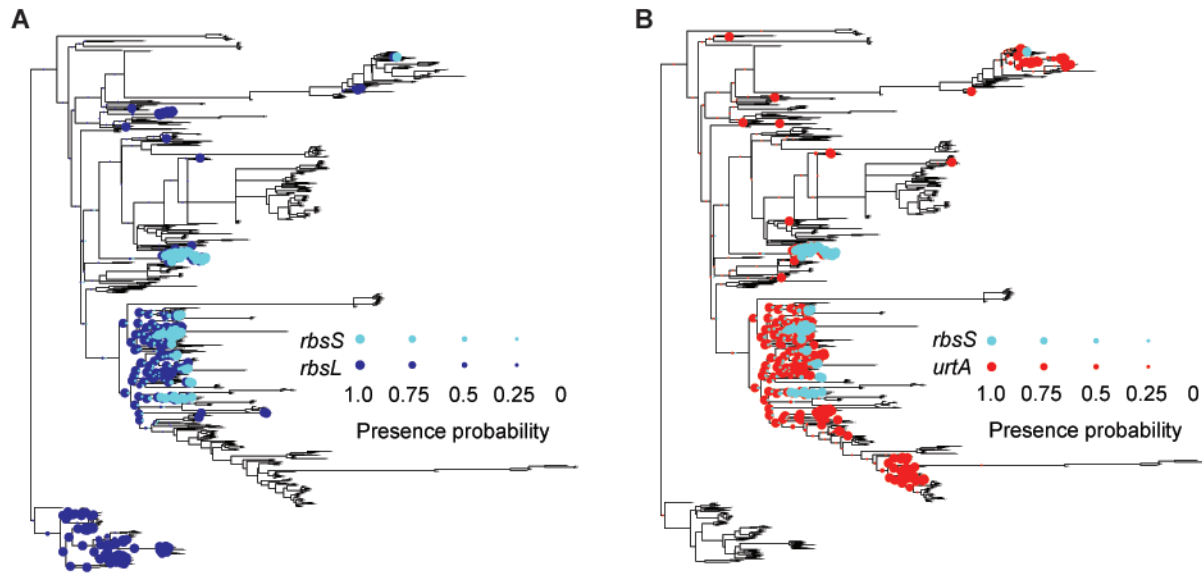
Figure 2



1  
2

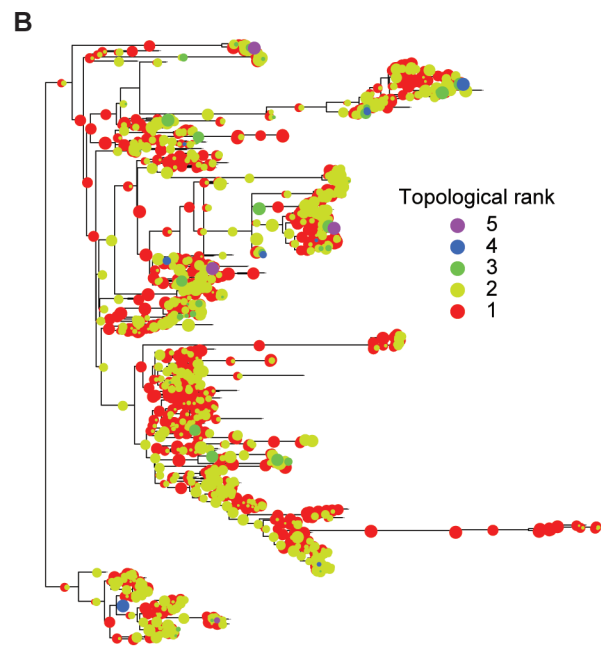
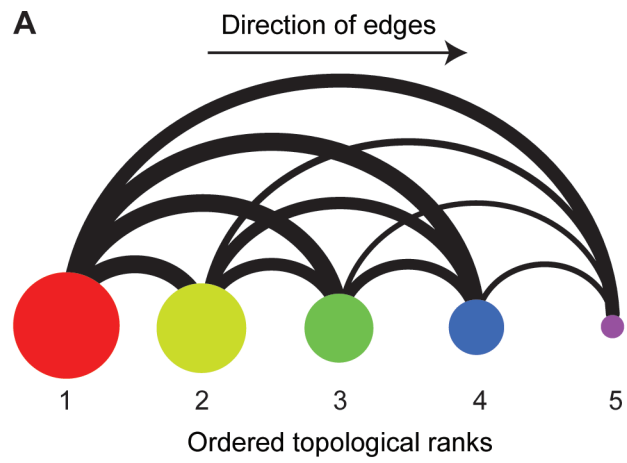


