

## Stepwise formation of the bacterial flagellar system

Renyi Liu, and Howard Ochman

*PNAS* published online Apr 16, 2007;  
doi:10.1073/pnas.0700266104

**This information is current as of June 2007.**

### Supplementary Material

Supplementary material can be found at:  
[www.pnas.org/cgi/content/full/0700266104/DC1](http://www.pnas.org/cgi/content/full/0700266104/DC1)

This article has been cited by other articles:  
[www.pnas.org#otherarticles](http://www.pnas.org#otherarticles)

### E-mail Alerts

Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

### Rights & Permissions

To reproduce this article in part (figures, tables) or in entirety, see:  
[www.pnas.org/misc/rightperm.shtml](http://www.pnas.org/misc/rightperm.shtml)

### Reprints

To order reprints, see:  
[www.pnas.org/misc/reprints.shtml](http://www.pnas.org/misc/reprints.shtml)

Notes:

# Stepwise formation of the bacterial flagellar system

Renyi Liu\* and Howard Ochman\*†‡

Departments of \*Biochemistry and Molecular Biophysics and †Ecology and Evolutionary Biology, University of Arizona, Tucson, AZ 85721

Edited by Francisco J. Ayala, University of California, Irvine, CA, and approved March 8, 2007 (received for review January 11, 2007)

**Elucidating the origins of complex biological structures has been one of the major challenges of evolutionary studies. The bacterial flagellum is a primary example of a complex apparatus whose origins and evolutionary history have proven difficult to reconstruct. The gene clusters encoding the components of the flagellum can include >50 genes, but these clusters vary greatly in their numbers and contents among bacterial phyla. To investigate how this diversity arose, we identified all homologs of all flagellar proteins encoded in the complete genome sequences of 41 flagellated species from 11 bacterial phyla. Based on the phylogenetic occurrence and histories of each of these proteins, we could distinguish an ancient core set of 24 structural genes that were present in the common ancestor to all Bacteria. Within a genome, many of these core genes show sequence similarity only to other flagellar core genes, indicating that they were derived from one another, and the relationships among these genes suggest the probable order in which the structural components of the bacterial flagellum arose. These results show that core components of the bacterial flagellum originated through the successive duplication and modification of a few, or perhaps even a single, precursor gene.**

bacterial evolution | biological complexity | gene duplication

**B**acterial flagella are complex and well honed organelles that provide swimming and swarming motilities and also play a central role in adhesion, biofilm formation, and host invasion (1). In the past several decades, extensive knowledge has accumulated about the structure, genetics, assembly, and regulation of flagella in widely diverse bacterial lineages (2–7). The typical bacterial flagellum consists of six components: a basal body (including MS ring, P ring, and L ring), a motor, a switch, a hook, a filament, and an export apparatus (2). In the best studied systems, those of *Escherichia coli* and *Salmonella enterica* sv. Typhimurium, >50 genes are involved in flagellar biosynthesis and function (3). Approximately half of these genes encode the structural components of the flagellum, and the rest are responsible for either the regulation of flagellar assembly or the detection and processing of environmental signals to which flagella respond.

Whereas *E. coli* and *Salmonella* have long served as the model organisms for studying flagellar assembly (2), there is extensive diversity among bacteria in the contents and organization of the gene complexes that specify flagella as well as structural variation in the flagellum itself (8, 9). For example, in Spirochaetes, flagella are located in the periplasm between the outer membrane sheath and cell cylinder (10); and, in accordance with their location, they have an enlarged C ring and rotor, and have a shape different from that seen in *Salmonella* (11). Furthermore, some bacteria, such as *Vibrio parahaemolyticus*, possess two flagellar systems (polar and lateral) that are encoded by distinct set of genes and use different motive forces (sodium and proton) but share a chemotaxis signal transduction system (12).

The bacterial flagellum has received attention as an exemplum of biological complexity; however, how this complexity and diversification have been achieved remains rather poorly understood. Although several scenarios have been posited to explain how this organelle might have been originated (13), the actual series of evolutionary events that have given rise to the flagellum,

as might be inferred from the relationships of all genes that contribute to the formation and expression of this organelle across taxa, has never been accomplished.

Insights into the evolution of the bacterial flagellum have been gained from the homologies between flagellar proteins and those functioning in other systems (13). For example, the sequence similarity between flagellum-specific ATPase FliI and the  $\beta$ -subunit of ATP synthase led to the speculation that flagellum possibly evolved from this highly conserved, membrane-bound enzyme, whose subunits rotate during catalysis of ATP from ADP (14). Because the flagellar motor proteins MotA/B are homologous to the motor proteins in the Tol-pal and TonB systems (15), the flagellum was hypothesized to have originated as a simple proton-driven secretion system (16). Most significantly, there are well established sequence and structural homologies between bacterial flagella and the type III secretion system (TTSS) demonstrating that the two apparatus derive from a common ancestor (17). Most evidence, including their much broader phylogenetic distribution, supports the view that the flagellum arose much earlier than the TTSS, which are largely limited to Proteobacteria (18–20).

