1 **Title**

2 **Polymorphism at four-fold degenerate site, but not at the intergenic regions, explains nucleotide**

3 **compositional strand asymmetry in bacteria**

4 Piyali Sen[a], Ruksana Aziz[b], Soumita Das[a], Nima Dondu Namsa[b], Ramesh Chandra Deka[c], Edward J Feil[d*],

5 Suvendra Kumar Ray[b*], Siddhartha Sankar Satapathy[a*]

6 [a]Department of Computer Science and Engineering, [b]Department of Molecular Biology and Biotechnology,

7 [c]Department of Chemical Sciences, Tezpur University, Napaam, Tezpur-784028, Assam, India

8 [d]The Milner Centre for Evolution, Department of Biology and Biochemistry, University of Bath, Bath BA2

9 7AY, UK

10 *Authors for Correspondence:   E. J. Feil: E.Feil@bath.ac.uk
11                                S. K. Ray: suven@tezu.ernet.in
12                                S. S. Satapathy: ssankar@tezu.ernet.in

13

14

15 **Running Title:** Base compositional strand asymmetry in genomes

16

17

18

19

20

21

22

23

24

25

26

27

28

1

**Abstract**

We investigated single nucleotide polymorphism in intergenic regions (IRs) and four-fold degenerate sites (FFS) in genomes of three γ-Proteobacteria and two Firmicutes to understand the mechanism of nucleotide compositional asymmetry between the leading and the lagging strands. Pattern of the polymorphism spectra were alike regarding transitions but variable regarding transversions in the IRs of these bacteria. Contrasting trends of complementary polymorphisms such as C→T vs G→A as well as A→G vs T→C in the IRs vindicated similar replication-associated strand asymmetry regarding cytosine and adenine deamination, respectively, across these bacteria. Surprisingly, the polymorphism pattern at FFS was different from that of the IRs and its frequency was always more than the IRs in these bacteria. Further, the polymorphism patterns within a bacterium were inconsistent across the five amino acids, which neither the replication nor the transcription-associated mutations could explain. However, the polymorphism at FFS coincided with amino acid specific codon usage bias in the five bacteria. Further, strand asymmetry in nucleotide composition could be explained by the polymorphism at FFS, not at the IRs. Therefore, polymorphisms at FFS might not be treated as nearly neutral unlike that in IRs in these bacteria.

**Key words**: nucleotide polymorphism, compositional strand asymmetry, intergenic regions, four-fold degenerate site, codon usage bias

**Introduction**

As each of the four DNA bases can mutate to any of the three other bases, there are twelve possible directional substitution mutation types that include four transitions and eight transversions. These directional substitution mutations do not occur at equal frequencies in bacterial genomes for mechanistic reasons such as unequal stability among different base pairs, differential propensity of bases to damages such as deamination, oxidation, and radiation as well as selective reasons such as differential impact on the structure and function of DNA, RNA, and proteins. There are twice the number of transversions than transitions, but the observed frequency of transitions are double the transversions (Seplyarskiy et al. 2012; Duchêne et al. 2015; Lewis et al. 2016; Stoltzfus and Norris 2016; Lyons and Lauring 2017; Schroeder et al. 2017). Further, the higher propensity of deamination of C and oxidation of G bases increase the frequency of C→T and G→T base

2

57   substitutions (Suzuki and Kamiya 2017), explaining the universal mutation bias towards A/T in genomes (Lind

58   and Andersson 2008; Balbi et al. 2009; Hershberg and Petrov 2010; Hildebrand et al. 2010; Rocha and Feil

59   2010; Van Leuven and McCutcheon 2012).

60       Transition and transversion frequencies also vary between the leading strand (LeS) and the lagging

61   strand (LaS) of replication known as strand substitution asymmetry (Lobry 1996; Frank and Lobry 1999),

62   which is explained with GC (or AT) skew in chromosomes (Grigoriev 1998; Karlin et al. 1998; Rocha et al.

63   1999; Lobry and Sueoka 2002; Lobry 1996). Rapid deamination of cytosine to uracil in single stranded DNA

64   leads to a higher rate of C→T mutation favouring a higher frequency of G and T in the LeS compared to the

65   LaS (Rocha et al. 2006; Francino and Ochman 1997). The resulting GC (or AT) skew violates Chargaff's

66   second parity rule which states that the frequency of G and C (or A and T) should be approximately equal

67   within individual DNA strands in genomes (Sueoka 1995; Forsdyke and Mortimer 2000; Powdel et al. 2009).

68   Firmicutes are exceptional as they exhibit higher frequencies of A than T in the LeS (Rocha et al. 2006)

69   implicating the selection and strong gene-orientation biases in the genomes of this bacterial group (Charneski

70   et al. 2011). Assuming four-fold degenerate sites (FFS) evolve under near neutrality, Rocha et al. (2006)

71   determined the polymorphism at four-fold degenerate sites (FFS) in seven different bacterial species to explain

72   patterns of GC skew in genomes. They noted that relative consistency between taxa in terms of base

73   compositional biases does not correspond with the underlying base substitution profiles. Although

74   transcription-associated mutation was known to occur, the emphasis was given towards replication-associated

75   mutation as highly expressed genes were not considered for analysis in their study (Davis 1989; Mugal et al.

76   2009; Kim and Jinks-Robertson 2012; Gaillard et al. 2013; Jinks-Robertson and Bhagwat 2014). Considering

77   FFS to be nearly neutral, several researchers have extensively used FFS as a reference to describe strand as

78   well as genome composition in bacteria (Muto and Osawa 1987; Rocha et al. 2006; Hershberg and Petrov

79   2010; Hildebrand et al. 2010), thereby neglecting the possible contribution of a selection mechanism. No cross-

80   verification studies have been carried out in support of the above assumption of near neutrality of FFS being

81   true. In case of nearly neutral nature, polymorphism patterns of different amino acids at FFS and at IRs are

82   expected to be similar considering the role of context dependent mutation being minimal. The assumption of

83   nearly neutral nature of FFS could be questioned based on the findings of earlier work in the areas of co-

84   translational protein folding, selection on synonymous codons and ribosome mediated gene regulation

85   (Satapathy et al. 2014; Sohmen et al. 2015; Ray and Goswami 2016; Ito 2016). Moreover, the findings by

3

86    Charneski et al. (2011) support the role of selection for the atypical AT skew in *S. aureus*, emphasising the

87    role of strand bias gene distribution in compositional skew (Charneski et al. 2011). A similar observation

88    suggesting the selection at the IRs has been reported recently in bacteria (Thorpe et al. 2017). Considering

89    mutation being AT biased in genomes (Hershberg and Petrov 2010; Hildebrand et al. 2010), it was previously

90    assumed that the entire bacterial genome is under selection (Rocha and Feil 2010; Raghavan et al. 2012). The

91    above findings have motivated us to relook at the polymorphism at FFS to determine the strand compositional

92    asymmetry.

93         We have investigated nucleotide polymorphisms by analysing many strains of *Escherichia coli* (*Ec*),

94    *Klebsiella pneumoniae* (*Kp*), *Salmonella enterica* (*Se*) belonging to γ-Proteobacteria and two members of

95    Firmicutes such as *Staphylococcus aureus* (*Sa*) and *Streptococcus pneumoniae* (*Sp*). The Firmicutes are known

96    to exhibit different nucleotide compositional asymmetry between strands as compared to γ-Proteobacteria,

97    which has been ascribed to replication-associated mutations due to the presence of two isoforms of DNA

98    polymerase III alpha subunit, PolC and DnaE in Firmicutes (Rocha 2004; Saha et al. 2014). However, a large-

99    scale analysis of diverse bacterial genomes indicates that purine asymmetry across two strands of replication,

100   and different DNA polymerase compositions are neither essential nor exclusive features of the Firmicutes(Saha

101   et al. 2014). It indicates a partial contribution of mutation in strand asymmetry composition in Firmicutes.

102   Considering IRs are nearly neutral regions in a genome, polymorphism at IRs can be attributed to the

103   replication-associated mutations (RAM) which are reflected by the contrasting pattern of complementary

104   polymorphisms (e.g., C→T vs G→A; T→G vs A→C etc.) between the LeS and the LaS. Whereas the

105   polymorphism at FFS can be attributed to RAM as well as transcription-associated mutations (TAM). In our

106   comparative study, while the polymorphism at the IRs could be explained by RAM, polymorphism at FFS

107   could not be explained by RAM or TAM in these bacteria, which was surprising.  Instead, the polymorphism

108   at FFS was observed to be influenced by the amino acid specific codon usage bias. Our work indicates that the

109   polymorphism at FFS, not at IRs, can explain the nucleotide compositional strand asymmetry in bacterial

110   genomes.