Figure 3



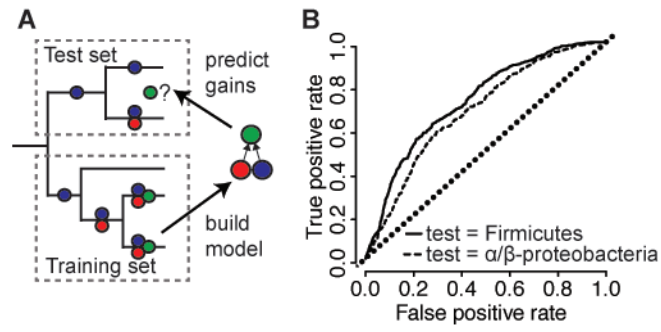
1  
2

Figure 4



1  
2

Figure 5



1  
2

1 **TABLES**

2 **Table 1.** Functional groups are enriched in different ranks of the topological sort.

Annotation label	P-value <sup>1</sup>	Enrichment Ratio <sup>2</sup>
<b>Rank 1 Enrichments</b>		
Cell motility	1.94E-07	1.40
Bacterial motility proteins	1.85E-11	1.41
Type II secretion system	2.61E-05	1.33
Two-component system	3.65E-04	1.25
Flagellar system	1.01E-09	1.43
Pilus system	2.11E-04	1.38
Metabolism <sup>3</sup>	3.37E-05	0.91
Xenobiotics biodegradation and metabolism <sup>3</sup>	1.07E-06	0.69
Carbohydrate metabolism <sup>3</sup>	0.00012	0.84
Type IV secretion system <sup>3</sup>	1.26E-09	0.20
<b>Rank 2 Enrichments</b>		
Metabolism	1.47E-04	1.23
Carbohydrate metabolism	3.08E-06	1.58
<b>Rank 4 Enrichments</b>		
Pathogenicity	1.88E-06	21.6
Conjugal transfer pilus assembly protein	1.08E-04	15.0
Type III protein secretion pathway protein	1.88E-06	21.6
ABC-2 type and other transporters	2.31E-04	12.5
Type IV secretion system	1.30E-03	8.04

3 1: from a hypergeometric test. All annotations displayed are significant at a 1% false discovery rate.

4 2: The ratio of the observed proportion of genes with this label in the indicated rank to the expected proportion  
5 based on all genes in the network.

6 3: These annotations are depleted (i.e. enrichment ratio significantly less than one) in the first rank.

