# Rearrangements between Differently Replicating DNA Strands in Asymmetric Bacterial Genomes

DOROTA MACKIEWICZ[1], PAWEŁ MACKIEWICZ[1], MARIA KOWALCZUK[1],
MAŁGORZATA DUDKIEWICZ[1], MIROSŁAW R. DUDEK[2] and STANISŁAW CEBRAT[1*]

[1] Institute of Genetics and Microbiology, Wrocław University,
ul. Przybyszewskiego 63/77, 51-148 Wrocław, Poland;
[2] Institute of Physics, University of Zielona Góra,
ul. Wojska Polskiego 69, 65-246 Zielona Góra, Poland,

Abstract

Many bacterial genomes are under asymmetric mutational pressure which introduces compositional asymmetry into DNA molecule resulting in many biases in coding structure of chromosomes. One of the processes affected by the asymmetry is translocation changing the position of the coding sequence on chromosome in respect to the orientation on the leading and lagging DNA strand. When analysing sets of paralogs in 50 genomes, we found that the number of observed genes which switched their positions on DNA strand is lowest for genomes with the highest DNA asymmetry. However, the number of orthologs which changed DNA strand increases with the phylogenetic distance between the compared genomes. Nevertheless, there is a fraction of coding sequences that stay on the leading strand in all analysed genomes, whereas there are no sequences that stay always on the lagging strand. Since sequences diverge very fast after switching the DNA strand, this bias in mobility of sequences is responsible, in part, for higher divergence rates among some of coding sequences located on the lagging DNA strand.

K e y   w o r d s:  DNA asymmetry, divergence, leading, lagging strand, mutation pressure, rearrangements

## Introduction

Rearrangements are common in bacterial genomes (M u s h e g i a n  and  K o o n i n, 1996; T a t u s o v  *et al.*, 1996; K o l s t o, 1997; W a t a n a b e  *et al.*, 1997; B e l l g a r d  *et al.*, 1999; I t o h  *et al.*, 1999; H u g h e s,  2001) but this phenomenon has not been analysed with respect to leading/lagging strand asymmetry of bacterial chromosomes which seems to be a characteristic (if not universal) feature of these genomes (*e.g.*

---

* Address for correspondence: cebrat@microb.uni.wroc.pl; tel. + 48-71-3756-303; fax: + 48-71-3252-151

L o b r y, 1996; F r e e m a n  *et al.*, 1998; G r i g o r i e v, 1998; M c L e a n  *et al.*, 1998; M a c k i e w i c z *et al.*, 1999a; R o c h a  *et al.*, 1999; T i l l i e r  and  C o l l i n s, 2000a; see for review: F r a n c i n o  and  O c h m a n, 1997; M r a z e k  and  K a r l i n, 1998; F r a n k  and  L o b r y, 1999; K o w a l c z u k  *et al.*, 2001a). Rearrangements of genes in bacterial chromosomes follow very specific rules. In very closely related genomes, many observed rearrangements are symmetric with respect to the origin or terminus of replication (E i s e n  *et al.*, 2000; R e a d  *et al.*, 2000; T i l l i e r  and C o l l i n s, 2000b; S u y a m a  and  B o r k, 2001). T i l l i e r  and  C o l l i n s (2000b) claim that such rearrangements are a result of higher frequency of recombination events at the replication forks which might be recombination hot spots. Another explanation involves the role of selection, and is supported by many genetic and experimental analyses (S c h m i d  and  R o t h, 1983; M a h a n  and  R o t h, 1988; 1991; R e b o l l o  *et al.*, 1988; S e g a l l  *et al.*, 1988; S e g a l  and  R o t h, 1989; F r a n c o i s  *et al.*, 1990; L i u  and  S a n d e r s o n, 1995; 1996; S a n d e r s o n  and L i u, 1998; A l o k a m  *et al.*, 2002). The distance from the origin of replication determines copy number of a gene (dosage effect). Thus, genes should be located in optimal distances from the origin, according to their required expression level. There is also a trend to keep the same size of both replichores which ensures the shortest time of chromosome replication. Furthermore, since inversions of sequences resulting in switching the position of the coding sequence with respect to leading/lagging role of DNA strand is connected with a higher mutational pressure (T i l l i e r  and C o l l i n s, 2000c; R o c h a  and  D a n c h i n, 2001; S z c z e p a n i k  *et al.*, 2001), there could be a higher probability that such a sequence will be eliminated by selection (M a c k i e w i c z  *et al.*, 2001a). (In the terminology, a coding sequence is supposed to be positioned on the leading strand if its sense strand is on the leading DNA strand, respectively the same for the lagging DNA strand). An inversion of a chromosome fragment which encompasses the origin or the terminus of replication does not change the positions of sequences in respect to the leading/lagging role of the DNA strand (M a c k i e w i c z  *et al.*, 2001b). This could lead to a bias in the observed rearrangements. Actually, experimental analyses have shown that permissive (viable) chromosome rearrangements include the origin or terminus of replication (S c h m i d and  R o t h, 1983; M a h a n  and  R o t h, 1991; A l o k a m  *et al.*, 2002). Nevertheless, this feature of keeping the same distance from the origin of replication disappears very fast with phylogenetic distance between analysed genomes which leaves an impression that there is no structural correlation between chromosomes of distant genomes (E i s e n  *et al.*, 2000; T i l l i e r  and  C o l l i n s, 2000b). On the other hand, there are some other phenomena, which could introduce some correlation or structural bias across genomes even at higher phylogenetic distances. Such a phenomenon is a differentiated mutational pressure for coding sequences located on the leading and the lagging strands (T i l l i e r  and  C o l l i n s, 2000c; R o c h a  and D a n c h i n, 2001; S z c z e p a n i k  *et al.*, 2001). There appears to be some preference in the accumulation of translocated coding sequences from the lagging to the leading strand rather than in the opposite direction (M c I n e r n e y, 1998; M a c k i e w i c z  *et al.*, 2001a). Again, mechanisms of selection are blamed for this

bias rather than bias in frequency of translocations themselves. A significant surplus of genes on the leading strand has been observed in many genomes (B r e w e r, 1988; F r a s e r *et al.*, 1995; K u n s t *et al.*, 1997; F r e e m a n *et al.*, 1998; M c L e a n *et al.*, 1998). Knowing that the divergence rate of coding sequences depends on their location on the leading/lagging DNA strand (S z c z e p a n i k *et al.*, 2001), we should expect also a correlation between the function of genes and their position on chromosome as well as differentiated frequency of switching the position of genes lying on the two DNA strands.

