

# Intergenic Transposable Elements Are Not Randomly Distributed in Bacteria

Gordon R. Plague\*

Louis Calder Center—Biological Field Station, Department of Biological Sciences, Fordham University, Armonk, New York

\*Corresponding author: E-mail: plague@fordham.edu.

Accepted: 3 July 2010

## Abstract

Insertion sequences (ISs) are mobile genetic elements in bacterial genomes. In general, intergenic IS elements are probably less deleterious for their hosts than intragenic ISs, simply because they have a lower likelihood of disrupting native genes. However, since promoters, Shine–Dalgarno sequences, and transcription factor binding sites are intergenic and upstream of genes, I hypothesized that not all neighboring gene orientations (NGOs) are selectively equivalent for IS insertion. To test this, I analyzed the NGOs of all intergenic ISs in 326 fully sequenced bacterial chromosomes. Of the 116 genomes with enough IS elements for statistical analysis, 68 have significantly more ISs between convergently oriented genes than expected, and 46 have significantly fewer ISs between divergently oriented genes. This suggests that natural selection molds intergenic IS distributions because they are least intrusive between convergent gene pairs and most intrusive between divergent gene pairs.

**Key words:** insertion sequence, IS element, natural selection.

## Are All Intergenic Regions Created Equally?

Insertion sequences (ISs) are common transposable elements in bacterial genomes. Although IS elements can generate beneficial mutations (Cooper et al. 2001; Safi et al. 2004), they are generally considered genomic parasites because they only code for the enzyme required for their own transposition (Doolittle and Sapienza 1980; Orgel and Crick 1980). While an IS element inhabits a chromosomal location, it is inherited along with its host's native genes, so its fitness is intimately tied to that of its host. Therefore, an IS that causes a deleterious mutation by disrupting an essential gene will probably be quickly eliminated from most natural populations, whereas an IS that inserts into a selectively neutral location will have a much greater chance of long-term survival (Lynch 2006). As a general rule, intergenic IS elements probably enjoy higher survival than those that integrate within genes, simply because they have a lower likelihood of disrupting native genes (Campbell 2002; Zaghloul et al. 2007). However, the question then arises: are all intergenic regions selectively equivalent for IS occupancy?

Bacterial genes can be transcribed from either the top ( $\rightarrow$ ) or bottom ( $\leftarrow$ ) DNA strand. Therefore, neighboring

genes on bacterial chromosomes can occur in three possible orientations: tandem ( $\rightarrow \rightarrow$  and  $\leftarrow \leftarrow$ ), convergent ( $\rightarrow \leftarrow$ ), and divergent ( $\leftarrow \rightarrow$ ). Because promoters, Shine–Dalgarno sequences, and transcription factor binding sites are upstream of genes, I hypothesized that the intergenic regions of the three neighboring gene orientations (NGOs) may not be selectively equivalent for IS insertion. Specifically, the intergenic region between: 1)  $\leftarrow \rightarrow$  neighbors will contain a promoter and a Shine–Dalgarno sequence for both genes, and possibly a transcription factor binding site for both, 2)  $\rightarrow \rightarrow$  and  $\leftarrow \leftarrow$  neighbors will contain a promoter (if the neighbors are not in the same operon) and a Shine–Dalgarno sequence for the respective downstream gene only, and possibly a transcription factor binding site for that gene, and 3)  $\rightarrow \leftarrow$  neighbors will contain no promoters, Shine–Dalgarno sequences, or transcription factor binding sites. Therefore, an IS that inserts between  $\leftarrow \rightarrow$  genes has a relatively high likelihood of disrupting the transcription or translation of its neighbors, an IS that inserts between  $\rightarrow \rightarrow$  or  $\leftarrow \leftarrow$  genes has a moderate likelihood of disrupting its neighbors, and an IS that inserts between  $\rightarrow \leftarrow$  genes will never disrupt its neighbors. Because of this discrepancy among intergenic regions, I hypothesized that intergenic ISs would be most common between  $\rightarrow \leftarrow$  oriented genes and