111

112

4

113 **Results**

114 **Strand asymmetry in intergenic regions of bacteria due to transition polymorphisms**

115        Prior to the analysis of nucleotide polymorphism at intergenic regions (IRs), we performed an

116 elaborate study on its nucleotide composition. G+C% difference between whole genome and IRs was more

117 prominent for the three γ-Proteobacteria (*Ec, Kp, Se*) than the two Firmicutes (*Sa*, *Sp*) (Supplementary Table

118 1). The G+C% in the IRs was similar between the LeS and the LaS within a bacterium. Abundance values

119 between complementary nucleotides were more similar than that between two non-complementary nucleotides

120 (Table 1). We found out different skews (AT/GC/KM/RY) in the IRs. AT skew was negative in the LeS and

121 positive in the LaS in all bacteria except *Sa* where the reverse AT skew pattern was observed. The atypical AT

122 skew in IR was in concordance with AT skew in *Sa* chromosome (Charneski et al. 2011). It is pertinent to note

123 that the AT skew patterns were not consistent in IRs of *Sa* and *Sp*, although both belong to Firmicutes. In

124 contrast to AT skew, GC skew values were found to be positive in the LeS and negative in the LaS of these

125 bacteria. The magnitude of GC skew was observed to be higher than that of AT skew across the five bacteria.

126 Among the five bacteria, the magnitude of both AT and GC skew was observed relatively high in *Sa*. The keto-

127 amino (KM) and purine-pyrimidine (RY) skews were in general positive in the LeS and negative in the LaS

128 of these bacteria (Table 1).

129 **Table 1: Compositional features of IRs in LeS and LaS in five bacteria**

| Nucleotide compositional features | *Ec* | | *Kp* | | *Se* | | *Sa* | | *Sp* | |
|---|---|---|---|---|---|---|---|---|---|---|
| | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS |
| A | 35099 | 37981 | 44856 | 25471 | 37686 | 39691 | 52274 | 48613 | 25008 | 21804 |
| T | 36538 | 36517 | 46257 | 23788 | 39566 | 37653 | 48314 | 52404 | 25032 | 21787 |
| G | 26040 | 23889 | 40820 | 19556 | 29022 | 26244 | 21509 | 17025 | 13852 | 9694 |
| C | 23398 | 25521 | 39052 | 21513 | 26341 | 28585 | 16212 | 21998 | 11371 | 11450 |
| AT Skew | -0.020 | 0.020 | -0.015 | 0.034 | -0.024 | 0.026 | 0.039 | -0.038 | -0.001 | 0.000 |
| GC Skew | 0.053 | -0.033 | 0.022 | -0.048 | 0.048 | -0.043 | 0.140 | -0.127 | 0.098 | -0.083 |
| KM Skew | 0.034 | -0.025 | 0.019 | -0.040 | 0.034 | -0.033 | 0.010 | -0.008 | 0.033 | -0.027 |
| RY Skew | 0.010 | -0.001 | 0.002 | -0.003 | 0.006 | -0.002 | 0.067 | -0.063 | 0.033 | -0.027 |

130 Table presents count of the nucleotides at IRs and nucleotide composition skews in LeS of IRs as well as LaS of IRs of
131 five bacteria of LeS and LaS of the five bacteria. AT Skew is defined as (A − T)/(A+T), GC skew is defined as (G −
132 C)/(G+C), KM skew is defined as [(G+T) − (A+C)]/[(G+T) + (A+C)], RY skew is defined as [(A+G) − (C+T)]/[(A+G)
133 + (C+T)]. Here the A, T, G, C represents respective nucleotide counts.
134

135        The disparity of nucleotide composition between the LeS and the LaS was analysed in the context of

136 nucleotide polymorphisms at the IRs. Frequencies of the twelve nucleotide polymorphisms were found out in

137  the LeS and the LaS in these bacteria (Table 2). In general, transitions were more frequent than transversions.

138  The *ti/tv* values varied from 1.5 to 2.3 among these bacteria (Table 2). In a more critical analysis of the

139  transition and transversion frequencies, for example C→T vs. C→G in LeS, it was observed that the former

140  was more than eight times than the latter in *Ec*, while the same was even ten time in *Se* and *Sp*. This fold

141  differences were considerably higher than the expected value four-fold. The frequency values of

142  complementary transition polymorphisms exhibited contrasting trends between the strands: C→T frequency

143  was more than that of G→A in the LeS while the reverse was the case in the LaS; A→G frequency was more

144  than that of T→C in LeS while the reverse was the case in the LaS (Figure 1). Frequency values of C→T and

145  G→A were about two times more than that of A→G and T→C in both the strands (*p*-value < 0.01). This

146  indicated that cytosine deamination is a major cause for the high frequency of C→T and G→A in genomes.

147  Further, the higher frequency of C→T than that of G→A in LeS can be attributed to the higher propensity of

148  single stranded DNA towards cytosine deamination over the double stranded DNA. Similarly, higher

149  deamination of adenine in the single stranded DNA than the double stranded DNA might attribute towards the

150  higher frequency of A→G than T→C in LeS. The difference between the frequencies of C→T and G→A

151  within a strand was significantly more than that between A→G and T→C (*p*-value < 0.01). This finding is in

152  concordance with both the following notions that cytosine is more prone to deamination than adenine in DNA

153  and single stranded DNA has greater impact on cytosine deamination over adenine deamination. The strand

154  asymmetry regarding transition polymorphisms observed in this study vindicates replication associated

155  asymmetry in cytosine and adenine deamination between the LeS and LaS.  As the pattern of transition

156  polymorphism was similar across these bacteria, regarding the strand asymmetry, there seems to be similar

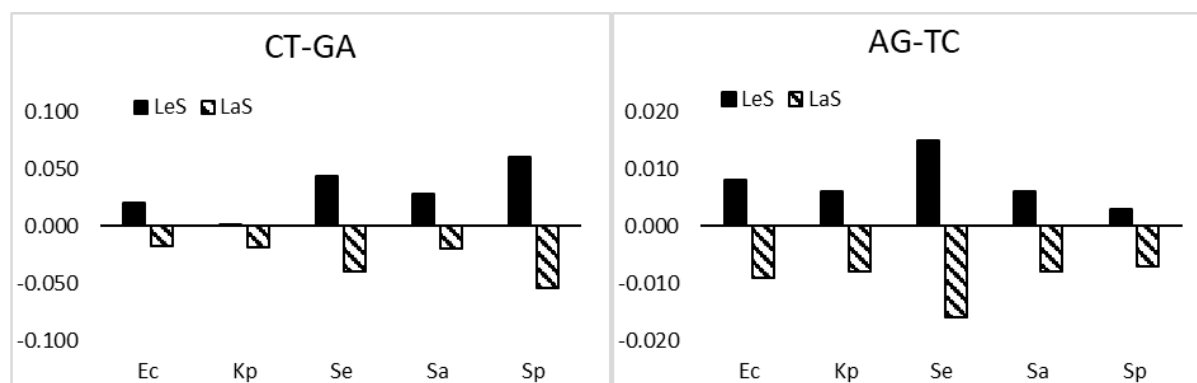157  impact of RAM in these two groups of bacteria.

158

159

160

161

162

163

164

6

165    **Table 2: Polymorphism spectra at IRs**

| Substitution spectra and features | *Ec* | | *Kp* | | *Se* | | *Sa* | | *Sp* | |
|---|---|---|---|---|---|---|---|---|---|---|
| | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS |
| A→T | 0.023 | 0.024 | 0.022 | 0.023 | 0.024 | 0.026 | 0.036 | 0.035 | 0.014 | 0.014 |
| A→G | 0.053 | 0.045 | 0.048 | 0.041 | 0.085 | 0.079 | 0.048 | 0.037 | 0.049 | 0.044 |
| A→C | 0.015 | 0.017 | 0.014 | 0.016 | 0.02 | 0.027 | 0.012 | 0.014 | 0.015 | 0.016 |
| T→A | 0.024 | 0.024 | 0.022 | 0.023 | 0.022 | 0.026 | 0.038 | 0.033 | 0.016 | 0.014 |
| T→G | 0.018 | 0.014 | 0.015 | 0.013 | 0.024 | 0.023 | 0.016 | 0.011 | 0.015 | 0.017 |
| T→C | 0.045 | 0.054 | 0.042 | 0.048 | 0.071 | 0.095 | 0.042 | 0.044 | 0.046 | 0.05 |
| G→A | 0.091 | 0.114 | 0.093 | 0.105 | 0.175 | 0.228 | 0.122 | 0.134 | 0.126 | 0.181 |
| G→T | 0.036 | 0.033 | 0.04 | 0.036 | 0.057 | 0.065 | 0.051 | 0.051 | 0.043 | 0.051 |
| G→C | 0.012 | 0.015 | 0.017 | 0.017 | 0.018 | 0.021 | 0.016 | 0.017 | 0.013 | 0.015 |
| C→A | 0.035 | 0.037 | 0.037 | 0.039 | 0.057 | 0.067 | 0.058 | 0.045 | 0.046 | 0.045 |
| C→T | 0.112 | 0.098 | 0.095 | 0.088 | 0.218 | 0.188 | 0.15 | 0.115 | 0.187 | 0.127 |
| C→G | 0.014 | 0.013 | 0.015 | 0.015 | 0.02 | 0.02 | 0.02 | 0.014 | 0.014 | 0.015 |
| ti/tv | 1.701 | 1.757 | 1.527 | 1.549 | 2.269 | 2.145 | 1.466 | 1.500 | 2.318 | 2.150 |
| AT bias (ti) | 2.071 | 2.141 | 2.089 | 2.169 | 2.519 | 2.391 | 3.022 | 3.074 | 3.295 | 3.277 |
| AT bias (tv) | 2.152 | 2.258 | 2.655 | 2.586 | 2.591 | 2.640 | 3.893 | 3.840 | 2.967 | 2.909 |

166    Table represents 12 different substitution frequencies in LeS and LaS, for example C→T represents total C→T mutations
167    divided by the total count of C in the gene. ti/tv is the ratio of transition to transversion. AT bias of transition is calculated
168    as (C→T + G→A) / (A→G + T→C). AT bias of transversion is calculated as (G→T + C→A) / (T→G + A→C).