One of the main mechanisms of genome evolution is gene duplication, which enables further independent evolution of the structure and function of the two copies (O h n o, 1970). These copies can be seen in genomes as paralogs – homologous sequences occurring in the same genome (F i t c h, 1970). It was found that both, duplication and elimination of paralogs should be ruled by some strict mechanisms, since the number of paralogs follows a very specific numerical law (H u y n e n and N i m w e g e n, 1998; S l o n i m s k i *et al.*, 1998; Q i a n *et al.*, 2001). What we observe is a final result of duplication itself and the paralogs elimination. Duplication of sequences could be connected with a transfer of a new copy into the other DNA strand (inversion) or the copy could stay at the same strand. The mutation rate in sequences after inversion is higher, thus there should be a higher elimination rate of inverted copies. We have already shown that it is true (M a c k i e w i c z *et al.*, 1999a). Genes which have switched DNA strand accommodate very quickly to a new mutational pressure and, in respect to their nucleotide composition, become similar to genes of the new strand (L a f a y *et al.*, 1999; T i l l i e r and C o l l i n s, 2000c; R o c h a and D a n c h i n, 2001).

In this paper we present the results of analysis of fully sequenced bacterial genomes which revealed asymmetry in frequency of translocations (viable inversions) of genes lying on the leading and the lagging DNA strands and we have shown how this affects the divergence rate of genes classified according to the criteria of their mobility.

## Experimental

### Materials and Methods

**Data for analysis.** Prokaryotic genomic sequences and gene annotations have been downloaded from the Genbank (ftp://www.ncbi.nlm.nih.gov). Boundaries between leading and lagging strands (positions of origins and termini of replication) and decisions concerning the location of genes on one of these strands were set on the basis of experimental results or on the basis of the results of DNA walks describing nucleotide compositional bias of differently replicating DNA strands (M a c k i e w i c z *et al.*, 1999b, see also: http://smorfland.microb.uni.wroc.pl). The asymmetry of the genomes was measured by the absolute value of the difference between the GC3 skews of the genes in the leading strand and the ones in the lagging strand:

$$\Delta GC3 \text{ skew} = |(G_d - C_d)/(G_d + C_d) - (G_g - C_g)/(G_g + C_g)|$$

where: $G_d$ and $C_d$ – numbers of guanine and cytosine in the third codon positions of the leading strand genes; $G_g$ and $C_g$ – numbers of guanine and cytosine in the third codon positions of the lagging strand genes. The AT skew and GC skew values proved to be good parameters describing asymmetry of DNA strands (L o b r y, 1996).

   Paralogs for 50 genomes (listed in Table I) showing leading/lagging strand asymmetry were extracted from the TIGR database (http://www.tigr.org). In the analysis only paralogs with minimum 50% identity were chosen.

   Classification of genes to orthologous groups and their amino acid sequences were extracted from Clusters of Orthologous Groups (COGs) downloaded from ftp://www.ncbi.nlm.nih.gov/pub/COG in September 2001. COGs contain protein sequences which are supposed to have evolved from one ancestral protein (K o o n i n  *et al.*, 1998; T a t u s o v  *et al.*, 2001). In the analyses only the best matches for each ortholog (the closest orthologs) have been chosen.

   Analyses of all orthologous sequences have been done on the two sets of bacterial genomes showing evident compositional asymmetry between leading and lagging strands.

   – 7 genomes belonging to γ-subdivision of Proteobacteria group compared with each other: *E. coli* K12-MG1655 (EcK), *E. coli* O157:H7 EDL933 (EcE), *H. influenzae* (Hi), *P. multocida* (Pm), *P. aeruginosa* (Pa), *V. cholerae* (Vc), *X. fastidiosa* (Xf);

   – 14 genomes compared with *E. coli* O157:H7 EDL933 (EcE): *E. coli* K12-MG1655 (EcK), *V. cholerae* (Vc), *P. multocida* (Pm), *P. aeruginosa* (Pa), *X. fastidiosa* (Xf), *N. meningitidis* MC58 (Nm), *B. subtilis* (Bs), *R. prowazekii* (Rp), *M. tuberculosis* H37Rv (Mt), *C. jejuni* (Cj), *T. pallidum* (Tp), *H. pylori* 26695 (Hp), *C. pneumoniae* CWL029 (Cp), *P. horikoshii* (Ph).

   Moreover, from the 7 genomes of the γ-Proteobacteria group, the 7 sets of 1521 orthologs present in all the genomes, being the "best hits" for *E. coli* EDL933 sequences (the closest orthologs), were withdrawn. Similarly, from the set of 14 genomes compared with *E. coli* EDL933, the 14 sets of 233 orthologs present in all the genomes, being the "best hits" for *E. coli* EDL933 sequences, were extracted.

   For each pair of genomes, orthologs and paralogs were classified into three groups according to their strand location: pairs of sequences lying on the leading strands, pairs of sequences lying on lagging strands, and pairs of sequences of which one is lying on the leading and the other on the lagging strand. For each case fractions of the three groups of sequences have been counted.

   **Phylogenetic analysis.** The amino acid sequences of each COG were aligned by the CLUSTAL W 1.8 v. program (T h o m p s o n  *et al.*, 1994). Pairwise evolutionary distances (expressed by the mean number of amino acid substitutions per site) between sequences of each COG were calculated using the WAG model of amino acid substitution (W h e l a n  and  G o l d m a n, 2001) as implemented in the TREE-PUZZLE program version 5.0 (S c h m i d t  *et al.*, 2002). The analyses of divergence of the three groups of orthologs were shown for the sets of 1521 orthologs present in all 7 γ-Proteobacteria genomes.

   For each of the three groups of orthologs a mean value of the evolutionary distances was calculated. Nonparametric analyses by Mann-Whitney U, Kolmogorov-Smirnov and ANOVA Kruskal-Wallis tests (S o k a l  and  R o h l f, 1995) were carried out to assess statistical significance of differences between these groups.

   Evolutionary distances between 16S rRNA sequences (measured by the number of substitutions per site) were calculated by the MEGA 2.1 program (K u m a r  *et al.*, 1993) assuming Tamura-Nei model of nucleotide substitutions (T a m u r a  and  N e i, 1993).

