

Supplementary Material

Computing low-cost rearrangements

For each new occurrence of an architecture, we use simple search or dynamic programming to identify rearrangement solutions with lowest cost. There is a separate procedure for each rearrangement class, as described below. The cost function is the total number of fusion and fission operations. In this section, we consider each architecture to be a sequence or ordered list of domains.

Fusion: Identify parent architectures that are subsequences of the new architecture. Use dynamic programming to select arrangements of parent architectures and new domains, when needed, which can be concatenated into the new architecture using the smallest number of fusion operations. The cost of a solution is equal to the number of sources used minus one.

Fission: Identify parent architectures that are super-sequences of the new architecture. The cost of a solution is one if the new architecture includes an endpoint of the parent architecture. Otherwise, cost is two. Choose the parent architectures with lowest cost.

Insertion: Search for one parent architecture and either a second parent architecture or a new domain, as required by the presence of new domains in the new architecture, such that the second can be inserted into the first. All such pairs are possible solutions with cost three.

Deletion: As in the Fission case, identify parent architectures that are super-sequences of the new architectures. In this case, the pairwise order of domains in the parent architecture should be the same as in the new architecture, but domains are not required to be consecutive in the parent architecture. Use dynamic programming to identify domains that can be removed from the parent architecture with minimum cost to produce the new architecture.

Other: This is the most general case, in which any rearrangement of parent architectures and new domains can be computed. First, we identify all parent architectures that share any domain with the new architecture. Then, we use dynamic programming as in the Fusion case to determine arrangements of domains that can be concatenated to form the new architecture, including all necessary fusion and fission operations into total cost. This procedure differs from Fusion in that for each domain, there may be a cost for splicing the domains from parent architectures. In theory, every domain in the new architecture can be taken from any parent architecture that contains the domain, yielding an exponentially growing number of solutions. To limit the number of solutions, after identifying lowest-cost solutions, we retain those with the smallest number of sources.

Alternative parsimony models

To probe the robustness of our results against changes in the parsimony rule for assigning architectures to internal tree nodes, we compare our MP setting with Dollo parsimony, which is a commonly assumed model, and two variations of MP. In the first variation which we call MP-2child, a parent node with at least two children labeled present for an algorithm will automatically be labeled present, otherwise defaults to standard MP. This variation addresses our concern that MP, typically envisioned with a binary or near-binary tree, is not the ideal rule to assign nodes with higher degree. In the second variation, which we call MP- X where X has a numeric value, we modify the first pass of the internal node assignment algorithm as follows:

- If architecture is present in more than $\max\{X\%, 50\%\}$ of labeled children \rightarrow label parent *present*.
- If architecture is present in less than $\min\{X\%, 50\%\}$ of labeled children \rightarrow label parent *absent*.
- If architecture is present in between $X\%$ and 50% of labeled children \rightarrow label parent *unknown*.

Standard maximum parsimony is equivalent to the case where $X=50$. Decreasing X makes it possible to identify an architecture as present in a parent node as present when fewer children contain that architecture. Increasing X has the opposite effect, making the rule for propagation of architectures up the tree more strict. Using a range of values rather than exact values circumvents automatic propagation of absent/present status. To illustrate this problem, if we were to use the rule that assigns architecture A to a node if more than 49% of children have architecture A present, starting with a leaf node containing A , every ancestor with only two children will be assigned A . This clearly violates the intent of MP to minimize tree cost while making reasonable assignments. We consider values of X between 25 and 75 that correspond to specific combinations at child nodes, e.g. $X=33$ refers to the case with 1 child labeled present and 2 children labeled absent.

The immediate result of modifying the parsimony scheme is new assignments of architectures to nodes.

Figure 1 provides a snapshot of the changes induced by various parsimony rules through the number of architectures assigned to selected nodes. Dollo parsimony is expected to yield the greatest changes because its underlying premise is differs from that of MP—it assigns each architecture to be gained at the Last Common Ancestor of all organisms containing that architecture without considering the tree cost, i.e. the number of gains and losses. Dollo parsimony yields five times as many architectures at the root node as MP and multiple times as many architectures at the Archaea, Bacteria, and Eukaryotes nodes. Nearly all architectures are assigned to internal nodes. These confirm that Dollo parsimony tends to label architectures as being more ancient than MP indicates.