169

170    **Figure 1: Difference between complementary transition polymorphisms in the LeS and the LaS at IRs**



172    The figure represents a two-panel histogram of difference between complementary transition polymorphisms (C→T *vs*
173    G→A and A→G *vs* T→C) in the LeS and LaS of IRs. The height of the vertical bar represents the polymorphism
174    frequency difference between a complementary polymorphism pair. Black bar and striped bar represent polymorphism
175    frequency differences in LeS and LaS, respectively. The X-axis represents the name of the five organisms: *Escherichia*
176    *coli* (*Ec*)*, Klebsiella pneumoniae* (*Kp*)*, Salmonella enterica* (*Se*)*, Staphylococcus aureus* (*Sa*) *and Streptococcus*
177    *pneumoniae* (*Sp*). Y-axis represents frequency difference values. In general, the pattern of complementary transition
178    polymorphism in LeS and LaS are found to be in opposite direction.

179

180          Among the transversion polymorphisms, G→T (C→A) was the highest across these bacteria (Table

181    2). But, the order of the other polymorphisms regarding their frequencies was variable. For example, while

182    A→T (T→A) frequency was the second highest among the transversions in *Ec*, *Kp*, *Se* and *Sa*; it was the

183     lowest in *Sp*. Though, frequencies were different among *tv* polymorphisms, the frequency of a specific *tv* was

184     similar between the strands. Therefore, the frequency values of complementary *tv* pairs were similar within a

185     strand, unlike *ti*. The most frequent transversions G→T (C→A) is known to be due to the oxidation of G to

186     form 8oxoG (Loon et al. 2010), which seems to be occurring equally in both the strands. In *Sa* and *Sp* the

187     frequency of G→T (C→A) was higher than the transitions polymorphisms A→G (T→C). This was

188     contradicting the general notion that frequency of a transition polymorphism is higher than that of a

189     transversion polymorphism. The reason behind the higher frequency of A→T (T→A) than that of A→C

190     (T→G) or C→G (G→C) is not known. In *ti* as well as in case of *tv*, polymorphisms were more than two times

191     biased towards A/T over G/C in three γ-Proteobacteria, and the same were more than three times biased

192     towards A/T over G/C in two Firmicutes. This was in concordance with their genome composition values that

193     polymorphism was more biased to AT in genome with low genome G+C composition.

194        In conclusion, polymorphism study at IRs has revealed that replication associated polymorphism

195     asymmetry between the LeS and the LaS, is mainly attributed to deamination of cytosine and adenine. This

196     polymorphism asymmetry is responsible for the positive GC or KM skew as well as the negative AT skew in

197     the IRs of LeS and the *vice versa* in the LaS. From the frequency values of these polymorphisms, it could be

198     found out that the polymorphism favours positive GC skew as well as negative AT skew in the LeS across

199     these bacteria. The atypical AT-skew in *Sa*, could not be explained by the polymorphism spectra at IRs.

200     **The polymorphism pattern at the four-fold degenerate site (FFS) is different from that at IRs**

201        Researchers had treated polymorphisms at FFS in bacteria as nearly neutral and had accordingly

202     derived their conclusion regarding strand compositional asymmetry and genome composition (Muto and

203     Osawa 1987; McLean et al. 1998; Reyes et al. 1998; Rocha and Danchin 2001; Rocha et al. 2006). However,

204     there was no comparative study of polymorphism at FFS with that at IRs. Also, the polymorphisms were not

205     compared across FFS of different amino acid codons within a bacterium. Here we found out polymorphism

206     frequencies at FFS in codons of amino acids such as Val, Pro, Thr, Ala and Gly in the γ-Proteobacteria and

207     Firmicutes (Supplementary Table 2). In general, *ti* was more frequent than *tv* across the five amino acids in

208     these bacteria. However, the *ti*/*tv* values were variable, which indicated differences across these amino acid

209     regarding polymorphism at FFS (Table 3). The transition as well as the transversion polymorphisms were more

210     biased towards A/T in the Firmicutes than the γ-Proteobacteria (*p*-value < 0.01). Further the magnitudes of

211 bias towards A/T over G/C were not consistent across the five amino acids either in the case of transition or in

212 the case of transversion polymorphisms (Table 3). It is pertinent to note that polymorphisms such as C→T and

213 G→A were not always more frequent than A→G and T→C. For example, in case of Thr in *Ec* and across the

214 five amino acids in case of *Kp* (Supplementary Table 2). Further, among the transversion polymorphisms,

215 G→T (C→A) were not the most frequent ones in *Ec* and *Kp*. So, the notion that C→T (G→A) is the most

216 frequent transition polymorphism and G→T (C→A) is the most frequent transversion polymorphism is

217 incorrect at FFS, unlike at IRs. Not only the frequency of a polymorphism was variable across the five amino

218 acids, but the order of different polymorphisms regarding their frequency at FFS was also not consistent across

219 the five amino acids (Supplementary Table 2, Figure 2). The difference among the amino acids regarding the

220 polymorphisms within a bacterium indicated that the polymorphisms at FFS not necessarily be treated as nearly

221 neutral.

222 **Table 3: Comparison between transition-transversion polymorphism at FFS of five bacteria**

| Bacteria | Polymorphism feature | Val | | Pro | | Thr | | Ala | | Gly | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS |
| *Ec* | ti/tv | 1.558 | 1.589 | 1.288 | 1.342 | 1.598 | 1.739 | 1.513 | 1.521 | 1.732 | 1.893 |
| | AT bias (ti) | 1.447 | 1.320 | 1.375 | 1.279 | 1.020 | 1.091 | 1.366 | 1.333 | 1.366 | 1.376 |
| | AT bias (tv) | 0.915 | 0.850 | 1.117 | 0.945 | 0.758 | 0.795 | 0.904 | 0.921 | 1.027 | 1.096 |
| *Kp* | ti/tv | 1.701 | 1.761 | 1.283 | 1.464 | 1.707 | 1.755 | 1.631 | 1.647 | 1.733 | 1.787 |
| | AT bias (ti) | 0.949 | 0.833 | 1.076 | 0.963 | 0.813 | 0.745 | 0.902 | 0.844 | 0.932 | 0.956 |
| | AT bias (tv) | 0.944 | 0.866 | 0.900 | 0.946 | 0.973 | 0.766 | 0.927 | 0.846 | 0.777 | 0.905 |
| *Se* | ti/tv | 2.042 | 2.185 | 1.604 | 1.624 | 1.818 | 1.982 | 1.741 | 1.846 | 2.462 | 2.541 |
| | AT bias (ti) | 2.292 | 2.165 | 2.202 | 2.242 | 1.843 | 2.037 | 2.195 | 2.098 | 1.559 | 1.508 |
| | AT bias (tv) | 1.168 | 1.078 | 1.403 | 1.278 | 1.176 | 1.161 | 1.091 | 1.065 | 1.109 | 1.220 |
| *Sa* | ti/tv | 1.383 | 1.274 | 1.383 | 1.296 | 1.088 | 1.339 | 1.359 | 1.311 | 1.476 | 1.406 |
| | AT bias (ti) | 3.362 | 2.922 | 5.130 | 5.917 | 4.041 | 4.786 | 3.928 | 4.491 | 3.077 | 2.713 |
| | AT bias (tv) | 3.800 | 3.280 | 5.122 | 7.641 | 3.770 | 4.511 | 3.943 | 4.347 | 3.427 | 3.259 |
| *Sp* | ti/tv | 1.627 | 1.669 | 1.477 | 1.472 | 1.435 | 1.417 | 1.570 | 1.398 | 1.408 | 1.484 |
| | AT bias (ti) | 2.097 | 2.225 | 3.967 | 4.051 | 2.663 | 2.852 | 2.619 | 2.817 | 2.698 | 3.168 |
| | AT bias (tv) | 1.677 | 1.563 | 3.059 | 4.677 | 2.039 | 2.604 | 2.371 | 3.325 | 2.520 | 3.064 |

223 ti/tv is the ratio of transition to transversion. AT bias of transition is calculated as (C→T + G→A) / (A→G + T→C). AT
224 bias of transversion is calculated as (G→T + C→A) / (T→G + A→C).