## Results and Discussion

   In highly asymmetric genomes, the mutational pressure after inversion should be relatively higher than for genomes with low asymmetry – there are stronger differences in substitution rates for the leading and lagging DNA strands in the asymmetric genomes (K o w a l c z u k  *et al.*, 2001b; R o c h a  and  D a n c h i n, 2001). Thus, we have anticipated and found a negative correlation between the chromosome asymmetry and the frequency of occurring paralogs in the trans-positions in the genome (one paralog on the leading strand, the other one on the lagging strand – we call these sequences "trans-paralogs"). In Table I we show data for each analysed genome and

in Fig. 1 we show the relation between the fraction of trans-paralogs in the genome and the asymmetry of chromosomes measured by ΔGC3 skew. The observed negative correlation (Spearman correlation coefficient, $r = -0.715$) is statistically significant with high confidence ($p = 5.6 \times 10^{-9}$). There are two possible explanations for the observed negative correlation. One, assuming a higher mutation rate and in consequence higher elimination rate of gene copies translocated to the other DNA strand in highly asymmetric genomes. The second, to us less plausible, refers to the influence of frequency of rearrangements on the maintenance of chromosomal asymmetry. If a global frequency of rearrangements in a genome is low, it does not disturb chromosomal asymmetry established by the mutational pressure. On the contrary, high frequency of rearrangements should diminish this asymmetry.

We have performed a pairwise analysis of orthologs found in compared genomes belonging to γ-Proteobacteria. For each pair of genomes, the orthologs were divided into three groups: i/ pairs of orthologs which are in both compared genomes on the leading strand, ii/ pairs of orthologs which are in both genomes on the lagging strand and iii/ pairs of orthologs of which one is located on the leading and the second on the
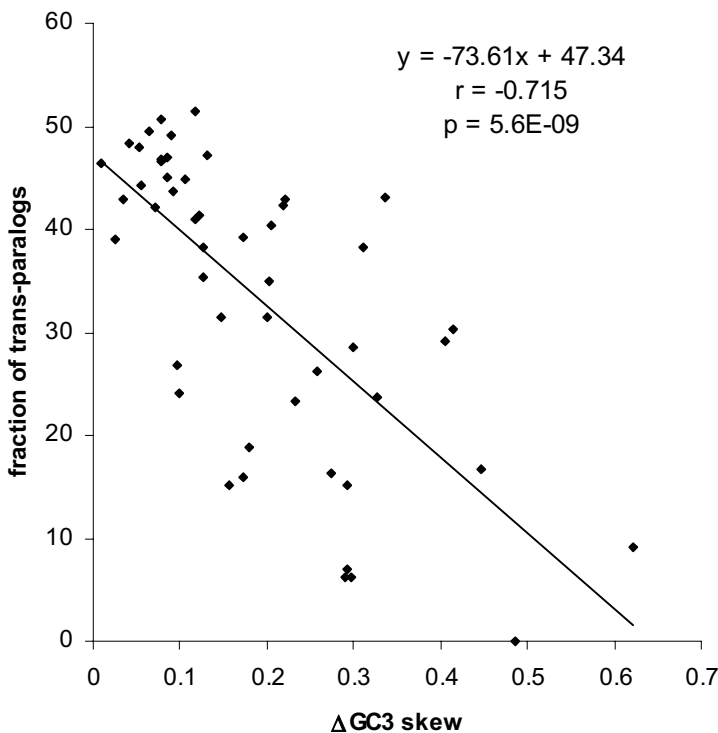


Fig. 1.  Relation between the fraction of trans-paralogs (one paralog on the
leading strand, the other one on the lagging strand) in 50 analysed genomes
and the asymmetry of chromosomes measured by ΔGC3 skew.
Spearman correlation coefficient (r) and its statistical significance (p) are shown.

Table I
Number of all paralogs, the fraction of trans-paralogs and DGC3 skew for 50 analysed genomes

| genome | number of all paralogs | fraction of trans-paralogs | ΔGC3 skew |
|---|---|---|---|
| *Agrobacterium tumefaciens C58 Cereon* | 1096 | 46.8 | 0.08 |
| *Agrobacterium tumefaciens C58 Uwash* | 1117 | 46.6 | 0.08 |
| *Bacillus halodurans C-125* | 2421 | 39.2 | 0.17 |
| *Bacillus subtilis 168* | 558 | 31.4 | 0.15 |
| *Borrelia burgdorferi B31* | 11 | 9.1 | 0.62 |
| *Brucella melitensis 16M* | 314 | 38.2 | 0.13 |
| *Campylobacter jejuni NCTC 11168* | 79 | 30.4 | 0.41 |
| *Caulobacter crescentus CB15* | 511 | 44.2 | 0.05 |
| *Chlamydia muridarum Nigg* | 19 | 0.0 | 0.49 |
| *Chlamydia pneumoniae AR39* | 111 | 6.3 | 0.29 |
| *Chlamydia pneumoniae CWL029* | 112 | 6.3 | 0.30 |
| *Chlamydia pneumoniae J138* | 100 | 7.0 | 0.29 |
| *Chlamydia trachomatis serovar D* | 6 | 16.7 | 0.45 |
| *Clostridium perfringens 13* | 217 | 29.0 | 0.41 |
| *Deinococcus radiodurans R1* | 282 | 46.5 | 0.01 |
| *Escherichia coli O157:H7 EDL933* | 3604 | 26.9 | 0.10 |
| *Escherichia coli K12-MG1655* | 919 | 47.0 | 0.09 |
| *Escherichia coli VT2-Sakai* | 4020 | 24.1 | 0.10 |
| *Haemophilus influenzae KW20* | 73 | 15.1 | 0.16 |
| *Helicobacter pylori 26695* | 198 | 51.5 | 0.12 |
| *Helicobacter pylori J99* | 109 | 41.3 | 0.12 |
| *Lactococcus lactis IL1403* | 811 | 42.9 | 0.22 |
| *Listeria innocua CLIP 11262* | 349 | 18.9 | 0.18 |
| *Listeria monocytogenes EGD-e* | 255 | 31.4 | 0.20 |
| *Mesorhizobium loti MAFF303099* | 1414 | 42.9 | 0.04 |
| *Mycobacterium leprae TN* | 121 | 47.1 | 0.13 |
| *Mycobacterium tuberculosis CDC1551* | 2417 | 43.7 | 0.09 |
| *Mycobacterium tuberculosis H37Rv* | 2279 | 45.1 | 0.08 |
| *Neisseria meningitidis MC58* | 595 | 35.0 | 0.20 |
| *Neisseria meningitidis Z2491* | 874 | 42.3 | 0.22 |
| *Pasteurella multocida PM70* | 86 | 23.3 | 0.23 |
| *Pseudomonas aeruginosa PAO1* | 786 | 44.8 | 0.11 |
| *Pyrococcus abyssi GE5* | 118 | 48.3 | 0.04 |
| *Pyrococcus horikoshii shinkaj OT3* | 147 | 42.2 | 0.07 |
| *Ralstonia solanacearum GMI1000* | 784 | 50.8 | 0.08 |
| *Rickettsia conorii Malish 7* | 478 | 40.4 | 0.21 |