© The Author(s) 2010. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

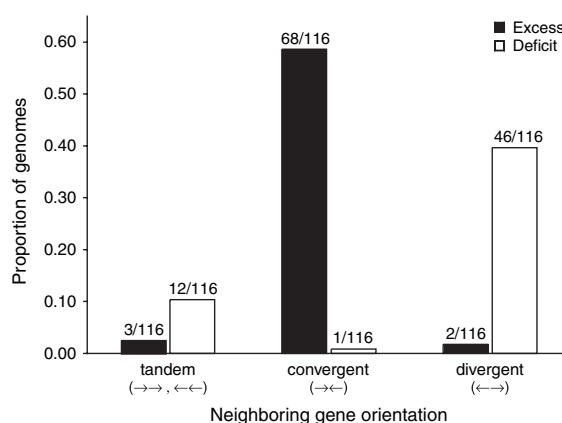
This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/2.5>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

least common between  $\leftarrow \rightarrow$  oriented genes in bacterial genomes.

## Intergenic IS Elements Are Not Randomly Distributed

I tested this hypothesis by analyzing the NGOs of all intergenic ISs from 326 fully sequenced bacterial chromosomes. Of these, 116 genomes have enough ISs to meet  $\chi^2$  test assumptions (Cochran 1954). Remarkably, 64% of these genomes ( $N = 74$ ) have observed intergenic IS quantities that deviate significantly ( $P \leq 0.05$ ) from expectations (under the null assumptions of random insertion and no natural selection) (table 1). These deviations are pervasive across the phylogenetic spectrum of Bacteria (table 1) and include a wide variety of IS families. Two NGOs exhibit extraordinary consistency in their contributions to these deviations:  $\rightarrow \leftarrow$  harbors significant IS excesses in 68 genomes and one significant deficit, and  $\leftarrow \rightarrow$  harbors two significant IS excesses and 46 significant deficits (fig. 1 and table 1). Overall, 105 of the 116 analyzed genomes contain more IS elements in the  $\rightarrow \leftarrow$  orientation than expected, and 104 contain fewer in the  $\leftarrow \rightarrow$  orientation than expected (the binomial probabilities of having distributions at least this skewed just by chance are  $1.1 \times 10^{-20}$  and  $9.3 \times 10^{-20}$ , respectively) (table 1). These nonrandom IS distributions also extend to bacterial chromosomes that contain relatively few IS elements. Specifically, of the 131 genomes that do not contain enough ISs for statistical analysis (Cochran 1954) but that have  $\geq 1$  expected IS in each NGO, 117 genomes contain more IS elements in the  $\rightarrow \leftarrow$  orientation than expected, and 108 contain fewer in the  $\leftarrow \rightarrow$  orientation than expected (the binomial probabilities of having distributions at least this skewed just by chance are  $1.0 \times 10^{-21}$  and  $1.1 \times 10^{-14}$ , respectively) (supplementary table S1, Supplementary Material online).

One possible explanation for these nonrandom IS distributions is a general insertion bias into  $\rightarrow \leftarrow$  and away from  $\leftarrow \rightarrow$  intergenic regions. I doubt that such a bias would result from target sequence specificity, largely because IS target site preferences are rarely very stringent or very long (Chandler and Mahillon 2002), so suitable insertion locations for many ISs occur thousands of times in each genome (Zaghloul et al. 2007). Instead, insertion bias could result from chromosomal differences between the three NGOs. For example, as bacterial genes are transcribed, DNA becomes positively supercoiled ahead of the polymerase and negatively supercoiled behind (Liu and Wang 1987). Consequently, the region between  $\rightarrow \leftarrow$  oriented genes may often be positively supercoiled, more so than between the other NGOs (and conversely, the region between  $\leftarrow \rightarrow$  genes may often be the most negatively supercoiled). If IS elements preferentially insert into positively supercoiled DNA, then this could explain the overabundances and



**FIG. 1.**—Proportion of fully sequenced bacterial chromosomes with a significant excess or deficit of IS elements in each NGO. Each bar is labeled with the number of excesses or deficits relative to the number of genomes analyzed (i.e., the number of genomes with enough IS elements for statistical analysis; see text).