225

226

227

228

229

230 **Figure 2: Polymorphism frequency at FFS across five amino acids in five bacteria**
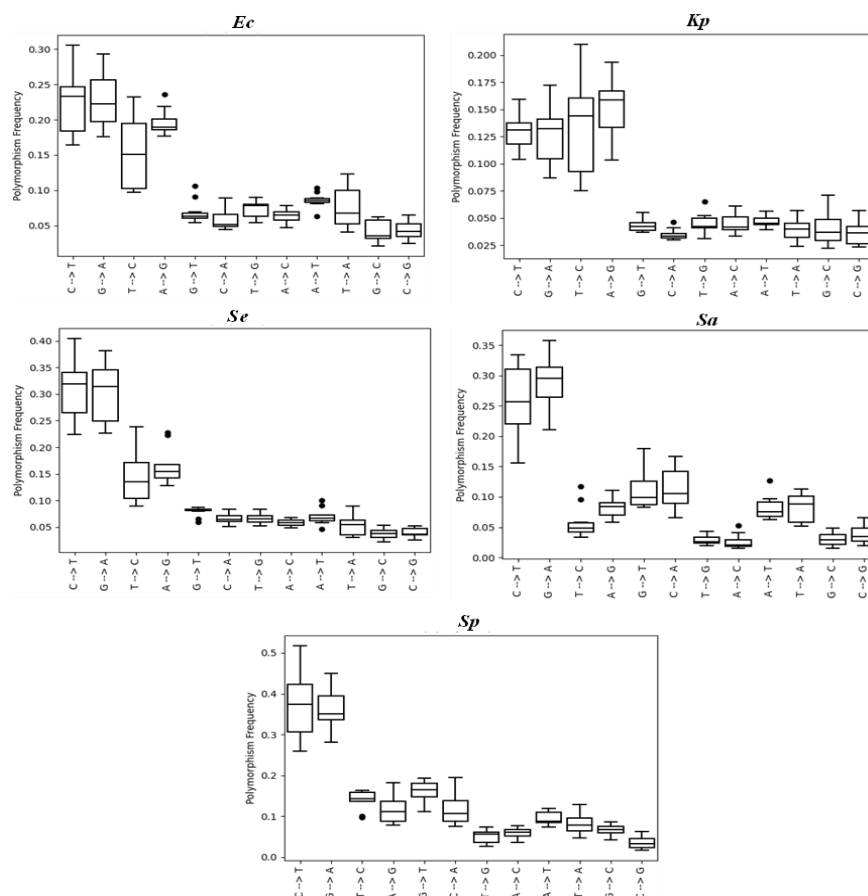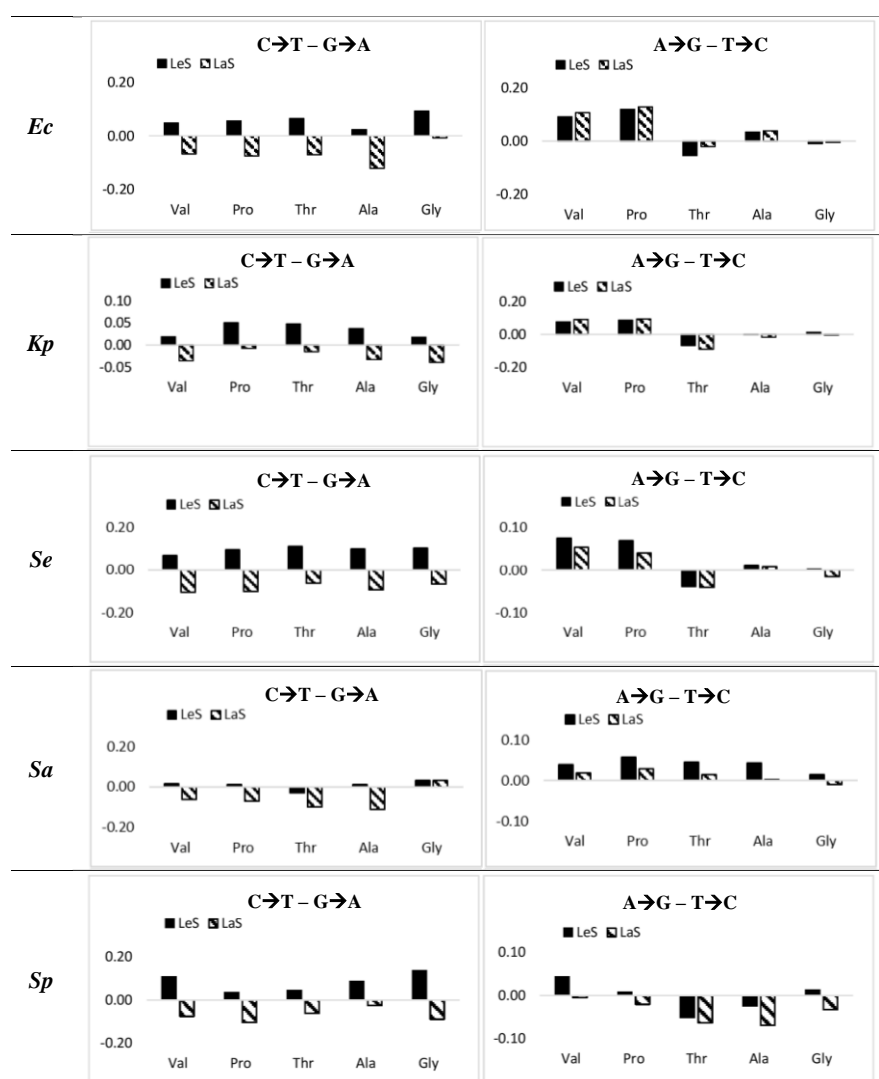


231

232 Figure presents distribution of 12 substitutions of FFS in terms of box-plot of Ec, Kp, Se, Sa and Sp. The x-axis presents
233 the 12 substitutions and the y-axis presents the polymorphism frequency. The frequency of polymorphism is not uniform
234 across the five amino acids.

235       We compared frequency values between complementary polymorphisms within a strand to find out

236 any possible role of strand asymmetry. In general, the difference between C→T and G→A were positive in

237 the LeS and negative in the LaS across the five amino acids in these five bacteria. This indicated replication

238 associated strand asymmetry regarding cytosine deamination. Regarding the other transition polymorphisms,

239 difference between A→G and T→C were high but non-uniform across the five amino acids within a bacterium:

240 the difference can be positive in case of an amino acid while negative in case of another amino acid (Figure

241 3). Further, the difference between A→G and T→C values were similar both in the LeS and the LaS. This

242 indicated that the difference was not generated due to replication associated strand asymmetry. Similarly, the

243 inconsistent pattern across the amino acids indicated that the pattern was not due to transcription associated

244 mutation. In case of complementary transversions, the difference values were high but not uniform across the

245 five amino acids within a bacterium. Further, the difference values remain similar both in the LeS and the LaS.

10

246      Unlike IRs, the difference between complementary polymorphisms were high at FFS and the difference value

247      reflected were more of amino acid specific, not strand specific.

248      **Figure 3: Difference between complementary transition polymorphisms in the LeS and the LaS at FFS**



249

250 The figure presents a two-panel histogram of difference between complementary transition polymorphisms in the LeS
251 and LaS at FFS of five bacteria. The height of the vertical bar represents the polymorphism frequency difference between
252 a complementary polymorphism pair. Black bar and striped bar represent polymorphism frequency differences in LeS
253 and LaS, respectively. The X-axis represents the name of the five amino acids namely Val, Pro, Thr, Ala and Gly. Y-axis
254 represents frequency values. In general, the pattern of complementary transition polymorphism of C→T and G→A in
255 LeS and LaS are found to be in opposite direction, whereas the other complementary transition polymorphism (A→G and
256 T→C) do not show contrasting pattern.