Table I continued

| genome | number of all paralogs | fraction of trans-paralogs | ΔGC3 skew |
|---|---|---|---|
| *Salmonella enterica Typhi CT18* | 679 | 35.3 | 0.13 |
| *Salmonella typhimurium LT2 SGSC1412* | 1280 | 41.0 | 0.12 |
| *Sinorhizobium meliloti 1021* | 460 | 48.0 | 0.05 |
| *Staphylococcus aureus Mu50* | 228 | 16.2 | 0.28 |
| *Staphylococcus aureus N315* | 482 | 15.1 | 0.29 |
| *Streptococcus pneumoniae R6* | 759 | 28.5 | 0.30 |
| *Streptococcus pneumoniae TIGR4* | 479 | 38.2 | 0.31 |
| *Streptococcus pyogenes SF370 M1* | 130 | 26.2 | 0.26 |
| *Thermoplasma acidophilum DSM 1728* | 41 | 39.0 | 0.03 |
| *Thermotoga maritima MSB8* | 216 | 49.5 | 0.06 |
| *Treponema pallidum Nichols* | 72 | 43.1 | 0.34 |
| *Vibrio cholerae El Tor N16961* | 868 | 16.0 | 0.17 |
| *Xylella fastidiosa 9a5c* | 779 | 23.7 | 0.33 |
| *Yersinia pestis CO92* | 8551 | 49.1 | 0.09 |

The set of paralogs (with minimum 50 % identity) was extracted from TIGR database.

lagging strand. If we assume that there is no bias in the frequency of inversions of genes located on the leading and on the lagging DNA strands, we should expect that the fractions of orthologs staying at the same strand in both genomes of the compared pair would decrease with the phylogenetic distance between genomes but the decrease should be proportional to the initial values on the two strands. The results of analyses do not follow these expected rules.

For each pair of compared genomes we have plotted (Fig. 2) the fractions of the three groups of orthologs against the evolutionary distance measured by divergence of 16S rRNA genes between the two compared genomes. The fraction of orthologs lying on the same strand decreases with evolutionary distance while fraction of orthologs which have switched their strands increases rapidly with divergence and become saturated for long evolutionary distances. The same results we have obtained for similar analysis when we compared the *E. coli* EDL933 genome with 14 other genomes belonging to different taxonomic groups (Fig. 3). Even at a short distance (up 0.22 of divergence of 16S rRNA), the total fraction of sequences which switched their strand reaches almost 50%. But there is a very biased input of sequences from the leading and the lagging DNA strands into this fraction. While the fraction of sequences which stay at the leading strands in both compared genomes drops to about 70% of the initial value in the most distant pair, the relative numbers for the lagging strand are up to 40%. These results suggest that the sequences lying on the lagging strand are much more prone to inversions than the sequences lying on the leading strand.
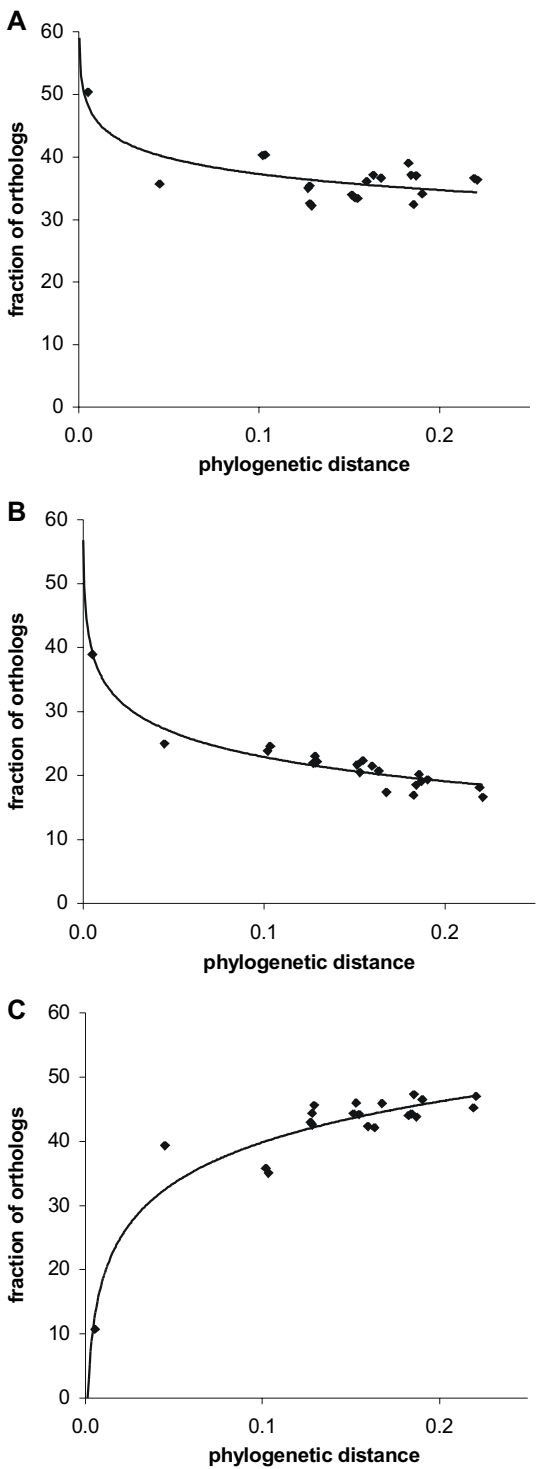
Fig. 2. Relation between the fractions of orthologs and the phylogenetic distance measured by 16S rRNA performed for three groups of orthologs: lying on the leading strand (A), lying on the lagging strand (B) and which changed DNA strand (C).