Here, we take advantage of complete genome sequence data to trace the history of each gene involved in the assembly and regulation of the bacterial flagellum. Our results show that flagellum originated very early, before the diversification of contemporary bacterial phyla, and evolved in a stepwise fashion through a series of gene duplication, loss and transfer events. In this article, we focus on the evolution of the core set of flagellar genes that is uniformly present in all flagellated bacteria. The later evolving and lineage-specific components of the flagellar gene complexes remain to be addressed.

## Results

**Defining the Core Set of Flagellar Genes.** By querying the genomes of flagellated bacteria for which complete genome sequences are available, we obtained the phylogenetic distribution of every gene known to be involved in the biosynthesis and regulation of flagella. To investigate the origin and evolution of the bacterial flagellar system, we then applied a phylogenetic profiling method (21) to assort genes into functional groups based on their co-occurrence and shared distributions across genomes. Genes with different functional roles have distinct phylogenetic distributions and profiles; however, most of genes whose protein products constitute the structural components of the flagellum are present in all bacterial phyla considered (Fig. 1). This distribution suggests this core set of structural genes originated before the divergence of the major bacterial lineages and includes 21 genes that specify proteins that form the filament (*fliC*,

Author contributions: R.L. and H.O. designed research; R.L. performed research; R.L. and H.O. analyzed data; and R.L. and H.O. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

Abbreviation: TTSS, type III secretion system.

†To whom correspondence should be addressed. E-mail: hochman@email.arizona.edu.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0700266104/DC1](http://www.pnas.org/cgi/content/full/0700266104/DC1).

© 2007 by The National Academy of Sciences of the USA



hook length control gene *fliK*), have highly variable distributions and are excluded from the core set, even though some of the genes are known to be essential for proper functioning of the flagellar system in a particular species. (The evolutionary histories of these regulatory genes, along with that of a second bacterial flagellar system remain to be described.)

**Phylogenetic Analysis of Flagellar Core Genes.** To ascertain whether the 24 genes that form the flagellar core set have congruent evolutionary histories with one another, we compared the phylogenetic tree inferred for each core gene to that based on concatenated alignments of proteins encoded by 14 of the core genes. (These 14 genes were selected because they were present in all species included in this study and encoded the proteins having a high proportion of alignable positions.) For each of the 24 genes, all branches with >75% bootstrap values agreed with those in the concatenated tree, indicating that no alternative branching orders show strong support and that each of these genes has followed a common history in bacteria since they originated.

**Congruence of Flagellar Genes with Organismal Phylogeny of Bacteria.** The distribution of the 24 core genes among divergent bacterial phyla is most consistent with an ancient origin, pre-dating the shared ancestor of Bacteria. However, the distribution could have been achieved through later horizontal transfer. We tested these alternatives by comparing the phylogeny of the flagellar core proteins with the phylogeny of the corresponding bacterial phyla based on 25 universally distributed genes. The phylogenies are largely congruent on branches that have >75% bootstrap support; however, there are two inconsistencies between the core-gene and the organismal phylogenies; in the placement of both the alphaproteobacterial *Zymomonas mobilis* and a clade of three Betaproteobacteria within the Gammaproteobacteria (Fig. 2). Because individual flagellar genes within the core set show the same evolutionary history (see above), these incongruities have likely resulted from the transfer of the entire flagellar gene complexes between proteobacterial lineages after their separation from other major bacteria groups.

**Core Flagellar Proteins Arose Through the Duplication and Diversification of a Single Precursor.** When each of the 24 core flagellar proteins of *E. coli* are compared (via BLAST) to all proteins encoded in the *E. coli* genome, their best and often only hits are to other core flagellar proteins. Pair-wise comparisons among these core proteins revealed that ten are homologous to other core proteins when applying an *e*-value cutoff of  $10^{-4}$  (Fig. 3). This pattern indicates that the structural genes specifying the portion of flagellum residing outside of cytoplasmic membrane (i.e., the rod, hook, and filament) are paralogs and were derived from one another through duplications.

Aside from these matches to other core proteins, pairwise comparisons of these flagellar proteins to the >4,000 nonflagellar proteins encoded by the entire *E. coli* genome recovered cumulatively a total of only 24 hits that reached the same level of significance. Among these matches, half (including some with *e*-values as low as  $3e^{-10}$  to the flagellar core proteins) are involved in other secretion systems, such as the P pilus and the Type V secretion system, which is consistent with the idea that the flagellum originated as a secretion system. An additional 10 of the 24 hits (with *e*-values ranging from  $10^{-5}$  to  $10^{-6}$ ) are membrane proteins, and the remaining two are prophage tail-fiber proteins. Thus, we conclude that despite their antiquity, the similarities among core proteins to one another are more common and, on average, stronger than to nonflagellar proteins.