257

258      **The polymorphism at the four-fold degenerate site coincides with codon usage bias**

259          The amino acid specific polymorphisms at FFS indicated that the polymorphisms were most probably

260      influenced by codon usage bias. So, we compared polymorphisms with nucleotide frequency at FFS of the five

261      amino acids, which represented the codon usage bias of individual amino acids. Wide variation regarding

11

262 nucleotide frequency at FFS was observed across the five amino acids within a genome (Table 4). In general,

263 codon usage bias was observed to be strand independent as the frequency values were similar between the two

264 strands for an amino acid at FFS. This was in concordance with observations by earlier researchers (Sharp et

265 al. 2005; Shah and Gilchrist 2011; Wald et al. 2012). However, a moderate impact of strand asymmetry was

266 observed on codon usage bias because in general the frequencies of G and T at the FFS of an amino acid in

267 LeS was more than that in LaS. The reverse was true regarding the frequencies of A and C.

268 **Table 4: Nucleotide frequency at FFS**

| Bacteria | Base | Val | | Pro | | Thr | | Ala | | Gly | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS | LeS | LaS |
| *Ec* | A | 0.146 | 0.151 | 0.172 | 0.190 | 0.104 | 0.116 | 0.201 | 0.207 | 0.093 | 0.096 |
| | T | 0.255 | 0.250 | 0.141 | 0.151 | 0.166 | 0.151 | 0.159 | 0.152 | 0.342 | 0.338 |
| | G | 0.397 | 0.366 | 0.587 | 0.525 | 0.282 | 0.258 | 0.400 | 0.333 | 0.158 | 0.129 |
| | C | 0.202 | 0.233 | 0.100 | 0.134 | 0.448 | 0.475 | 0.240 | 0.308 | 0.406 | 0.437 |
| *Kp* | A | 0.090 | 0.082 | 0.068 | 0.074 | 0.031 | 0.034 | 0.063 | 0.058 | 0.062 | 0.066 |
| | T | 0.141 | 0.123 | 0.084 | 0.085 | 0.085 | 0.075 | 0.091 | 0.080 | 0.169 | 0.149 |
| | G | 0.497 | 0.445 | 0.713 | 0.661 | 0.254 | 0.223 | 0.463 | 0.381 | 0.203 | 0.166 |
| | C | 0.272 | 0.349 | 0.134 | 0.181 | 0.630 | 0.668 | 0.383 | 0.481 | 0.566 | 0.620 |
| *Se* | A | 0.155 | 0.155 | 0.112 | 0.122 | 0.082 | 0.087 | 0.116 | 0.119 | 0.107 | 0.105 |
| | T | 0.216 | 0.209 | 0.152 | 0.147 | 0.120 | 0.106 | 0.126 | 0.115 | 0.236 | 0.222 |
| | G | 0.388 | 0.334 | 0.605 | 0.555 | 0.371 | 0.328 | 0.492 | 0.416 | 0.176 | 0.130 |
| | C | 0.241 | 0.303 | 0.130 | 0.177 | 0.427 | 0.478 | 0.266 | 0.351 | 0.481 | 0.544 |
| *Sa* | A | 0.348 | 0.314 | 0.528 | 0.482 | 0.507 | 0.492 | 0.475 | 0.454 | 0.229 | 0.195 |
| | T | 0.410 | 0.412 | 0.318 | 0.396 | 0.277 | 0.311 | 0.304 | 0.348 | 0.562 | 0.548 |
| | G | 0.145 | 0.112 | 0.134 | 0.083 | 0.182 | 0.130 | 0.161 | 0.099 | 0.072 | 0.045 |
| | C | 0.098 | 0.161 | 0.021 | 0.039 | 0.035 | 0.066 | 0.060 | 0.099 | 0.137 | 0.213 |
| *Sp* | A | 0.202 | 0.171 | 0.490 | 0.453 | 0.344 | 0.305 | 0.269 | 0.234 | 0.295 | 0.330 |
| | T | 0.398 | 0.410 | 0.340 | 0.384 | 0.320 | 0.377 | 0.404 | 0.435 | 0.436 | 0.395 |
| | G | 0.188 | 0.122 | 0.092 | 0.050 | 0.113 | 0.070 | 0.114 | 0.061 | 0.134 | 0.101 |
| | C | 0.212 | 0.298 | 0.078 | 0.113 | 0.223 | 0.247 | 0.212 | 0.270 | 0.135 | 0.173 |

269 Table presents amino acid wise nucleotide frequencies in both LeS and LaS at FFS of five bacteria

270

271 To understand the relation of the polymorphism with codon usage bias, if any, we did a detailed

272 comparative study in individual bacteria. The polymorphism spectra at FFS for individual bacterium is given

273 in Supplementary Table 2. In case of *Ec*, the polymorphism pattern between complementary nucleotide pairs

274 revealed that Pro and Gly often behaved opposite to each other. Considering the previous knowledge that G-

275 ending codon in Pro (CCG) is the most preferred whereas the G-ending codon (GGG) in Gly is the least

276 preferred (Satapathy et al. 2016), we analysed both forward and reverse nucleotide polymorphisms. In case of

12

277    Gly, G→C and G→T transversions were more frequent than C→G and T→G respectively, that supported the

278    higher abundance of GGT/GGC over GGG in the genome. Whereas in case of Pro, G→C and G→T

279    transversions were found to be less frequent than C→G and T→G respectively, that supported the lower

280    abundance of CCT/CCC than CCG in the genome. A→T was more frequent than T→A in Val and Gly that

281    favoured higher abundance of GTT/GGT over GTA/GGA. A→T was less frequent than T→A in Pro that

282    favoured CCA to be more frequent than CCT. T→C was the most frequent in case of Thr among the five amino

283    acids that made ACC the most preferred codon. Low frequency of T→C in case of Pro and Val, favoured the

284    low frequency of GTC/CCC. While C→T and T→C were of similar frequency in case of Thr, C→T was more

285    than the frequency of T→C in case of Val, Pro and Ala that favoured the higher abundance of T-ending codon

286    (GTT/CCT/GCT) over the C-ending codons. G→A transition was more frequent than A→G in case of Gly,

287    which corresponds to the higher frequency of GGG over GGA. These observations suggested that

288    polymorphism at FFS was influenced by codon usage bas in *Ec*.

289    A similar comparative study of the polymorphism at FFS and codon usage bias in these amino acids

290    were studied in the other two γ-proteobacteria *Kp* and *Se*. In *Kp*, as the genome is G+C high, G/C-ending

291    codons are more abundant over the A/T-ending codons. In general A→G was more frequent than C→T and

292    G→A in all amino acids that favoured higher abundance of G-ending codons in the genome. A→C was more

293    frequent than C→A in all amino acids except Pro that supported CCC being preferred low here. Further, we

294    noticed correlation in preference of C-ending codons and higher frequency of G→C in Thr and Gly, whereas

295    the G-ending codons were preferred in Val and Pro that corresponded to the higher frequency of C→G. A→T

296    was more frequent than T→A in Val and Gly, which supported the T-ending codons being preferred over the

297    A-ending codons in these amino acids. T→G was more frequent than G→T in Pro while T→G was less

298    frequent than G→T in Gly which supported G-ending codon being preferred in Pro but not in Gly. Therefore,

299    the impact of codon usage bias was observed on the polymorphism at FFS of *Kp*. Similarly, in *Se*, G→C

300    transversion was more frequent than C→G in Thr and Gly while the reverse pattern was true in Pro in both the

301    strands. This was in concordance with the preference of C-ending codon in Thr, avoidance of C-ending codon

302    in Pro and G-ending codon in Gly. Polymorphism at A→T was more frequent than T→A in Val and Gly which

303    supported the higher abundance of GTT/GGT in the genome. Therefore, polymorphism pattern at FFS was

304    influenced by codon usage bias in *Se*.

305     In the two Firmicutes, A- and T-ending codons were the most frequent codons across the five amino

306     acids. In *Sa*, A-ending codons were more frequent than T-ending codons in Pro, Thr and Ala. In concordance,

307     A→T was more frequent than T→A in Val and Gly, while the reverse was the case in Pro, Thr and Ala. In

308     addition, G→A was more frequent than C→T in case of Thr because ACA is the most frequent codon here. In

309     *Sp*, A-ending codon was more frequent than T-ending codons in Pro. It was obvious to observe that A→T was

310     more frequent than T→A in case of Val, Ala and Gly but the reverse was true in case of Pro, which was in

311     concordance with the codon usage bias. Therefore, polymorphism at FFS in *Sa* as well as in *Sp* was influenced

312     by codon usage bias in these bacteria.

313     **Discussion**

314     Analysing nucleotide polymorphism in genome sequences of several strains belonging to a single

315     species provides avenue to study mechanisms of molecular evolution. Further the leading and the lagging

316     strands in chromosomes can be easily segregated using computational method facilitating researcher to

317     investigate the replication-associated mutation asymmetry between the strands in bacteria (Frank and Lobry

318     1999). Despite of low abundance of IRs, these regions can be distinctly identified because of the simplicity of

319     the coding sequences in bacterial chromosomes. In this study, nucleotide polymorphism has been analysed at

320     the intergenic regions (IRs) and four-fold degenerate site (FFS) of three bacterial species namely *E. coli* (*Ec*),

321     *K. pneumoniae* (*Kp*), *S. enterica* (*Se*) belonging to γ-proteobacteria and two other bacterial species namely *S.*

322     *aureus* (*Sa*) and *S. pneumoniae* (*Sp*) belonging to Firmicutes. Nucleotide compositional asymmetry between

323     the strands is well studied in bacterial chromosomes (Lobry 1996; Lobry and Sueoka 2002). The higher

324     abundance of the keto nucleotides (G, T) in the LeS than the LaS has been explained based on frequent cytosine

325     deamination in single stranded DNA (Reyes et al. 1998; Frank and Lobry 1999; Rocha 2004) which is

326     corroborated by the observation of positive GC skew and negative AT skew in the LeS in bacteria. However,

327     the higher magnitude of GC skew than AT skew in the LeS of most bacterial genomes could not be explained

328     based on the cytosine deamination theory. Rocha et al. (2006) had explained different mutation bias that might

329     lead to similar skew patterns in bacterial genomes. Similarly, the positive AT skew observed in the LeS of *S.*

330     *aureus* could not be explained by the cytosine deamination theory. Therefore, a detailed investigation has been

331     done in this study to understand the nucleotide compositional asymmetry between the strands in bacteria.