The indirect consequence of new architecture assignments at nodes is changes in the solutions we compute for each architecture. While new architectures will appear at different nodes, our main conclusions regarding the distribution of rearrangements over our defined cases does not change over any parsimony method (Figure 2). Even Dollo parsimony produces fairly small changes in the number of architectures at each rearrangement case.

Variants of MP show little difference in node assignment or distribution of rearrangement cases. The one exception is MP-25 and MP-35 assigning more architectures to the root. This outcome is due to these methods allowing an architecture to be assigned to the root if it is present

in only one of Archaea, Bacteria, and Eukaryotes, which in fact is not ideal and suggests that MP-X with X closer to 50 is preferable.

Figure 1: Number of architectures present in each node or group of nodes. To simplify our analysis, we reduce the TOL into Cellular Organisms (root), Archaea, Bacteria, Eukaryotes, all other interior nodes (interim nodes).

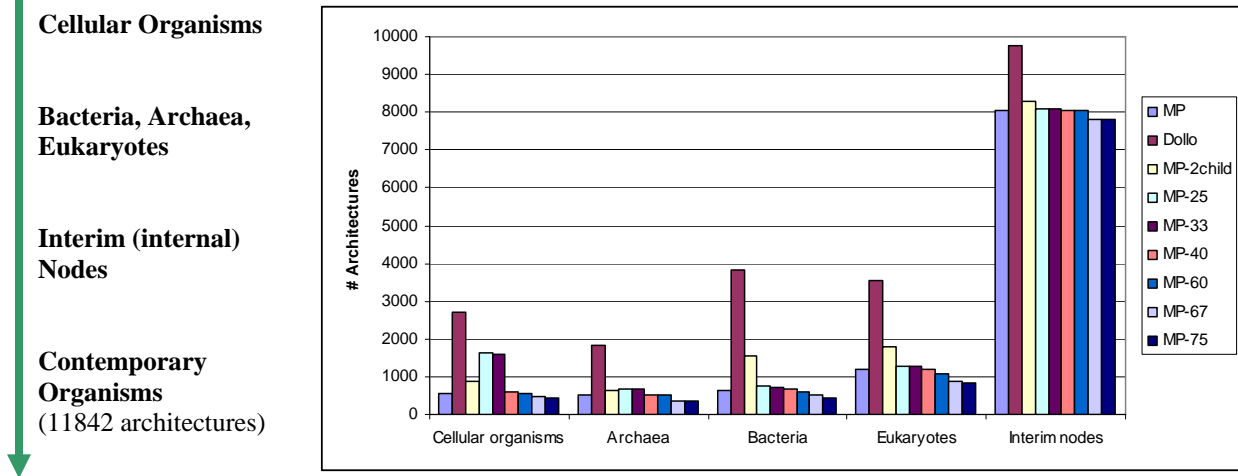
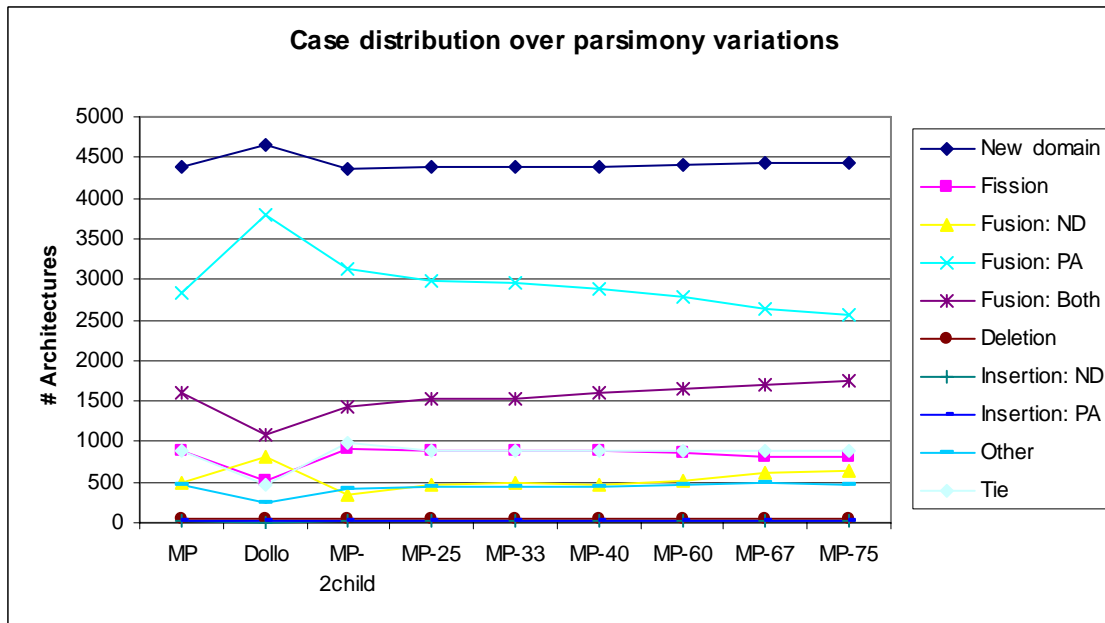