Data obtained from pairwise comparison of 7 genomes belonging to g- Proteobacteria.
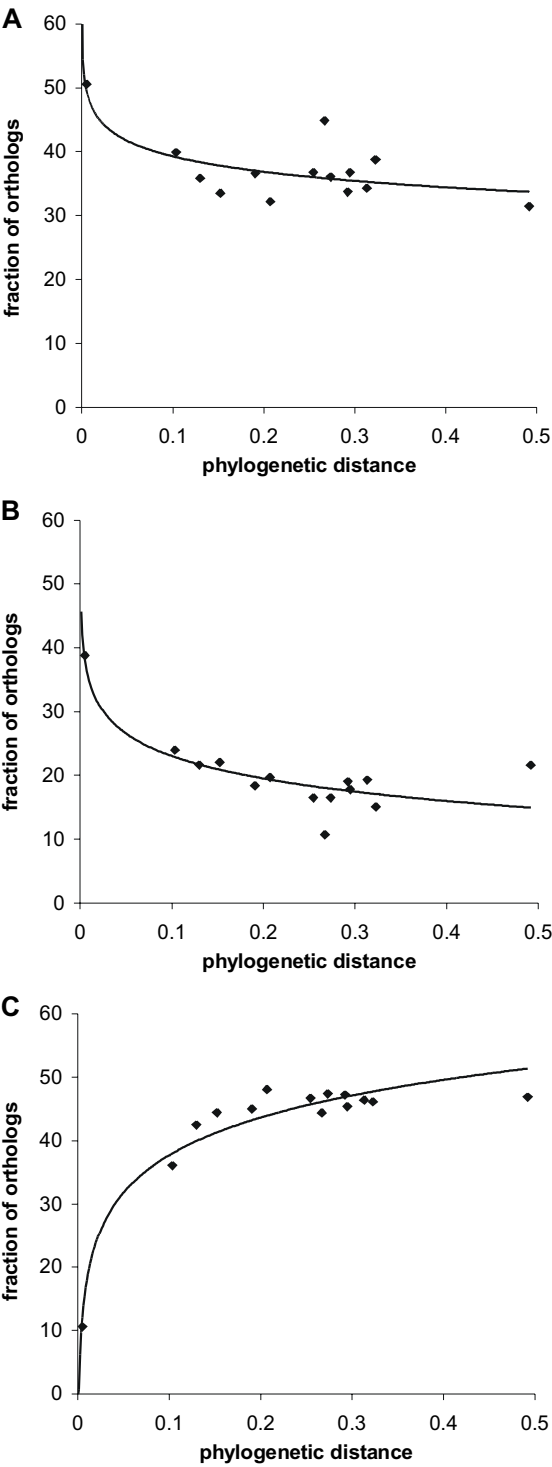
**A**



**B**



**C**



Fig. 3. Relation between the fractions of orthologs and the phylogenetic distance measured by 16S rRNA performed for three groups of orthologs: lying on the leading strand (A), lying on the lagging strand (B) and which changed DNA strand (C).

Data obtained from comparison of the *E. coli* EDL933 genome with 14 other genomes belonging to different taxonomic groups.

This observation implies also that there are some sequences which "are used to" staying on the leading DNA strand and they have lower probability of being inverted than sequences which "are used to" staying on the lagging strand. As a consequence, the set of coding sequences found on the leading strand should be not uniform. It should consist of a set of sequences which permanently or preferentially stay on the leading strand and a set of mobile sequences which are only transiently transferred from the lagging strand. To test this hypothesis we analysed the sets of 233 orthologs represented in all 15 genomes. In the first step we compared the most closely related genomes in the analysed set – two *E. coli* strains – and we counted the fractions of orthologs which stayed at the same DNA strands (leading or lagging) and the fraction of orthologs which switched their strands. In the next step we added to the comparison the third genome (the closest to the *E. coli* EDL933 genome according to the 16S rRNA phylogenetic distance) and again counted sequences which stayed at the same DNA strand in all the three genomes and sequences which switched their strand at least in one genome and so on, adding new, more distant genome to the analysed group. In Fig. 4, in the diagram, we have presented the results of analysis; values on y-axis correspond to the fraction of sequences of a given group of orthologs, while at the bottom the name of a new genome added to the comparison is shown. The fraction
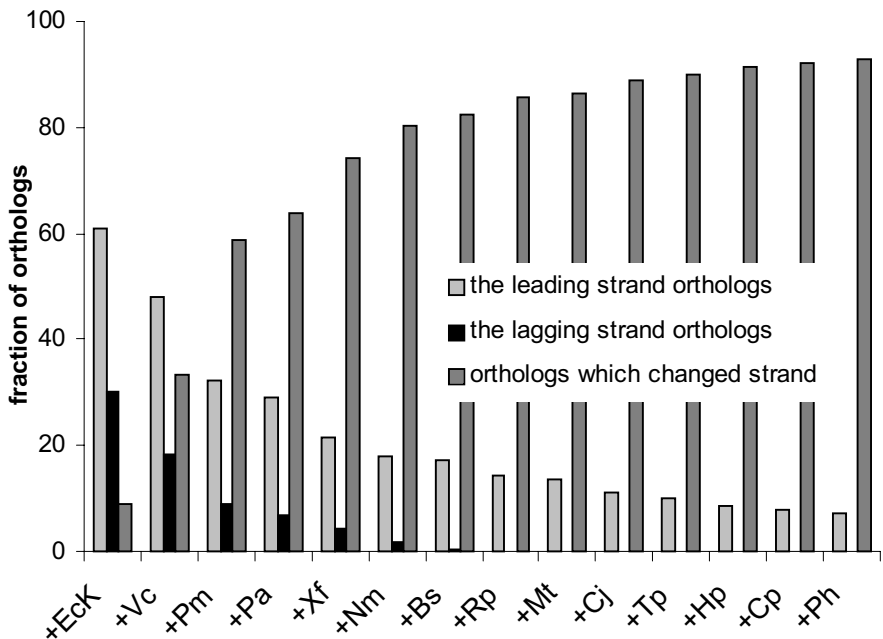


Fig. 4. The fractions of three groups of orthologs counted for the comparisons of *E. coli* EDL933 with successively added genomes to the comparison.

The group of the leading strand orthologs contains sequences which stay on the leading strand in all analysed genomes in a given comparison. Analogously for the lagging strand orthologs. The third group of orthologs includes sequences which switched their strand at least in one genome in a given comparison. Data were obtained for the sets 233 orthologs present in all 15 genomes. For genomes name abbreviations see Materials and Methods.