Because the genes that constitute the core set are ancient and highly diverged, it is possible that some of the relationships among genes might not be recognized from analyses limited to

the *E. coli* flagellar complex. We repeated this analysis and compared the core gene set of each other flagellated bacterium to all proteins encoded in the corresponding genomes and among themselves, and we obtained a similar result, i.e., the best (and often the only) hits of the flagellar core genes were to other flagellar core genes. However, by extending this analysis beyond *E. coli*, the similarity-relationships and links among several other core genes were resolved. For example, a highly significant match between *fliM* and *fliN* (that was not detected for *E. coli* homologs) was evident in 15 genomes from diverse bacterial subdivisions (Fig. 3). In addition, the interacting export components encoded by *fliP*, *fliR*, and *fliQ* are related based on their protein sequences within several taxa. And even among the 10 *E. coli* core genes that originally showed similarity to one another, there were several new interconnections (e.g., *flgB* to both *flgE* and *flgG*, and between *flgL* and *flgK*) revealed by performing the analysis on other genomes. Cumulatively, each of the 24 core genes shows significant similarity to one or more of the other core genes (Fig. 3), a pattern that would result from their successive origination from one another by independent gene duplications and/or gene fusions.

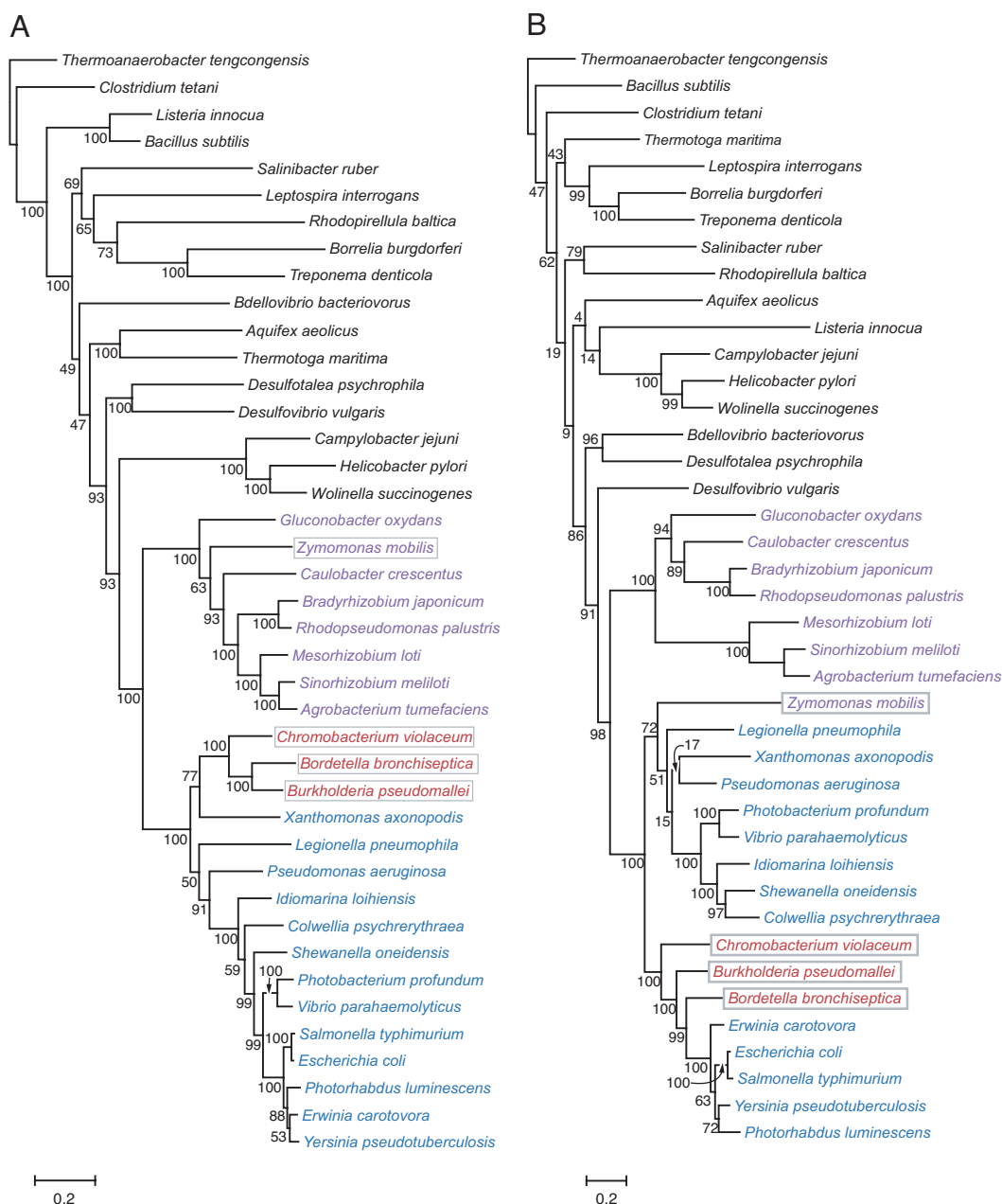
The similarity among the proximal rod protein FlgF, the distal rod protein FlgG, and the hook protein FlgE exemplifies the relationships among these flagellar proteins (Fig. 4). FlgF and FlgG are of similar size (251 aa vs. 260 aa in *E. coli*) and show 31% amino acid identity over their entire lengths. In contrast, the *flgE* gene is much longer and appears to have evolved from *flgG* through an intragenic duplication that added a 160-aa domain to the N terminus of its encoded protein. PSI-BLAST searches reveal two significant alignments between FlgE and FlgG in *E. coli*: one with 24% identity between whole length of FlgG and the C terminus of FlgE (156–401 aa), and the other with 29% identity between the N terminus of two proteins ( $\approx 160$  aa). That *flgE* evolved by a duplication is also supported by the fact that there are two versions of *flgE* in the genus *Bacillus*: among sequenced genomes, four species (*B. subtilis*, *B. clausii*, *B. licheniformis*, and *B. halodurans*) contain a shorter version, which is similar in length to *flgG*, and three species (*B. thuringiensis*, *B. cereus*, and *B. anthracis*) have the longer version.