Figure 2: Number of architectures whose most common solution is of each case, with multiple cases among most common solutions counted among the ties. Each line illustrates the values for one rearrangement case.



Supplementary Table 1: Table of species

Anopheles gambiae	Helicobacter hepaticus ATCC 51449
Arabidopsis thaliana	Idiomarina loihiensis L2TR
Caenorhabditis elegans	Lactobacillus plantarum WCFS1
Drosophila melanogaster	Lactococcus lactis subsp. lactis II1403
Encephalitozoon cuniculi	Legionella pneumophila subsp. pneumophila str. Philadelphia 1
Saccharomyces cerevisiae	Leifsonia xyli subsp. xyli str. CTCB07
Plasmodium falciparum	Leptospira interrogans serovar Copenhageni str. Fiocruz L1-130
Schizosaccharomyces pombe	Listeria monocytogenes str. 4b F2365
Homo sapiens	Mannheimia succiniciproducens MBEL55E
Mus musculus	Mesoplasma florum L1
Rattus norvegicus	Mesorhizobium loti MAFF303099
Danio rerio	Methanocaldococcus jannaschii DSM 2661
Oryza sativa	Methanococcus maripaludis S2
Filobasidiella neoformans	Methanopyrus kandleri AV19
Entamoeba histolytica	Methanosarcina acetivorans C2A
Eremothecium gossypii	Methanothermobacter thermautotrophicus str. Delta H
Kluyveromyces lactis	Methylococcus capsulatus str. Bath
Leishmania major	Mycobacterium avium subsp. paratuberculosis K-10
Trypanosoma cruzi	Mycoplasma penetrans HF-2
Yarrowia lipolytica	Nanoarchaeum equitans Kin4-M
Plasmodium berghei	Natronomonas pharaonis DSM 2160
Plasmodium chabaudi	Neisseria meningitidis Z2491
Trypanosoma brucei	Nitrobacter winogradskyi Nb-255
Caenorhabditis briggsae	Nitrosococcus oceani ATCC 19707
Bos taurus	Nitrosomonas europaea ATCC 19718
Canis familiaris	Nocardia farcinica IFM 10152
Dictyostelium discoideum	Nostoc sp. PCC 7120
Apis mellifera	Oceanobacillus iheyensis HTE831
Strongylocentrotus purpuratus	Onion yellows phytoplasma OY-M
Gallus gallus	Parachlamydia sp. UWE25
Xenopus laevis	Pasteurella multocida subsp. multocida str. Pm70
Acinetobacter sp. ADP1	Pelobacter carbinolicus DSM 2380
Aeropyrum pernix K1	Photobacterium profundum SS9
Agrobacterium tumefaciens str. C58	Photorhabdus luminescens subsp. laumondii TTO1
Anabaena variabilis ATCC 29413	Picrophilus torridus DSM 9790
Anaplasma marginale str. St. Maries	Porphyromonas gingivalis W83
Aquifex aeolicus VF5	Prochlorococcus marinus str. MIT 9313
Archaeoglobus fulgidus DSM 4304	Propionibacterium acnes KPA171202
Azoarcus sp. EbN1	Pseudoalteromonas haloplanktis TAC125
Bacillus cereus ATCC 14579	Pseudomonas fluorescens Pf-5
Bacteroides fragilis NCTC 9343	Psychrobacter arcticus 273-4
Bartonella henselae str. Houston-1	Pyrobaculum aerophilum str. IM2
Bdellovibrio bacteriovorus HD100	Pyrococcus furiosus DSM 3638
Bifidobacterium longum NCC2705	Ralstonia eutropha JMP134
Bordetella bronchiseptica RB50	Rhodobacter sphaeroides 2.4.1
Borrelia garinii PBI	Rhodopirellula baltica SH 1
Bradyrhizobium japonicum USDA 110	Rhodopseudomonas palustris CGA009
Brucella suis 1330	Rickettsia felis URRWXCal2
Buchnera aphidicola str. APS (Acyrtosiphon pisum)	