Table II
Orthologs found in all 15 analysed genomes
on the leading strand

| COG number | description |
|---|---|
| COG0051 | Ribosomal protein S10 |
| COG0087 | Ribosomal protein L3 |
| COG0088 | Ribosomal protein L4 |
| COG0090 | Ribosomal protein L2 |
| COG0091 | Ribosomal protein L22 |
| COG0092 | Ribosomal protein S3 |
| COG0093 | Ribosomal protein L14 |
| COG0094 | Ribosomal protein L5 |
| COG0096 | Ribosomal protein S8 |
| COG0097 | Ribosomal protein L6 |
| COG0098 | Ribosomal protein S5 |
| COG0185 | Ribosomal protein S19 |
| COG0186 | Ribosomal protein S17 |
| COG0197 | Ribosomal protein L16/L10E |
| COG0198 | Ribosomal protein L24 |
| COG0200 | Ribosomal protein L15 |
| COG0256 | Ribosomal protein L18 |

of sequences which stay in all analysed genomes on the lagging strand drops very fast and after adding the eighth genome it reaches zero, which means that there are no orthologous coding sequences located on the lagging strands in all compared genomes. For this group of compared genomes, there are still some orthologs which stay on the leading strand in all the genomes and this fraction seems to approximate asymptotically about 7% of all compared coding sequences, even after adding the most distant genome belonging to Archaea. These orthologs code for ribosomal proteins commonly considered highly conserved (Table II). The position of these genes on the leading strand seems to be conserved even across the two kingdoms (Bacteria and Archaea). It was observed that their operons are well preserved even in divergent species (W a t a n a b e  *et al.*, 1997; I t o h  *et al.*, 1999; N i k o l a i c h i k  and  D o n a c h i e, 2000;  T a m a m e s,  2001). Moreover, it was found that ribosomal genes are preferentially located in many genomes on the leading strand (M c L e a n  *et al.* 1998) probably (what is important for highly expressed genes) to avoid head-on collisions between replication and transcription complexes (B r e w e r, 1988; F r e n c h, 1992).

In the next studies we have analysed the divergence measured by the mean number of amino acid substitutions per site in groups of sequences classified according to their mobility between differently replicating DNA strands. Analyses were performed with the sets of 1521 orthologs present in all 7 genomes belonging to γ-Proteobacteria. We compared the *E. coli* EDL933 genome with six other genomes.
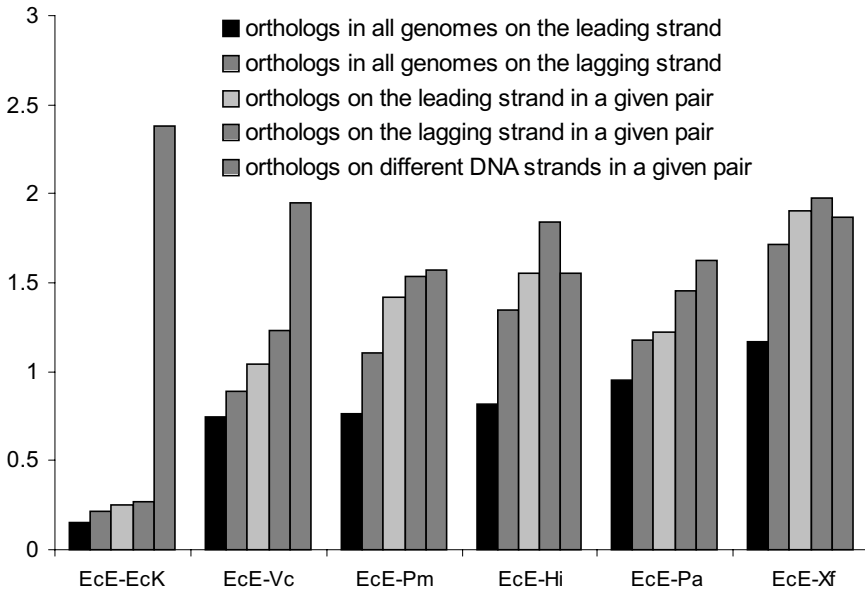
Fig. 5. The divergence measured by the mean number of amino acid substitutions per site according to WAG model (Whelan and Goldman, 2001) in five groups of sequences classified according to their mobility between differently replicating DNA strands for comparisons of *E. coli* EDL933 with other genomes. Analyses were performed with the sets of 1521 orthologs present in all 7 genomes belonging to γ-Proteobacteria. For genomes name abbreviations see Materials and Methods.

We divided all orthologs into five sets: 1 – genes staying in all analysed genomes on the leading strand, 2 – genes staying in all analysed genomes on the lagging strand, 3 – genes which are located in *E. coli* EDL933 and in the compared genome on the leading strand but can be found in at least one of the other genomes of γ-Proteobacteria on the lagging strand, 4 – genes which are located in *E. coli* EDL933 and in the compared genome on the lagging strand but can be found in at least one of the other genomes on the leading strand and, 5 – sequences which are located on different DNA strands in the compared genomes. The divergence values between genes of the *E. coli* EDL933 genome and other genomes of γ-Proteobacteria are shown in Fig. 5. We have found that there are statistically significant differences in the relative divergence between genes classified according to their position and mobility. The differences between set 1 and set 5 are statistically significant (with $p < 0.01$) for all comparisons. It is clear that the divergence of the orthologs which switched strand (set 5) is especially high for the closest genomes, which was already reported (T i l l i e r  and  C o l l i n s, 2000c; S z c z e p a n i k  *et al.*, 2001; R o c h a and  D a n c h i n,  2001) and decreases for pairs of distant genomes. Differences in divergence between set 5 and all other sets are statistically significant (with $p < 0.01$) for pairs: EcE-EcK, EcE-Vc and EcE-Pa.

In all compared pairs of genomes the lowest divergence is observed for the orthologs which permanently stay at the leading strand and do not change their strand

even at long evolutionary distances (set 1). If we eliminate this set of conserved genes from the set of all orthologs found on the leading strand (receiving set 3), the rest still seems to be less prone to accumulate substitutions than the genes from the lagging strand. However, we have found only one statistically significant difference in divergence (5.6% of all comparisons) when we compared sets 2, 3 and 4 with each other for all pairs of genomes. Furthermore, the divergence in these three sets is significantly different (with $p < 0.01$) when analysed by the ANOVA Kruskal-Wallis test only for one pair EcE-Pa. It indicates that these three sets form rather uniform group.

## Conclusions

The observed rearrangements in bacterial chromosomes are not random. Mutational pressure, responsible for the observed asymmetry in DNA composition, affects especially the copies of genes translocated to other DNA strand. According to the mobility (frequency of translocations between leading and lagging strand) it is possible to classify genes into two groups: highly conserved genes permanently or preferentially lying on the leading strand and genes switching their position between the leading and lagging DNA strands.

## Literature

A l o k a m  S., S.L. L i u, K. S a i d  and K.E. S a n d e r s o n. 2002. Inversions over the terminus region in *Salmonella* and *Escherichia coli*: IS200s as the sites of homologous recombination inverting the chromosome of *Salmonella enterica* serovar typhi. *J. Bacteriol.* **184**: 6190–6197.

B e l l g a r d  M.I., T. I t o h, H. W a t a n a b e, T. I m a n i s h i  and T. G o j o b o r i. 1999. Dynamic evolution of genomes and the concept of genome space. *Ann. N. Y. Acad. Sci.* **870**: 293–300.

B r e w e r  B.J. 1988. When polymerases collide: replication and the transcriptional organization of the *E. coli* chromosome. *Cell* **53**: 679–686.