underabundances of ISs between  $\rightarrow \leftarrow$  and  $\leftarrow \rightarrow$  oriented genes, respectively (fig. 1). However, no evidence exists for such an insertion bias, and some transposons prefer the opposite: negatively supercoiled DNA (Lodge and Berg 1990). Another possibility is that IS elements generally preferentially insert downstream of genes; for example, near transcription termination sequences. At least one IS element exhibits such a preference (Tetu and Holmes 2008), although this is not a ubiquitous tendency among ISs because some exhibit the opposite preference, inserting upstream of genes between Shine–Dalgarno sequences and start codons (Doran et al. 1997; Inglis et al. 2003). Therefore, insertion bias may affect the distribution of some IS elements in some bacterial genomes, although it is unlikely to explain the widespread bias exhibited across Bacteria (table 1).

Without any evidence for systematic IS insertion bias to explain these nonrandom IS distributions (table 1), the most likely explanation at present is that natural selection molds intergenic IS distributions. From a host bacterium's perspective, all potential IS insertion locations are not equally viable, and natural selection eventually eliminates disadvantageous genotypes from most populations. In fact, few IS elements are probably truly selectively neutral because at the very least they appropriate host resources for transposase expression (Nuzhdin 1999). So unless a particular IS element beneficially impacts its host (Safi et al. 2004), the likely fate of most ISs is eventual extinction from their host population (Wagner 2006). For an individual IS locus, the likelihood of extinction is largely correlated to its fitness cost, with the most deleterious ISs eliminated most quickly, and those inserting in innocuous locations having the greatest potential for long-term survival (Lynch 2006). Therefore, the most innocuous ISs will be overrepresented in bacterial genomes, and the most deleterious will be underrepresented. The

**Table 1.**Observed (O) and Expected (E) Quantities of Intergenic IS Elements in Fully Sequenced Bacterial Chromosomes, and the  $\chi^2$  Test Statistic for Each