332     Polymorphism analysis at IRs of these five bacteria has revealed that C→T and G→A transition polymorphism

14

333　are the most frequent ones and display the main difference between the strands. This is in favour of the notion

334　that cytosine deamination is the major cause of polymorphism in genomes and the process is more frequent in

335　the LeS than the LaS. In a small magnitude, we have also observed that A→G and T→C contributes towards

336　strands asymmetry at IRs. In parallel with cytosine deamination theory, it may be hypothesized that more

337　frequent adenine deamination in LeS might result in higher A→G transition in the strand than the

338　complementary strand. It is known that cytosine deamination and adenine deamination have an opposite impact

339　on genome G+C%. Regarding transition polymorphisms at IRs, the five bacteria behave similar in this study.

340　However, in transversions, frequency of a polymorphism is observed to be similar between the strands and the

341　difference value between complementary transversions within strand is very low. Therefore, contribution of

342　transversion polymorphism in strand asymmetry is very low in the IRs of these bacteria. It is pertinent to note

343　that frequency of G→T (C→A) are higher than other transversions. Transversion polymorphisms increases

344　A/T at IRs like transition polymorphisms. But the bias towards A/T of these polymorphisms are more in the

345　two Firmicutes than the γ-Proteobacteria. The G→T (C→A) value is higher as well as A→G (T→C) value is

346　lower in Firmicutes in comparison to the γ-Proteobacteria for which A/T bias is observed to be more in the

347　former than the latter. The polymorphism study at the IRs suggests that, the replication associated strand

348　asymmetry is indifferent between the two groups of bacteria. Therefore, the atypical AT skew in the

349　chromosome of *Sa* is not supported by the polymorphism at IRs.

350　　　　In the two Firmicutes, *Sa* and *Sp* exhibit opposite patterns of nucleotide composition at FFS. The

351　nucleotide A is more frequent than T in *Sa*, while T is more frequent than A in *Sp*. The coding sequence is

352　more abundant in the leading strand than the lagging strand of the Firmicutes (Rocha 2004). Therefore, the

353　abundance of A is more than T in the LeS of *Sa* and the abundance of T is more than A in the LeS of *Sp*. Codon

354　usage bias at FFS of an amino acid is similar between the strands indicating the weak influence of strand

355　specific polymorphism. Therefore, the atypical AT skew in *Sa* can be attributed to codon usage bias, which is

356　due to the selection on codon usage bias. Our findings are in concordance with the earlier observation that

357　selection and gene distribution asymmetry between the strands was associated with the atypical AT skew in

358　*Sa* (Charneski et al. 2011).  Hence, our observations in this study suggest that selection on codon usage bias

359　influences the polymorphism at FFS in the Firmicutes. In conclusion, we would like to state that, our

360 understanding regarding the influence of codon usage bias on the compositional strand asymmetry become

361 clear in this study because of the polymorphism study done separately in IRs and FFS.

362  Genome G+C% in bacteria varies from 13.0 to 75.0% which accounts a difference of 62% between

363 the minimum and maximum genome composition (Raghavan et al. 2012). Directional mutation bias in support

364 of the neutral theory of evolution has been proposed to explain genome G+C% in organisms (Sueoka 1988).

365 In support of directional mutation theory, Muto and Osawa (1987) demonstrated that synonymous codon usage

366 varies between high and low G+C bacterial genomes. But theoretical analysis of the G+C% of the synonymous

367 codons suggests that the maximum G+C composition difference between two synonymous codons (e.g., GGU

368 and GGC) of an amino acid can be 33.33% with exceptions only in Arg (e.g., CGG, AGU) and Leu (e.g.,

369 UUA, CUG) where the difference can be up to 66.67%. Hence, the synonymous codon usage range should be

370 held accountable to a value around 33.33% instead of 62.0% (75.0 – 13.0). It can be argued that the directional

371 mutation theory inadequately explains the genome G+C composition in bacteria as IRs contribute a minor

372 portion of the genome size. It is pertinent to note that the results of earlier mutation analysis in bacterial

373 genomes could not provide substantial evidence in support of the directional mutation theory (Hershberg and

374 Petrov 2010; Hildebrand et al. 2010; Rocha and Feil 2010). Hence, the possible existence of an unknown

375 selection mechanism responsible for genome G+C% has been hypothesized (Rocha and Feil 2010; Raghavan

376 et al. 2012). The role of recombination in determining the genome composition of bacteria have also been

377 implicated (Bobay and Ochman 2017). G+C% variation at FFS across the amino acids in different genomes

378 reported in this work is an indication of codon usage bias influencing the observed difference in genome

379 composition. Notably, we found that Pro, Thr and Ala codons having C at the second position differ from each

380 other in terms of codon usage bias. Similarly, Val, Ala and Gly codons having G at the first position also differ

381 from each other in terms of codon usage bias. Therefore, the possibility of context-dependent polymorphism

382 causing codon usage difference across the amino acids is not anticipated. The polymorphism difference

383 observed across the amino acids at FFS suggest that the difference is due to the amino acid specific translational

384 selection. Our comparative analysis of polymorphism and codon usage bias in the five studied bacteria led us

385 to believe that the selection on codon usage bias responsible for the observed polymorphism. It is already

386 known that GGG and CCC codons are prone to translational frameshift (O'Connor 1998, 2002). Interestingly,

387 GGG and CCC codons were not preferred in both the strands of five studied bacteria. Transversions such as

388    G→C is more than C→G in case of Gly while the same is less in case of Pro. GGC codon has been reported

389    to be selected positively in bacterial genomes (Satapathy et al. 2014, 2016). Further amino acid specific codon

390    usage bias is similar in the two strands indicating a weak influence of the strand specific mutation bias, in

391    comparison to the translational selection in genomes. This observation supports that G+C% variations across

392    different amino acids at FFS is due to selection on codon usage bias and may not be due to the directional

393    mutation. Further, considering a limited set of high expression genes, we observed that the polymorphism at

394    FFS of the high expression genes is in line with that at FFS of whole genome. It is interesting to note that the

395    earlier notion of genome composition determining the codon usage bias (Muto and Osawa 1987) is found

396    inconsistent in this study. Now our analysis using large number of genomes of γ-Proteobacteria and Firmicutes

397    have suggested the role of codon usage bias in determining the genome G+C% supporting the selection theory

398    of evolution. Assuming that the selection theory is true, it is speculated that in an AT rich genome, A/T ending

399    synonymous codons are likely to be selected over the G/C ending synonymous codons and *vice versa* is true

400    for GC rich genomes (Hershberg and Petrov 2009). We anticipated that future research on translation rate of

401    synonymous codons is expected to uncover the mechanism of genome composition in AT and GC rich bacteria.

402    Large range of genome G+C% is a classic example of the neutral theory of evolution in bacteria which

403    means that there is no specific advantage that could be linked to genome composition (Lassalle et al. 2015).

404    Under this assumption, the low genomic G+C content of endosymbiotic bacteria were considered in favour of

405    the neutral theory. However, recently Dietel et al. (2019) have discussed that selective advantages favour high

406    genomic AT-contents in intracellular genetic elements. Genome composition is known to be associated with

407    bacterial phylogeny such as Firmicutes with low genome G+C%, Actinomycetes and β-Proteobacteria with

408    high genome G+C% (Satapathy et al. 2010). However, the reason for the high and low G+C% in these phyla

409    is not clearly understood. But phylogeny specific optimal codon selection has been reported recently

410    (Satapathy et al. 2016). Future understanding of translational decoding by the ribosome might explain the

411    phylogeny specific codon usage bias and genome composition. It is pertinent to note that ribosome mediated

412    gene regulation by co-translational protein folding has been demonstrated to be species specific in *E. coli* and

413    *B. subtilis* (Sohmen et al. 2015). It has been reported that genome size and G+C% are positively correlated

414    (Satapathy et al. 2010). However, we have observed a lack of correlation between the genome size and genome

415    G+C% in different phylogeny (Supplementary Table 3). It has been reported that the strength of selection on

416    codon usage bias is variable among the bacteria (Sharp et al. 2005; Satapathy et al. 2012, 2014, 2016). In that

417    case bacteria with poor selection on codon usage bias should exhibit low genome G+C%, while bacteria with

418    high genome G+C% must exhibit high selection on codon usage bias. In a different study, it has been shown

419    that bacteria with high genome G+C% indeed exhibited selection on codon usage bias (Satapathy et al. 2014).

420    Hence, this further supports that the selection of codon usage bias is responsible for genome composition in

421    bacteria. As codon usage bias is universal in bacteria, it may be possible that the difference between two

422    genomes regarding codon usage bias may act as a selection against lateral gene transfer in bacteria.