E i s e n  J.A., J.F. H e i d e l b e r g, O. W h i t e  and S.L. S a l z b e r g. 2000. Evidence for symmetric chromosomal inversions around the replication origin in bacteria. *Genome Biol.* **1**(6): research0011.

F i t c h  W.M. 1970. Distinguishing homologous from analogous proteins. *Sys Zool.* **19**: 99–113.

F r a n c i n o  M.P. and H. O c h m a n. 1997. Strand asymmetries in DNA evolution. *Trends Genet.* **13**: 240–245.

F r a n c o i s  V., J. L o u a r n, J. P a t t e, J.E. R e b o l l o  and J.M. L o u a r n. 1990. Constraints in chromosomal inversions in *Escherichia coli* are not explained by replication pausing at inverted terminator-like sequences. *Mol. Microbiol.* **4**: 537–542.

F r a n k  A.C. and J.R. L o b r y. 1999. Asymmetric substitution patterns: a review of possible underlying mutational or selective mechanisms. *Gene* **238**: 65–77.

F r a s e r  C.M., J.D. G o c a y n e, O. W h i t e, M.D. A d a m s, R.A. C l a y t o n, R.D. F l e i s c h m a n n, C.J. B u l t, A.R. K e r l a v a g e, G.G. S u t t o n, J.M. K e l l e y, *et al.* 1995. The minimal gene complement of *Mycoplasma genitalium*. *Science* **270**: 397–403.

F r e e m a n  J.M., T.N. P l a s t e r e r, T.F. S m i t h  and S.C. M o h r. 1998. Patterns of genome organization in bacteria. *Science* **279**: 1827.

F r e n c h  S. 1992. Consequences of replication fork movement through transcription units *in vivo*. *Science* **258**: 1362–1365.

G r i g o r i e v  A. 1998. Analyzing genomes with cumulative skew diagrams. *Nucleic Acids Res.* **26**: 2286–2290.

H u g h e s  D. 2000. Evaluating genome dynamics: the constraints on rearrangements within bacterial genomes. *Genome Biol.* **1**(6):reviews0006.

H u y n e n  M.A. and E. v a n  N i m w e g e n. 1998. The frequency distribution of gene family sizes in complete genomes. *Mol. Biol. Evol.* **15**: 583–589.

I t o h  T., K. T a k e m o t o, H. M o r i  and T. G o j o b o r i. 1999. Evolutionary Instability of Operon Structures Disclosed by Sequence Comparisons of Complete Microbial Genomes. *Mol. Biol. Evol.* **16**: 332–346.

K o l s t o  A.B. 1997. Dynamic bacterial genome organization. *Mol. Microbiol.* **24**: 241–248.

K o o n i n  E.V., R.L. T a t u s o v  and M.Y. G a l p e r i n. 1998. Beyond complete genomes: from sequence to structure and function. *Curr. Opin. Struct. Biol.* **8**: 355–363.

K o w a l c z u k  M., P. M a c k i e w i c z, D. M a c k i e w i c z, A. N o w i c k a, M. D u d k i e w i c z, M.R. D u d e k  and S. C e b r a t. 2001a. DNA asymmetry and the replicational mutational pressure. *J. Appl. Genet.* **42**: 553–577.

K o w a l c z u k  M., P. M a c k i e w i c z, D. M a c k i e w i c z, A. N o w i c k a, M. D u d k i e w i c z, M.R. D u d e k  and S. C e b r a t. 2001b. High correlation between the turnover of nucleotides under mutational pressure and the DNA composition. *BMC Evol. Biol.* **1**(1):13.

K u m a r  S., K. T a m u r a  and M. N e i. 1993. MEGA: Molecular Evolutionary Genetics Analysis. Pennsylvania State University, University Park, PA.

K u n s t  F., N. O g a s a w a r a, I. M o s z e r, A.M. A l b e r t i n i, G. A l l o n i, V. A z e v e d o, M.G. B e r t e r o, P. B e s s i e r e s, A. B o l o t i n, S. B o r c h e r t, *et al.* 1997. The complete genome sequence of the gram-positive bacterium *Bacillus subtilis*. *Nature* **390**: 249–256.

L a f a y  B., A.T. L l o y d, M.J. M c L e a n, K.M. D e v i n e, P.M. S h a r p  and K.H. W o l f e. 1999. Proteome composition and codon usage in spirochaetes: species-specific and DNA strand-specific mutational biases. *Nucleic Acids Res.* **27**: 1642–1649.

L i u  S.L. and K.E. S a n d e r s o n. 1995. Rearrangements in the genome of the bacterium *Salmonella typhi*. *Proc. Natl. Acad. Sci. USA* **92**: 1018–1022.

L i u  S.L. and K.E. S a n d e r s o n. 1996. Highly plastic chromosomal organization in *Salmonella typhi*. *Proc. Natl. Acad. Sci. USA* **93**: 10303–10308.

L o b r y  J.R. 1996. Asymmetric substitution patterns in the two DNA strands of bacteria. *Mol. Biol. Evol.* **13**: 660–665.

M a c k i e w i c z  P., A. G i e r l i k, M. K o w a l c z u k, M.R. D u d e k  and S. C e b r a t. 1999a. Asymmetry of nucleotide composition of prokaryotic chromosomes. *J. Appl. Genet.* **40**: 1–14.

M a c k i e w i c z  P., A. G i e r l i k, M. K o w a l c z u k, M.R. D u d e k  and S. C e b r a t. 1999b. How does replication-associated mutational pressure influence amino acid composition of proteins? *Genome Res.* **9**: 409–416.

M a c k i e w i c z  P., D. S z c z e p a n i k, A. G i e r l i k, M. K o w a l c z u k, A. N o w i c k a, M. D u d k i e w i c z, M.R. D u d e k  and S. C e b r a t. 2001a. The differential killing of genes by inversions in prokaryotic genomes. *J. Mol. Evol.* **53**: 615–621.

M a c k i e w i c z  P., D. S z c z e p a n i k, M. K o w a l c z u k  and S. C e b r a t. 2001b. Flip-flop around the origin and terminus of replication in prokaryotic genomes. *Genome Biol.* **2**(12): inter-actions1004.

M a h a n  M.J. and J.R. R o t h. 1988. Reciprocality of recombination events that rearrange the chromosome. *Genetics* **120**: 23–35.

M a h a n  M.J. and J.R. R o t h. 1991. Ability of a bacterial chromosome segment to invert is dictated by included material rather than flanking sequence. *Genetics* **129**: 1021–1032.