	NGO <sup>a</sup>						$\chi^2$ <sup>b</sup>
	$\rightarrow \rightarrow$ , $\leftarrow \leftarrow$		$\rightarrow \leftarrow$		$\leftarrow \rightarrow$		
	O	E	O	E	O	E	
Actinobacteria							
<i>Corynebacterium efficiens</i> YS-314	36	39.8	21	14.5	16	18.6	3.6
<i>Corynebacterium glutamicum</i> ATCC 13032	21	20.5	11	7.4	7	11.1	3.3
<i>Corynebacterium jeikeium</i> K411	38	38.1	<b>18</b>	<b>8.5</b>	12	21.5	14.9***
<i>Frankia</i> sp. Ccl3	88	79.1	<b>27</b>	<b>16.6</b>	21	40.4	16.9*** <sup>c</sup>
<i>Mycobacterium avium</i> paratuberculosis	29	33.8	<b>14</b>	<b>7.6</b>	14	15.6	6.3*
<i>Mycobacterium bovis</i> AF2122/97	22	22.7	<b>15</b>	<b>6.6</b>	6	13.7	15.1***
<i>Mycobacterium smegmatis</i> MC2	33	37.9	<b>24</b>	<b>7.1</b>	11	23.0	47.3*** <sup>c</sup>
<i>Mycobacterium tuberculosis</i> CDC1551	29	26.6	<b>13</b>	<b>6.7</b>	6	14.7	11.3**
<i>M. tuberculosis</i> H37Rv	31	27.9	<b>16</b>	<b>8.1</b>	6	17.0	15.2***
<i>Streptomyces avermitilis</i> MA-4680	36	38.1	<b>25</b>	<b>11.6</b>	8	19.3	22.3*** <sup>c</sup>
<i>Streptomyces coelicolor</i> A3(2)	19	20.6	9	6.1	11	12.3	1.7
Bacteroidetes							
<i>Bacteroides thetaiotaomicron</i> VPI-5482	22	30.9	<b>24</b>	<b>7.2</b>	6	13.8	45.7*** <sup>c</sup>
<i>Porphyromonas gingivalis</i> W83	21	27.3	10	8.7	13	8.0	4.7
<i>Prevotella intermedia</i> 17	26	27.1	14	10.1	9	11.8	2.3
<i>Salinibacter ruber</i> DSM 13855	24	23.8	4	5.6	11	9.6	0.6
Chlamydiae							
<i>Protochlamydia amoebophila</i> UWE25	26	32.3	14	8.0	17	16.7	5.7
Cyanobacteria							
<i>Anabaena variabilis</i> ATCC 29413	35	31.5	<b>16</b>	<b>8.1</b>	7	18.4	15.0***
<i>Gloeobacter violaceus</i> PCC7421	41	35.0	<b>18</b>	<b>10.1</b>	6	19.9	17.0*** <sup>c</sup>
<i>Nostoc</i> sp. PCC 7120	39	37.1	<b>17</b>	<b>8.7</b>	11	21.1	12.8**
<i>Synechococcus</i> sp. JA-2-3Ba(2-13)	42	37.6	21	27.7	16	13.7	2.5
<i>Synechococcus</i> sp. JA-3-3Ab	44	40.4	17	29.9	<b>21</b>	<b>11.7</b>	13.2***
<i>Synechocystis</i> sp. PCC6803	30	30.4	<b>17</b>	<b>6.9</b>	5	14.7	21.0*** <sup>c</sup>
<i>Thermosynechococcus elongatus</i> BP-1	<b>41</b>	<b>30.5</b>	6	6.6	17	26.9	7.3*
Deinococcus							
<i>Deinococcus radiodurans</i> R1	17	18.1	10	6.0	4	7.0	4.1
Firmicutes							
<i>Bacillus anthracis</i> A0039	28	26.1	8	5.2	6	10.7	3.7
<i>B. anthracis</i> Ames	26	23.1	6	5.1	6	9.8	2.0
<i>B. anthracis</i> Ames Ancestor	26	24.3	7	5.4	7	10.3	1.7
<i>B. anthracis</i> CNEVA-9066	27	24.4	7	5.1	6	10.4	2.8
<i>B. anthracis</i> USA6153	27	25.3	8	5.2	6	10.5	3.6
<i>B. anthracis</i> Vollum	27	24.4	7	5.2	6	10.4	2.8
<i>Bacillus cereus</i> 10987	37	32.3	6	6.6	8	12.2	2.2
<i>B. cereus</i> ATCC 14579	27	28.6	9	5.0	9	11.4	3.8
<i>B. cereus</i> Zk	29	25.6	8	5.0	5	11.4	5.8
<i>Bacillus halodurans</i> C-125	73	74.0	<b>22</b>	<b>13.0</b>	14	22.0	9.1**
<i>Bacillus thuringiensis</i> konkukian	39	38.6	<b>16</b>	<b>7.7</b>	8	16.6	13.4***
<i>Clostridium perfringens</i> SM101	39	39.9	9	7.7	12	12.4	0.2
<i>Desulfobacterium hafniense</i> Y51	66	60.9	6	11.3	18	17.8	2.9
<i>Geobacillus kaustophilus</i> HTA426	<b>65</b>	<b>51.4</b>	10	13.7	8	17.9	10.1**
<i>Staphylococcus epidermidis</i> ATCC 12228	24	34.8	<b>30</b>	<b>9.1</b>	2	12.2	60.3*** <sup>c</sup>
<i>S. epidermidis</i> RP62A	24	31.7	<b>23</b>	<b>9.6</b>	5	10.7	23.6*** <sup>c</sup>
<i>Staphylococcus haemolyticus</i> JCSC1435	36	57.0	<b>44</b>	<b>10.6</b>	5	17.4	121.4*** <sup>c</sup>
<i>Streptococcus pneumoniae</i> G54	39	40.2	11	8.3	8	9.6	1.2
<i>S. pneumoniae</i> R6	30	39.0	<b>15</b>	<b>7.2</b>	9	7.8	10.7**
<i>S. pneumoniae</i> TIGR4	37	43.7	<b>22</b>	<b>11.7</b>	6	9.7	11.6**
<i>Thermoanaerobacter tengcongensis</i> MB4(T)	39	37.5	7	5.8	7	9.8	1.1
Spirochaetes							
<i>Leptospira interrogans</i> lai 56601	34	35.4	17	16.9	15	13.7	0.2