423    **Materials and Methods**

424    **Segregating the leading and the lagging strands in bacterial chromosomes**

425        We carried out a detailed single nucleotide polymorphism study using computational analysis of

426    genomes of total 157 *Escherichia coli* (*Ec*) strains (Thorpe et al. 2017), 208 *Klebsiella pneumoniae* (*Kp*) strains

427    (Holt et al. 2015), 366 *Salmonella enterica* (*Se*) strains (Thorpe et al. 2017), 132 *Staphylococcus aureus* (*Sa*)

428    strains (Reuter et al. 2016) and 264 *Streptococcus pneumoniae* (*Sp*) strains (Chewapreecha et al. 2014).

429    Considering cumulative GC skew diagram for each bacterium, we segregated respective chromosomes into

430    the leading strand (LeS) and the lagging strand (LaS) as has been described earlier by different researchers

431    (Lobry 1996; Grigoriev 1998). The method is mentioned below in brief. We found out abundance values of

432    each nucleotide along a genome sequence using non-overlapping moving window of size 1.0 kb. GC skew was

433    calculated as (G-C)/(G+C). Similarly AT skew was calculated as (A-T)/(A+T), RY skew was calculated as

434    [{(A+G)-(C+T)}/(A+G+C+T)], and KM skew was calculated as [{(G+T)-(A+C)}/(A+G+C+T)]. For each of

435    the bacteria, cumulative skew diagrams were generated from these deviations, which was used to identify the

436    leading (LeS) and the lagging (LaS) strands regions in a chromosome (Supplementary Figure 1a and 1b). GC

437    skew is positive in the LeS whereas the same is negative in the LaS. Similarly, KM skew is positive in the LeS

438    whereas the same is negative in the LaS. A schematic view of the LeS and the LaS in a double stranded DNA

439    is presented in Figure 4. The LeS and the LaS regions of the Watson strand are aligned with the LaS and the

440    LeS, respectively, of the Crick strand in chromosomes. Using the coordinates of protein coding genes (CDS)

441    from the genome annotation, the CDS were mapped to the LeS and the LaS. Sequences other than those coding

442    for the rRNA, tRNA,protein genes and misc_RNAs were considered as intergenic regions (IRs). IRs in the

443    LeS and the LaS of the Watson strand are aligned opposite to IRs in the LaS and the LeS, respectively, of the

18

444     Crick strand (Figure 4). For polymorphism analysis, IRs in either the Watson or the Crick strand were

445     considered. IRs belonging to either the Watson or the Crick strand were segregated into the LeS and the LaS

446     (Supplementary Figure 1a, 1b and Figure 4). In our analysis, we included only large IRs (size greater than 100

447     bases) from where 35 bases from both the ends of each IR were ignored to minimize the inclusion of any

448     regulatory regions and considered only the nearly neutral polymorphisms.

449     **Figure 4: A schematic view of the distribution of IRs, CDS, tRNA and rRNA in the leading and lagging**

450     **strands in double stranded DNA**



451

452     Figure represents a schematic view of the LeS/LaS and distribution of IRs and CDS, tRNA and rRNA in a double stranded
453     DNA. Considering nucleotide composition of the Watson strand, the skew diagram is generated. In the strand, the region
454     between point (a) and point (b) with positive GC skew is designated as the LeS and the sequence between point (b) and
455     point (c) as the LaS. The LeS and the LaS regions of the Watson strand are aligned with the LaS and the LeS, respectively,
456     of the Crick strand in chromosomes.

457

458     **Polymorphism analysis at IRs and at FFS in CDS of the bacterial chromosomes**

459        For each bacterium, we extracted alignments of intergenic regions (IRs) and protein-coding regions

460     (CDS) using computer programs written in Python script. Considering the most frequent nucleotide at a

461     position in the sequence alignment, we computed a consensus sequence which was used to identify

462     polymorphisms at different positions. A detailed description of the approach used for analysing polymorphisms

463     in this study is provided in the Supplementary Document 1. Phylogenetic relationships among the five bacteria

464     (Supplementary Figure 2a) were obtained using the reference sequence of *rpoB* and *rpoC* genes by the

465     MEGAX software (Kumar et al. 2018). The three γ-proteobacteria (*Ec*, *Kp* and *Se*) and two Firmicutes (*Sa* and

466     *Sp*) made different clusters. Further, phylogenetic relationships among the population (Supplementary Figure

467     2b-2f) were constructed using the *rpoB* sequence of all strains of individual organisms (Kumar et al. 2018).

468     The distribution of polymorphisms among the strains in reference to *rpoB* and *rpoC* genes (Supplementary

469     Figure 3) shows that, in general, polymorphism observed in this study is not because of any specific strain but

470     mutations accumulated among all the strains (Supplementary Table 4). A known set of previously published

471     high expression genes (Sharp et al. 2005; Sen et al. 2020) were considered in this work for the analysis of

472     nucleotide composition and polymorphism at FFS in high expression genes (HEGs) (Supplementary Table 5

473     and 6).

474     In coding sequences (CDS), we considered polymorphism at FFS of the amino acids such as Val, Pro,

475     Thr, Ala and Gly. For example, if a nucleotide change (suppose A$\rightarrow$T) observed at the $3^{rd}$ position of codon,

476     the corresponding codon was found out in the reference sequence (considering the preceding two nucleotides).

477     If the codon codes for Val (i.e., the codon is GTT/GTC/GTA/GTG), then we increase A$\rightarrow$T change for Val

478     by 1. Using this approach, we calculated polymorphism at FFS of the five amino acids. Polymorphism

479     frequencies were normalized by dividing the total count of a given change by the total count of the nucleotide

480     in which polymorphism has occurred in the reference sequence. For example, if the total number of C$\rightarrow$T

481     change is 20 and the total number of C in the reference sequence (either at IRs or at FFS of that amino acid) is

482     100, then the normalized frequency is calculated as 20/100 = 0.2. The frequencies of different nucleotide

483     polymorphisms were calculated accordingly. For statistical analysis and determining p-value for significance

484     test, Mann Whitney test is used (Mann and Whitney 1947).

485     **Competing interest statement**

486     The authors declare no competing interests.

487     **Acknowledgements**

499

**References**

501    Balbi KJ, Rocha EPC, Feil EJ. 2009. The temporal dynamics of slightly deleterious mutations in Escherichia

502        coli and Shigella spp. *Mol Biol Evol* **26**: 345–355.

503    Bobay L-M, Ochman H. 2017. Impact of recombination on the base composition of bacteria and archaea.

504        *Mol Biol Evol* **34**: 2627–2636.

505    Charneski CA, Honti F, Bryant JM, Hurst LD, Feil EJ. 2011. Atypical AT skew in Firmicute genomes results

506        from selection and not from mutation. *PLoS Genet* **7**: e1002283.

507    Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, Mather AE, Hanage WP, Goldblatt D,

508        Nosten FH, Turner C. 2014. Comprehensive identification of single nucleotide polymorphisms

509        associated with beta-lactam resistance within pneumococcal mosaic genes. *PLoS Genet* **10**: e1004547.

510    Davis BD. 1989. Transcriptional bias: a non-Lamarckian mechanism for substrate-induced mutations. *Proc*

511        *Natl Acad Sci* **86**: 5005–5009.

512    Duchêne S, Ho SYW, Holmes EC. 2015. Declining transition/transversion ratios through time reveal

513        limitations to the accuracy of nucleotide substitution models. *BMC Evol Biol* **15**: 1–10.

514    Forsdyke DR, Mortimer JR. 2000. Chargaff's legacy. *Gene* **261**: 127–137.

515    Francino MP, Ochman H. 1997. Strand asymmetries in DNA evolution. *Trends Genet* **13**: 240–245.

516    Frank AC, Lobry JR. 1999. Asymmetric substitution patterns: a review of possible underlying mutational or

517        selective mechanisms. *Gene* **238**: 65–77.

518    Gaillard H, Herrera-Moyano E, Aguilera A. 2013. Transcription-associated genome instability. *Chem Rev*

519     **113**: 8638–8661.

520     Grigoriev A. 1998. Analyzing genomes with cumulative skew diagrams. *Nucleic Acids Res* **26**: 2286–2290.

521     Hershberg R, Petrov DA. 2010. Evidence that mutation is universally biased towards AT in bacteria. *PLoS*
522         *Genet* **6**: e1001115.

523     Hershberg R, Petrov DA. 2009. General rules for optimal codon choice. *PLoS Genet* **5**: e1000556.

524     Hildebrand F, Meyer A, Eyre-Walker A. 2010. Evidence of selection upon genomic GC-content in bacteria.
525         *PLoS Genet* **6**: e1001107.

526     Holt KE, Wertheim H, Zadoks RN, Baker S, Whitehouse CA, Dance D, Jenney A, Connor TR, Hsu LY,
527         Severin J. 2015. Genomic analysis of diversity, population structure, virulence, and antimicrobial
528         resistance in Klebsiella pneumoniae, an urgent threat to public health. *Proc Natl Acad Sci* **112**: E3574–
529         E3581.

530     Ito K. 2016. Strolling Toward New Concepts. *Annu Rev Microbiol* **70**: 1–23.

531     Jinks-Robertson S, Bhagwat AS. 2014. Transcription-associated mutagenesis. *Annu Rev Genet* **48**: 341–359.

532     Karlin S, Campbell AM, Mrazek J. 1998. Comparative DNA analysis across diverse genomes. *Annu Rev*
533         *Genet* **32**: 185–225.