M c I n e r n e y  J.O. 1998. Replicational and transcriptional selection on codon usage in *Borrelia burgdorferi*. *Proc. Natl. Acad. Sci. USA* **95**: 10698–10703.

M c L e a n  M.J., K.H. W o l f e  and K.M. D e v i n e. 1998. Base composition skews, replication orientation, and gene orientation in 12 prokaryote genomes. *J. Mol. Evol.* **47**: 691–696.

M r a z e k  J. and  S.  K a r l i n.  1998. Strand compositional asymmetry in bacterial and large viral genomes. *Proc. Natl. Acad. Sci. USA* **95**: 3720–3725.

M u s h e g i a n  A.R. and E.V.  K o o n i n.  1996. Gene order is not conserved in bacterial evolution. *Trends Genet.* **12**: 289–290.

N i k o l a i c h i k  Y.A. and W.D.  D o n a c h i e.  2000. Conservation of gene order amongst cell wall and cell division genes in Eubacteria, and ribosomal genes in Eubacteria and Eukaryotic organelles. *Genetica* **108**: 1–7.

O h n o  S.  1970. Evolution by gene duplication. George Allen and Unwin, London.

Q i a n  J., N.M.  L u s c o m b e  and M.  G e r s t e i n.  2001. Protein family and fold occurrence in genomes: power-law behaviour and evolutionary model. *J. Mol. Biol.* **313**: 673–681.

R e a d  T.D., R.C.  B r u n h a m,  C.  S h e n,  S.R.  G i l l,  J.F.  H e i d e l b e r g,  O.  W h i t e,  E.K. H i c k e y,  J.  P e t e r s o n,  T.  U t t e r b a c k,  K.  B e r r y,  *et al.* 2000. Genome sequences of *Chlamydia trachomatis* MoPn and *Chlamydia pneumoniae* AR39. *Nucleic Acids Res.* **28**: 1397–1406.

R e b o l l o  J.E., V.  F r a n c o i s  and J.M.  L o u a r n.  1988. Detection and possible role of two large nondivisible zones on the *Escherichia coli* chromosome. *Proc. Natl. Acad. Sci. USA* **85**: 9391–9395.

R o c h a  E.P., A.  D a n c h i n  and A.  V i a r i.  1999. Universal replication biases in bacteria. *Mol. Microbiol.* **32**: 11–16.

R o c h a  E.P. and A.  D a n c h i n.  2001. Ongoing evolution of strand composition in bacterial genomes. *Mol. Biol. Evol.* **18**: 1789–1799.

S a n d e r s o n  K.E. and S.L.  L i u.  1998. Chromosomal rearrangements in enteric bacteria. *Electrophoresis* **19**: 569–572.

S c h m i d  M.B., J.R.  R o t h.  1983. Selection and endpoint distribution of bacterial inversion mutations. *Genetics* **105**: 539–55.7

S c h m i d t  H.A., K.  S t r i m m e r,  M.  V i n g r o n  and A.  v o n  H a e s e l e r.  2002. TREE-PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics* **18**: 502–504.

S e g a l l  A., M.J.  M a h a n  and J.R.  R o t h.  1988. Rearrangement of the bacterial chromosome: forbidden inversions. *Science* **241**: 1314–1318.

S e g a l l  A.M. and J.R.  R o t h.  1989. Recombination between homologies in direct and inverse orientation in the chromosome of Salmonella: intervals which are nonpermissive for inversion formation. *Genetics* **122**: 737–747.

S ł o n i m s k i  P.P., M.O.  M o s s e,  P.  G o l i k,  A.  H e n a u t,  Y.  D i a z,  J.L.  R i s l e r,  J.P.  C o m e t, J.C.  A u d e,  A.  W o ź n i a k,  E.  G l e m e t,  *et al.* 1998. The first laws of genomics. *Microb. Comp. Genomics* **3**: 46.

S o k a l  R.R. and F.J.  R o h l f.  1995. Biometry. Freeman, New York.

S u y a m a  M. and P.  B o r k.  2001. Evolution of prokaryotic gene order: genome rearrangements in closely related species. *Trends Genet.* **17**: 10–13.

S z c z e p a n i k  D., P.  M a c k i e w i c z,  M.  K o w a l c z u k,  A.  G i e r l i k,  A.  N o w i c k a,  M.R. D u d e k  and S.  C e b r a t.  2001. Evolution rates of genes on leading and lagging DNA strands. *J. Mol. Evol.* **52**: 426–433.

T a m a m e s  J.  2001. Evolution of gene order conservation in prokaryotes. *Genome Biol.* **2** (6): research 0020.

T a m u r a  K. and M.  N e i.  1993. Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Mol. Biol. Evol.* **10**: 512–526.

T a t u s o v  R.L., A.R.  M u s h e g i a n,  P.  B o r k,  N.P.  B r o w n,  W.S.  H a y e s,  M.  B o r o d o v s k y, K.E.  R u d d  and E.V.  K o o n i n.  1996. Metabolism and evolution of *Haemophilus influenzae* deduced from a whole-genome comparison with *Escherichia coli*. *Curr. Biol.* **6**: 279–291.

T a t u s o v  R.L., D.A.  N a t a l e,  I.V.  G a r k a v t s e v,  T.A.  T a t u s o v a,  U.T.  S h a n k a v a r a m, B.S.  R a o,  B.  K i r y u t i n,  M.Y.  G a l p e r i n,  N.D.  F e d o r o v a  and E.V.  K o o n i n.  2001. The COG database: new developments in phylogenetic classification of proteins from complete genomes. *Nucleic Acids Res.* **29**: 22–28.

T h o m p s o n  J.D., D.G. H i g g i n s  and T.J. G i b s o n. 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**: 4673–4680.

T i l l i e r  E.R. and R.A. C o l l i n s. 2000a. The contributions of replication orientation, gene direction, and signal sequences to base-composition asymmetries in bacterial genomes. *J. Mol. Evol.* **50**: 249–257.

T i l l i e r  E.R. and R.A. C o l l i n s. 2000b. Genome rearrangement by replication-directed translocation. *Nat. Genet.* **26**: 195–197.

T i l l i e r  E.R. and R.A. C o l l i n s. 2000c. Replication orientation affects the rate and direction of bacterial gene evolution. *J. Mol. Evol.* **51**: 459–463.

W a t a n a b e  H., H. M o r i, T. I t o h  and T. G o j o b o r i. 1997. Genome plasticity as a paradigm of eubacterial evolution. *J. Mol. Evol.* **44** (Suppl. 1):S57–S64.

W h e l a n  S. and N. G o l d m a n. 2001. A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach. *Mol. Biol. Evol.* **18**: 691–699.