**Table 1**  
Continued

	NGO <sup>a</sup>						$\chi^2$ <sup>b</sup>
	$\rightarrow \rightarrow$ , $\leftarrow \leftarrow$		$\rightarrow \leftarrow$		$\leftarrow \rightarrow$		
	O	E	O	E	O	E	
Unclassified proteobacteria							
<i>Magnetococcus</i> sp. MC-1	41	38.4	16	18.1	15	15.5	0.4
Alphaproteobacteria							
<i>Bradyrhizobium japonicum</i> USDA 110	58	65.3	31	21.7	28	30.0	5.0
<i>Caulobacter crescentus</i> CB15	12	15.7	<b>14</b>	<b>5.0</b>	2	7.4	21.2*** <sup>c</sup>
<i>Gluconobacter oxydans</i> 621H	28	28.7	<b>12</b>	<b>5.8</b>	10	15.5	8.4*
<i>Magnetospirillum magneticum</i> AMB-1	<b>33</b>	<b>24.0</b>	9	8.6	4	13.4	9.9**
<i>Mesorhizobium loti</i> MAFF303099	31	29.9	<b>13</b>	<b>7.1</b>	11	17.9	7.6*
<i>Nitrobacter winogradskyi</i> Nb-255	50	51.4	<b>34</b>	<b>13.8</b>	11	29.7	41.2*** <sup>c</sup>
<i>Rhodopseudomonas palustris</i> BisB18	15	23.4	<b>19</b>	<b>7.7</b>	8	10.9	20.7*** <sup>c</sup>
<i>Rickettsia bellii</i> RML369-C	24	24.0	6	6.5	10	9.5	0.1
<i>Sinorhizobium melliloti</i> 1021	28	32.4	<b>22</b>	<b>10.9</b>	7	13.8	15.3***
<i>Wolbachia pipientis</i> wMel	23	24.7	<b>15</b>	<b>6.6</b>	3	9.7	15.4***
Betaproteobacteria							
<i>Azoarcus</i> sp. EbN1	63	61.4	<b>29</b>	<b>16.8</b>	11	24.8	16.7*** <sup>c</sup>
<i>Bordetella pertussis</i> Tohama I	68	79.9	<b>52</b>	<b>13.1</b>	15	42.0	134.2*** <sup>c</sup>
<i>Burkholderia cenocepacia</i> AU 1054	30	40.5	<b>23</b>	<b>9.3</b>	18	21.2	23.3*** <sup>c</sup>
<i>Burkholderia mallei</i> ATCC 23344	46	56.6	<b>39</b>	<b>20.7</b>	15	22.7	20.8*** <sup>c</sup>
<i>Burkholderia pseudomallei</i> 1710b	31	33.9	22	15.4	9	12.7	4.1
<i>B. pseudomallei</i> K96243	26	29.8	<b>22</b>	<b>10.2</b>	6	13.9	18.6*** <sup>c</sup>
<i>Burkholderia thailandensis</i> E264	40	38.8	<b>21</b>	<b>13.6</b>	9	17.6	8.2*
<i>Burkholderia</i> sp. 383	20	20.0	10	5.6	7	11.4	5.2
<i>Neisseria meningitidis</i> MC58	21	29.3	<b>21</b>	<b>10.2</b>	8	10.5	14.2***
<i>N. meningitidis</i> Z2491	14	22.3	<b>17</b>	<b>8.1</b>	8	8.6	12.8**
<i>Nitrosomonas europaea</i> ATCC 19718	37	52.6	<b>30</b>	<b>9.7</b>	19	23.7	48.4*** <sup>c</sup>
<i>Nitrosospora multiformis</i> ATCC 25196	32	32.1	<b>15</b>	<b>7.8</b>	7	14.0	10.1**