534     Kim N, Jinks-Robertson S. 2012. Transcription as a source of genome instability. *Nat Rev Genet* **13**: 204–
535         214.

536     Kumar S, Stecher G, Li M, Knyaz C, Tamura K. 2018. MEGA X: molecular evolutionary genetics analysis
537         across computing platforms. *Mol Biol Evol* **35**: 1547–1549.

538     Lassalle F, Périan S, Bataillon T, Nesme X, Duret L, Daubin V. 2015. GC-content evolution in bacterial
539         genomes: the biased gene conversion hypothesis expands. *PLoS Genet* **11**: e1004941.

540     Lewis CA, Crayle J, Zhou S, Swanstrom R, Wolfenden R. 2016. Cytosine deamination and the precipitous
541         decline of spontaneous mutation during Earth's history. *Proc Natl Acad Sci* **113**: 8194–8199.

542     Lind PA, Andersson DI. 2008. Whole-genome mutational biases in bacteria. *Proc Natl Acad Sci* **105**: 17878–

543      17883.

544  Lobry JR. 1996. Asymmetric substitution patterns in the two DNA strands of bacteria. *Mol Biol Evol* **13**:

545      660–665.

546  Lobry JR, Sueoka N. 2002. Asymmetric directional mutation pressures in bacteria. *Genome Biol* **3**: 1–14.

547  Lyons DM, Lauring AS. 2017. Evidence for the selective basis of transition-to-transversion substitution bias

548      in two RNA viruses. *Mol Biol Evol* **34**: 3205–3215.

549  Mann HB, Whitney DR. 1947. On a test of whether one of two random variables is stochastically larger than

550      the other. *Ann Math Stat* 50–60.

551  McLean MJ, Wolfe KH, Devine KM. 1998. Base composition skews, replication orientation, and gene

552      orientation in 12 prokaryote genomes. *J Mol Evol* **47**: 691–696.

553  Mugal CF, von Grünberg H-H, Peifer M. 2009. Transcription-induced mutational strand bias and its effect

554      on substitution rates in human genes. *Mol Biol Evol* **26**: 131–142.

555  Muto A, Osawa S. 1987. The guanine and cytosine content of genomic DNA and bacterial evolution. *Proc

556      Natl Acad Sci* **84**: 166–169.

557  O'Connor M. 2002. Imbalance of tRNAPro isoacceptors induces+ 1 frameshifting at near-cognate codons.

558      *Nucleic Acids Res* **30**: 759–765.

559  O'Connor M. 1998. tRNA imbalance promotes− 1 frameshifting via near-cognate decoding. *J Mol Biol* **279**:

560      727–736.

561  Powdel BR, Satapathy SS, Kumar A, Jha PK, Buragohain AK, Borah M, Ray SK. 2009. A study in entire

562      chromosomes of violations of the intra-strand parity of complementary nucleotides (Chargaff's second

563      parity rule). *DNA Res* **16**: 325–343.

564  Raghavan R, Kelkar YD, Ochman H. 2012. A selective force favoring increased G+ C content in bacterial

565      genes. *Proc Natl Acad Sci* **109**: 14504–14507.

566  Ray SK, Goswami I. 2016. Synonymous codons are not the same with respect to the speed of translation

567      elongation. *Curr Sci* **110**: 1612–1614.

23

568  Reuter S, Török ME, Holden MTG, Reynolds R, Raven KE, Blane B, Donker T, Bentley SD, Aanensen DM,
569       Grundmann H. 2016. Building a genomic framework for prospective MRSA surveillance in the United
570       Kingdom and the Republic of Ireland. *Genome Res* **26**: 263–270.

571  Reyes A, Gissi C, Pesole G, Saccone C. 1998. Asymmetrical directional mutation pressure in the
572       mitochondrial genome of mammals. *Mol Biol Evol* **15**: 957–966.

573  Rocha EPC. 2004. The replication-related organization of bacterial genomes. *Microbiology* **150**: 1609–1627.

574  Rocha EPC, Danchin A. 2001. Ongoing evolution of strand composition in bacterial genomes. *Mol Biol Evol*
575       **18**: 1789–1799.

576  Rocha EPC, Danchin A, Viari A. 1999. Universal replication biases in bacteria. *Mol Microbiol* **32**: 11–16.

577  Rocha EPC, Feil EJ. 2010. Mutational patterns cannot explain genome composition: are there any neutral
578       sites in the genomes of bacteria? *PLoS Genet* **6**: e1001104.

579  Rocha EPC, Touchon M, Feil EJ. 2006. Similar compositional biases are caused by very different mutational
580       effects. *Genome Res* **16**: 1537–1547.

581  Saha SK, Goswami A, Dutta C. 2014. Association of purine asymmetry, strand-biased gene distribution and
582       PolC within Firmicutes and beyond: a new appraisal. *BMC Genomics* **15**: 1–26.

583  Satapathy SS, Dutta M, Buragohain AK, Ray SK. 2012. Transfer RNA gene numbers may not be completely
584       responsible for the codon usage bias in asparagine, isoleucine, phenylalanine, and tyrosine in the high
585       expression genes in bacteria. *J Mol Evol* **75**: 34–42.

586  Satapathy SS, Dutta M, Ray SK. 2010. Variable correlation of genome GC% with transfer RNA number as
587       well as with transfer RNA diversity among bacterial groups: α-Proteobacteria and Tenericutes exhibit
588       strong positive correlation. *Microbiol Res* **165**: 232–242.

589  Satapathy SS, Powdel BR, Buragohain AK, Ray SK. 2016. Discrepancy among the synonymous codons with
590       respect to their selection as optimal codon in bacteria. *DNA Res* **23**: 441–449.

591  Satapathy SS, Powdel BR, Dutta M, Buragohain AK, Ray SK. 2014. Selection on GGU and CGU codons in
592       the high expression genes in bacteria. *J Mol Evol* **78**: 13–23.

593    Schroeder JW, Randall JR, Hirst WG, O'Donnell ME, Simmons LA. 2017. Mutagenic cost of

594        ribonucleotides in bacterial DNA. *Proc Natl Acad Sci* **114**: 11733–11738.

595    Sen P, Tula D, Ray SK, Satapathy SS. 2020. Estimating RNA Secondary Structure by Maximizing Stacking

596        Regions. In *Applications of Internet of Things*, pp. 165–176, Springer.

597    Seplyarskiy VB, Kharchenko P, Kondrashov AS, Bazykin GA. 2012. Heterogeneity of the

598        transition/transversion ratio in Drosophila and Hominidae genomes. *Mol Biol Evol* **29**: 1943–1955.

599    Shah P, Gilchrist MA. 2011. Explaining complex codon usage patterns with selection for translational

600        efficiency, mutation bias, and genetic drift. *Proc Natl Acad Sci* **108**: 10231–10236.

601    Sharp PM, Bailes E, Grocock RJ, Peden JF, Sockett RE. 2005. Variation in the strength of selected codon

602        usage bias among bacteria. *Nucleic Acids Res* **33**: 1141–1153.

603    Sohmen D, Chiba S, Shimokawa-Chiba N, Innis CA, Berninghausen O, Beckmann R, Ito K, Wilson DN.

604        2015. Structure of the Bacillus subtilis 70S ribosome reveals the basis for species-specific stalling. *Nat*

605        *Commun* **6**: 1–10.

606    Stoltzfus A, Norris RW. 2016. On the causes of evolutionary transition: transversion bias. *Mol Biol Evol* **33**:

607        595–602.

608    Sueoka N. 1988. Directional mutation pressure and neutral molecular evolution. *Proc Natl Acad Sci* **85**:

609        2653–2657.

610    Sueoka N. 1995. Intrastrand parity rules of DNA base composition and usage biases of synonymous codons.

611        *J Mol Evol* **40**: 318–325.

612    Suzuki T, Kamiya H. 2017. Mutations induced by 8-hydroxyguanine (8-oxo-7, 8-dihydroguanine), a

613        representative oxidized base, in mammalian cells. *Genes Environ* **39**: 1–6.

614    Thorpe HA, Bayliss SC, Hurst LD, Feil EJ. 2017. Comparative analyses of selection operating on

615        nontranslated intergenic regions of diverse bacterial species. *Genetics* **206**: 363–376.

616    Van Leuven JT, McCutcheon JP. 2012. An AT mutational bias in the tiny GC-rich endosymbiont genome of

617        Hodgkinia. *Genome Biol Evol* **4**: 24–27.

618    Loon BV, Markkanen E, Hübscher U. 2010. Oxygen as a friend and enemy: How to combat the mutational

619        potential of 8-oxo-guanine. *DNA Repair (Amst)* **9**: 604–616.

620    Wald N, Alroy M, Botzman M, Margalit H. 2012. Codon usage bias in prokaryotic pyrimidine-ending

621        codons is associated with the degeneracy of the encoded amino acids. *Nucleic Acids Res* **40**: 7074–

622        7083.