<p> <i>Burkholderia</i> sp. 383 <i>Campylobacter jejuni</i> RM1221 <i>Candidatus Pelagibacter ubique</i> HTCC1062 <i>Carboxydotherrnus hydrogenoformans</i> Z-2901 <i>Caulobacter crescentus</i> CB15 <i>Chlamydia trachomatis</i> D/UW-3/CX <i>Chlamydophila pneumoniae</i> TW-183 <i>Chlorobium tepidum</i> TLS <i>Chromobacterium violaceum</i> ATCC 12472 <i>Clostridium acetobutylicum</i> ATCC 824 <i>Colwellia psychrerythraea</i> 34H <i>Corynebacterium glutamicum</i> ATCC 13032 <i>Coxiella burnetii</i> RSA 493 <i>Dechloromonas aromatica</i> RCB <i>Dehalococcoides ethenogenes</i> 195 <i>Deinococcus radiodurans</i> R1 <i>Desulfotalea psychrophila</i> LSv54 <i>Desulfovibrio vulgaris</i> subsp. <i>vulgaris</i> str. Hildenborough <i>Ehrlichia ruminantium</i> str. Welgevonden <i>Enterococcus faecalis</i> V583 <i>Erwinia carotovora</i> subsp. <i>atroseptica</i> SCRI1043 <i>Escherichia coli</i> CFT073 <i>Francisella tularensis</i> subsp. <i>tularensis</i> SCHU S4 <i>Fusobacterium nucleatum</i> subsp. <i>nucleatum</i> ATCC 25586 <i>Geobacillus kaustophilus</i> HTA426 <i>Geobacter metallireducens</i> GS-15 <i>Gloeobacter violaceus</i> PCC 7421 <i>Gluconobacter oxydans</i> 621H <i>Haemophilus influenzae</i> 86-028NP <i>Haloarcula marismortui</i> ATCC 43049 <i>Halobacterium</i> sp. NRC-1 </p>	<p> <i>Salmonella typhimurium</i> LT2 <i>Shewanella oneidensis</i> MR-1 <i>Shigella flexneri</i> 2a str. 301 <i>Silicibacter pomeroyi</i> DSS-3 <i>Sinorhizobium meliloti</i> 1021 <i>Staphylococcus aureus</i> subsp. <i>aureus</i> MSSA476 <i>Streptococcus pneumoniae</i> TIGR4 <i>Streptomyces coelicolor</i> A3(2) <i>Sulfolobus tokodaii</i> str. 7 <i>Symbiobacterium thermophilum</i> IAM 14863 <i>Synechococcus elongatus</i> PCC 6301 <i>Synechocystis</i> sp. PCC 6803 <i>Thermoanaerobacter tengcongensis</i> MB4 <i>Thermobifida fusca</i> YX <i>Thermococcus kodakarensis</i> KOD1 <i>Thermoplasma volcanium</i> GSS1 <i>Thermosynechococcus elongatus</i> BP-1 <i>Thermotoga maritima</i> MSB8 <i>Thermus thermophilus</i> HB8 <i>Thiobacillus denitrificans</i> ATCC 25259 <i>Thiomicrospira denitrificans</i> ATCC 33889 <i>Treponema denticola</i> ATCC 35405 <i>Tropheryma whipplei</i> TW08/27 <i>Ureaplasma parvum</i> serovar 3 str. ATCC 700970 <i>Vibrio vulnificus</i> YJ016 <i>Wigglesworthia glossinidia</i> endosymbiont of <i>Glossina brevipalpis</i> <i>Wolbachia</i> endosymbiont of <i>Drosophila melanogaster</i> <i>Xanthomonas campestris</i> pv. <i>vesicatoria</i> str. 85-10 <i>Xylella fastidiosa</i> 9a5c <i>Yersinia pestis</i> biovar <i>Medievalis</i> str. 91001 <i>Zymomonas mobilis</i> subsp. <i>mobilis</i> ZM4 </p>
---	---