<i>Ralstonia solanacearum</i> GMI1000	21	24.0	<b>11</b>	<b>5.2</b>	8	10.8	7.7*
Deltaproteobacteria							
<i>Desulfovibrio desulfuricans</i> G20	33	28.7	12	10.8	5	10.5	3.6
<i>Geobacter metallireducens</i> GS-15	48	50.0	<b>16</b>	<b>8.9</b>	14	19.0	7.0*
<i>Myxococcus xanthus</i> DK 1622	23	21.8	<b>10</b>	<b>5.5</b>	7	12.7	6.2*
<i>Pelobacter carbinolicus</i> DSM 2380	20	26.5	<b>12</b>	<b>5.5</b>	10	10.0	9.4**
Gammaproteobacteria							
<i>Acidithiobacillus ferrooxidans</i> ATCC 23270	33	34.4	<b>14</b>	<b>8.0</b>	8	12.6	6.3*
<i>Coxiella burnetii</i> RSA 493	17	16.9	3	5.0	10	8.1	1.2
<i>Escherichia coli</i> CFT073	44	44.1	<b>26</b>	<b>12.5</b>	11	24.4	22.0*** <sup>c</sup>
<i>E. coli</i> K12 MG1655	33	34.1	<b>18</b>	<b>8.0</b>	9	17.9	17.0*** <sup>c</sup>
<i>E. coli</i> O157:H7 EDL933	24	28.6	<b>16</b>	<b>5.7</b>	6	11.7	22.3*** <sup>c</sup>
<i>E. coli</i> O157:H7 VT2-Sakai	38	44.1	<b>26</b>	<b>9.0</b>	8	18.9	39.0*** <sup>c</sup>
<i>E. coli</i> UT189	19	23.5	<b>17</b>	<b>5.6</b>	5	11.9	27.8*** <sup>c</sup>
<i>Francisella tularensis</i> holarctica	60	56.5	<b>22</b>	<b>11.2</b>	16	30.3	17.2*** <sup>c</sup>
<i>F. tularensis</i> tularensis	28	25.6	<b>14</b>	<b>5.8</b>	5	15.6	18.9*** <sup>c</sup>
<i>Hahella chejuensis</i> KCTC 2396	25	22.0	5	5.4	8	10.6	1.1
<i>Legionella pneumophila</i> Paris	19	22.5	12	6.4	12	14.1	5.8
<i>Methylococcus capsulatus</i> Bath	16	17.2	10	5.7	6	9.1	4.3
<i>Nitrosococcus oceani</i> ATCC 19707	48	52.3	17	11.5	23	24.2	3.0
<i>Photobacterium profundum</i> SS9	121	109.1	38	30.3	35	54.6	10.3**
<i>Photorhabdus luminescens</i> TTO1	79	71.0	<b>28</b>	<b>16.7</b>	14	33.3	19.7*** <sup>c</sup>
<i>Pseudomonas putida</i> KT2440	29	31.6	<b>19</b>	<b>10.4</b>	9	15.1	9.8**
<i>Pseudomonas syringae</i> DC3000	61	66.7	<b>36</b>	<b>18.0</b>	24	36.2	22.5*** <sup>c</sup>
<i>P. syringae</i> pv B728a	16	21.2	<b>15</b>	<b>6.5</b>	8	11.3	13.6***
<i>P. syringae</i> pv phaseolicola	58	57.6	<b>28</b>	<b>16.7</b>	19	30.7	12.1**
<i>Psychrobacter arcticum</i> 273-4	20	28.7	<b>16</b>	<b>5.9</b>	12	13.3	19.9*** <sup>c</sup>
<i>Salmonella enterica</i> Choleraesuis	18	28.4	<b>19</b>	<b>7.7</b>	13	13.9	20.4*** <sup>c</sup>