7

## REFERENCES

- Andam CP, Gogarten JP. 2011. Biased gene transfer in microbial evolution. *Nat Rev Microbiol* **9**: 543–55.
- Andersson I, Backlund A. 2008. Structure and function of Rubisco. *Plant Physiol Biochem* **46**: 275–91.
- Baltrus DA. 2013. Exploring the costs of horizontal gene transfer. *Trends Ecol Evol* **28**: 489–95.
- Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J R Stat Soc B* **57**: 289–300.
- Chen HD, Jewett MW, Groisman E a. 2011. Ancestral genes can control the ability of horizontally acquired loci to confer new traits. *PLoS Genet* **7**: e1002184.
- Christin P-A, Arakaki M, Osborne CP, Edwards EJ. 2015. Genetic enablers underlying the clustered evolutionary origins of C4 photosynthesis in angiosperms. *Mol Biol Evol* **32**: 846–58.
- Cohen O, Ashkenazy H, Burstein D, Pupko T. 2012. Uncovering the co-evolutionary network among prokaryotic genes. *Bioinformatics* **28**: i389–i394.
- Cohen O, Gophna U, Pupko T. 2011. The complexity hypothesis revisited: connectivity rather than function constitutes a barrier to horizontal gene transfer. *Mol Biol Evol* **28**: 1481–9.
- Cohen O, Pupko T. 2010. Inference and characterization of horizontally transferred gene families using stochastic mapping. *Mol Biol Evol* **27**: 703–13.
- Conte GL, Arnegard ME, Peichel CL, Schluter D. 2012. The probability of genetic parallelism and convergence in natural populations. *Proc Biol Sci* **279**: 5039–47.
- Csardi G, Nepusz T. 2006. The igraph Software Package for Complex Network Research. *InterJournal Complex Sy*.
- Dehal PS, Joachimiak MP, Price MN, Bates JT, Baumohl JK, Chivian D, Friedland GD, Huang KH, Keller K, Novichkov PS, et al. 2010. MicrobesOnline: an integrated portal for comparative and functional genomics. *Nucleic Acids Res* **38**: D396–400.
- Delwiche CF, Palmer JD. 1996. Rampant horizontal transfer and duplication of rubisco genes in eubacteria and plastids. *Mol Biol Evol* **13**: 873–882.
- Eppley RW, Peterson BJ. 1979. Particulate organic matter flux and planktonic new production in the deep ocean. *Nature* **282**: 677–680.

- Gogarten JP, Doolittle WF, Lawrence JG. 2002. Prokaryotic evolution in light of gene transfer. *Mol Biol Evol* **19**: 2226–38.
- Gong LI, Suchard MA, Bloom JD. 2013. Stability-mediated epistasis constrains the evolution of an influenza protein. *Elife* **2**: e00631.
- Gophna U, Ofra Y. 2011. Lateral acquisition of genes is affected by the friendliness of their products. *Proc Natl Acad Sci U S A* **108**: 343–8.
- Gould SJ, Lewontin RC. 1979. The Spandrels of San Marco and the Panglossian Paradigm: A Critique of the Adaptationist Programme. *Proc R Soc B Biol Sci* **205**: 581–598.
- Hagberg A, Schult D, Swart P. 2013. NetworkX. High productivity software for complex networks. <https://networkx.lanl.gov/>.
- Harms MJ, Thornton JW. 2014. Historical contingency and its biophysical basis in glucocorticoid receptor evolution. *Nature* **512**: 203–7.
- Hsu HT. 1975. An Algorithm for Finding a Minimal Equivalent Graph of a Digraph. *J ACM* **22**: 11–16.
- Huelsenbeck JP, Nielsen R, Bollback JP. 2003. Stochastic Mapping of Morphological Characters. *Syst Biol* **52**: 131–158.
- Iwasaki W, Takagi T. 2009. Rapid pathway evolution facilitated by horizontal gene transfers across prokaryotic lineages. ed. I. Matic. *PLoS Genet* **5**: e1000402.
- Jain R, Rivera MC, Lake JA. 1999. Horizontal gene transfer among genomes: the complexity hypothesis. *Proc Natl Acad Sci U S A* **96**: 3801–6.
- Jain R, Rivera MC, Moore JE, Lake JA. 2003. Horizontal gene transfer accelerates genome innovation and evolution. *Mol Biol Evol* **20**: 1598–602.
- Johnson CM, Grossman AD. 2014. Identification of host genes that affect acquisition of an integrative and conjugative element in *Bacillus subtilis*. *Mol Microbiol*.
- Kahn AB. 1962. Topological Sorting of Large Networks. *Commun ACM* **5**: 558–562.
- Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M. 2012. KEGG for integration and interpretation of large-scale molecular data sets. *Nucleic Acids Res* **40**: D109–14.
- Lerat E, Daubin V, Ochman H, Moran NA. 2005. Evolutionary origins of genomic repertoires in bacteria. ed. D. Hillis. *PLoS Biol* **3**: e130.
- Letunic I, Bork P. 2011. Interactive Tree Of Life v2: online annotation and display of phylogenetic trees made easy. *Nucleic Acids Res* **39**: W475–8.