From the matrix of relationships and protein sequence alignments of the flagellar core genes of *E. coli*, it is also possible to infer the order in which many of these genes and their corresponding structures originated. The low levels of protein identity among these paralogs, paralogous pairs are between 18% and 32% identical, required that we apply a method that combines the output of series of multiple alignment programs to derive a consensus alignment. The alignments on the terminal regions of the proteins, especially at the C terminus, offer the highest confidence. An unrooted neighbor-joining tree and a maximum-likelihood tree [supporting information (SI) Fig. 5] show that the rod proteins originated with either FlgB or FlgC, which are both short proteins, and then generated FlgF and FlgG (and hook protein FlgE) through a series of duplication events. The evolutionary relationships of these flagellar genes parallel the locations of their encoded proteins in contemporary flagella. The proximal, then distal, rod proteins precede (both evolutionarily and physically) the hook proteins, which preceded the hook-filament junction and filament proteins.

## Discussion

Comparisons of the complete genome sequences of flagellated bacteria revealed that the flagellum is based on an ancestral set of 24 core genes for which homologs are present in genomes of all bacterial phyla. The most striking finding from our analysis is that these core genes originated from one another through a series of duplications, an inference based on the fact that they still retain significant sequence homology. The individual core genes show phylogenetic histories congruent with one another,





**Fig. 2.** Congruence between species tree and flagellar protein tree. (A) Species tree based on concatenated protein alignment of 25 single-copy proteins. (B) Flagellar protein tree based on concatenated protein alignment of 14 flagellar core proteins. Bacterial groups are shaded to highlight incongruencies resulting from gene transfer events.

and this core flagellar phylogeny is largely consistent in its deepest branches with the phylogenetic relationships as currently resolved for Bacteria. Taken together, these results indicate that the core set of flagellar genes arose and was assembled from a single or few ancestral sequences, and that the individual genes diversified, before the shared ancestor of Bacteria.

Although sequence similarities among some of the rod and hook proteins were noted in early analyses (24), the degree of paralogy for the ancestral set of flagellar genes, and its implications for the origins of the bacterial flagellum, have gone unrecognized. From a phylogeny of these core proteins, it is possible to reconstruct the order in which they appeared, which in turn, can help elucidate the progression by which the flagellum was originally formed. Based on their relationships and on the

physical locations of proteins forming the flagellum, the rod, hook, and filament proteins originated in an order that mirrors the “inside-out” flagellar assembly process (2, 3). The earliest proteins are proximate to the cytoplasmic membrane with later proteins situated distally, first spanning the outer membrane and then giving rise to structures (i.e., the hook, junction, filament, and capping proteins) that extend outside of the bacterial cell. Thus, the flagellum represents a case whereby its order of assembly recapitulates its evolutionary history.

The structural features of the flagellum, along with the evidence of homology between FliH and ATP synthase subunits and between MotA/B and the secretion proteins TolQ-TolR, suggests that it originated as a primitive secretion system (16), first involving ATPase and then adding the rod, hook, and