**Table 1**  
Continued

	NGO <sup>a</sup>						$\chi^{2b}$
	$\rightarrow \rightarrow$ , $\leftarrow \leftarrow$		$\rightarrow \leftarrow$		$\leftarrow \rightarrow$		
	O	E	O	E	O	E	
<i>Shewanella oneidensis</i> MR-1	68	73.9	<b>33</b>	<b>21.9</b>	33	38.2	6.7*
<i>Shigella boydii</i> Sb227	100	114.6	<b>55</b>	<b>24.0</b>	48	64.4	46.1*** <sup>c</sup>
<i>Shigella dysenteriae</i> Sd197	156	177.4	<b>72</b>	<b>33.0</b>	78	95.6	51.9*** <sup>c</sup>
<i>Shigella flexneri</i> 2a 301	116	113.9	<b>51</b>	<b>38.3</b>	34	48.9	8.8*
<i>S. flexneri</i> 2a 2457T	60	61.1	<b>27</b>	<b>11.1</b>	20	34.8	28.9*** <sup>c</sup>
<i>Shigella sonnei</i> Ss046	103	100.8	<b>37</b>	<b>23.1</b>	36	52.1	13.3***
<i>Sodalis glossinidius</i> morsitans	12	18.3	7	6.8	<b>14</b>	<b>7.9</b>	6.9*
<i>Vibrio cholerae</i> El Tor N16961	15	12.7	6	5.0	3	6.3	2.3
<i>Vibrio vulnificus</i> YJ016	24	25.2	<b>15</b>	<b>6.9</b>	6	13.0	13.4***
<i>Xanthomonas axonopodis</i> pv. citri 306	28	28.2	11	8.9	8	9.9	0.9
<i>Xanthomonas campestris</i> 8004	25	28.5	16	10.8	9	10.7	3.3
<i>X. campestris</i> ATCC 33913	37	36.4	16	13.6	10	13.0	1.1
<i>X. campestris</i> pv. armoraciae 756C	24	22.6	14	12.5	4	6.9	1.5
<i>X. campestris</i> pv. vesicatoria 85-10	33	34.4	12	11.5	14	13.1	0.1
<i>Xanthomonas oryzae</i> KACC10331	179	189.6	93	77.2	44	49.2	4.3
<i>X. oryzae</i> pv. oryzae MAFF 311018	155	170.2	<b>91</b>	<b>58.5</b>	45	62.3	24.2*** <sup>c</sup>
<i>X. oryzae</i> pv. oryzicola BLS256	98	95.9	62	57.2	23	29.9	2.0
<i>Yersinia pestis</i> biovar Medievalis 91001	37	38.4	<b>29</b>	<b>10.3</b>	2	19.3	49.3*** <sup>c</sup>
<i>Y. pestis</i> CO92	44	48.7	<b>35</b>	<b>13.0</b>	7	24.3	50.0*** <sup>c</sup>
<i>Y. pestis</i> KIM	57	62.7	<b>45</b>	<b>19.9</b>	11	30.3	44.4*** <sup>c</sup>
<i>Yersinia pseudotuberculosis</i> IP32593	17	19.4	<b>12</b>	<b>5.0</b>	5	9.6	12.3**

<sup>a</sup> NGOs in bold contribute a significant excess of observed ISs to significant  $\chi^2$  deviations, and those in gray contribute a significant deficit of observed ISs.

<sup>b</sup> Asterisks indicate significant *P* values: \**P* ≤ 0.05, \*\**P* ≤ 0.01, \*\*\**P* ≤ 0.001.

<sup>c</sup> *P* value is significant following a sequential Bonferroni (Rice 1989).

remarkable consistency with which intergenic IS elements are overrepresented and underrepresented between  $\rightarrow \leftarrow$  and  $\leftarrow \rightarrow$  oriented genes, respectively (fig. 1), suggests that these are generally relatively innocuous and deleterious insertion locations, thus supporting the hypothesis that differential selection pressure molds global intergenic IS distributions. Further fine-scale analyses of intergenic IS distributions (e.g., ISs may be less common between  $\rightarrow \rightarrow$  and  $\leftarrow \leftarrow$  neighbors when they are members of the same operon; ISs may be relatively rare next to highly expressed genes, no matter what their orientation) may shed additional light on the fate and impact of IS elements in bacterial genomes.

## Materials and Methods

I obtained the primary annotations of all fully sequenced bacterial chromosomes from the Comprehensive Microbial Resource database (data releases 1.0–20.0) at The Institute for Genomic Research (<http://cmr.tigr.org/tigr-scripts/CMR/CmrHomePage.cgi>). Specifically, I obtained the locus name (i.e., the locus number), the common name, the nucleotide sequence, and the nucleotide positions of the 5' and 3' ends of all annotated proteins on each chromosome. My goal for each genome was to assess whether the observed quantities of intergenic IS elements located within each of the

three NGOs differ from the quantities expected if insertion is random and not subsequently influenced by natural selection. This required four steps for each fully sequenced genome.