- Levy R, Borenstein E. 2013. Metabolic modeling of species interaction in the human microbiome elucidates community-level assembly rules. *Proc Natl Acad Sci U S A* **110**: 12804–9.
- Maddison WP. 1990. A Method for Testing the Correlated Evolution of Two Binary Characters: Are Gains or Losses Concentrated on Certain Branches of a Phylogenetic Tree? *Evolution (N Y)* **44**: 539–557.
- Nowell RW, Green S, Laue BE, Sharp PM. 2014. The extent of genome flux and its role in the differentiation of bacterial lineages. *Genome Biol Evol* **6**: 1514–29.
- Ord TJ, Summers TC. 2015. Repeated evolution and the impact of evolutionary history on adaptation. *BMC Evol Biol* **15**: 137.
- Orr HA. 2005. The Probability of Parallel Evolution. *Evolution (N Y)* **59**: 216–220.
- Pal C, Papp B, Lercher MJ. 2005. Horizontal gene transfer depends on gene content of the host. *Bioinformatics* **21**: ii222–ii223.
- Paradis E, Claude J, Strimmer K. 2004. APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics* **20**: 289–290.
- Van Passel MWJ, Marri PR, Ochman H. 2008. The emergence and fate of horizontally acquired genes in Escherichia coli. *PLoS Comput Biol* **4**: e1000059.
- Popa O, Dagan T. 2011. Trends and barriers to lateral gene transfer in prokaryotes. *Curr Opin Microbiol* **14**: 615–23.
- Popa O, Hazkani-Covo E, Landan G, Martin W, Dagan T. 2011. Directed networks reveal genomic barriers and DNA repair bypasses to lateral gene transfer among prokaryotes. *Genome Res* **21**: 599–609.
- Press MO, Li H, Creanza N, Kramer G, Queitsch C, Sourjik V, Borenstein E. 2013. Genome-scale co-evolutionary inference identifies functions and clients of bacterial Hsp90. *PLoS Genet* **9**: e1003631.
- Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin E V. 2014a. Genomes in turmoil: Quantification of genome dynamics in prokaryote supergenomes. *BMC Biol* **12**: 66.
- Puigbò P, Lobkovsky AE, Kristensen DM, Wolf YI, Koonin E V. 2014b. Genomes in turmoil: Quantification of genome dynamics in prokaryote supergenomes. *BMC Biol* **12**: 66.
- Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res* **13**: 2498–504.

Smillie CS, Smith MB, Friedman J, Cordero OX, David LA, Alm EJ. 2011. Ecology drives a global network of gene exchange connecting the human microbiome. *Nature* **480**: 241–4.

Sohm JA, Webb EA, Capone DG. 2011. Emerging patterns of marine nitrogen fixation. *Nat Rev Microbiol* **9**: 499–508.

Sorek R, Zhu Y, Creevey CJ, Francino MP, Bork P, Rubin EM. 2007. Genome-wide experimental determination of barriers to horizontal gene transfer. *Science* **318**: 1449–52.

Soucy SM, Huang J, Gogarten JP. 2015. Horizontal gene transfer: building the web of life. *Nat Rev Genet* **16**: 472–482.

Venables WN, Ripley BD. 2002. *Modern Applied Statistics with S*. Fourth Edi. Springer, Springer.

De Visser JAGM, Krug J. 2014. Empirical fitness landscapes and the predictability of evolution. *Nat Rev Genet* **15**: 480–490.

Weinreich DM, Delaney NF, Depristo MA, Hartl DL. 2006. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science* **312**: 111–4.

2012. R: A language and environment for statistical computing. R Development Core Team.