The first step was to find all chromosomal copies of intergenic IS elements. I used the BlastX program in the IS-finder database (<http://www-is.biotoul.fr/is.html>) (Siguier et al. 2006) to identify all coding sequences (CDSs) in each genome that exhibit homology to IS elements in the database. I considered a CDS with a best BlastX hit *E* value ≤ 10<sup>-10</sup> to be an IS element (Touchon and Rocha 2007). Because I was only interested in the distribution of ISs between functional native bacterial genes, I took a relatively conservative approach when identifying intergenic IS elements (i.e., it is better to exclude some intergenic ISs than to include any intragenic ISs). Specifically, I eliminated the following IS elements from the analysis: 1) all intragenic ISs, including elements with at least one neighboring gene annotated as being truncated (or similar synonyms), conservatively assuming that the neighboring gene became degenerate following IS insertion into the gene; 2) all ISs bordered by genes with annotated frameshift or point mutations that introduce premature stop codons, conservatively assuming that these mutations preceded IS insertion; that is, the IS was never exposed to selection from two functional neighboring genes; 3) all ISs bordered by

nonconsecutively numbered and therefore presumably nonneighboring genes (e.g., some are bordered by nonannotated gene remnants, which may have become degenerate following IS insertion); and 4) all ISs bordered by a phage-annotated gene, and those annotated as being or bordering an integron or an integrative genetic element (for the quantities of ISs eliminated for each of these reasons in each genome, see [supplementary table S2, Supplementary Material](#) online). Conversely, I included IS elements with both functional and nonfunctional transposases because ISs can affect their neighboring genes even if they are no longer mobile (e.g., by displacing promoters). Also, multiple IS insertions into the same intergenic space were included only once in the analysis.

The second step was to calculate the observed quantity of intergenic IS elements within each NGO (i.e., assessing whether the two neighboring genes are coded on the top or bottom DNA strand for each IS element). I did this by simply subtracting the nucleotide position of the 5' end from that of the 3' end for each neighbor, which produces a positive number for top strand genes and a negative number for bottom strand genes.

The third step was to calculate the expected quantity of intergenic IS elements within each NGO, assuming that IS insertion is random and not subsequently affected by natural selection. I calculated these expected quantities based on the premise that large and abundant NGO intergenic regions should receive more ISs than small and rare ones, all things being equal. Therefore, the expected quantities were calculated individually for each genome using the product of 1) the mean intergenic distance between neighboring native bacterial genes in the three NGOs and 2) the global proportion of each native gene pair NGO; for an example of this calculation, see table S3 ([Supplementary Material](#) online).

Finally, the fourth step was to use a  $\chi^2$  goodness-of-fit test to assess whether the observed quantities of intergenic IS elements within each NGO deviate from the expected quantities. The assumptions of the  $\chi^2$  test are that no cell has an expected value  $<1.0$  and that  $\leq 20\%$  of cells have expected values  $<5.0$  (Cochran 1954). Therefore, many fully sequenced genomes do not contain enough intergenic IS elements for statistical analysis (all 116 genomes with enough intergenic ISs are included in table 1, and the remaining 210 genomes are included in table S1, [Supplementary Material](#) online). I did not Bonferroni-adjust the  $\chi^2$  test *P* values (Moran 2003), although all  $\chi^2$  values that would be significant with a Bonferroni correction are indicated in table 1. To identify the NGOs contributing to each significant  $\chi^2$  deviation, I performed cell-by-cell comparisons of observed and expected quantities using an adjusted residual method, considering any adjusted residual with an absolute value  $>2$  to contribute significantly to the overall  $\chi^2$  deviation (Agresti 1996).

## Supplementary Material

Supplementary tables S1–S3 are available at *Genome Biology and Evolution* online ([http://www.oxfordjournals.org/our\\_journals/gbe/](http://www.oxfordjournals.org/our_journals/gbe/)).

## Acknowledgments

I thank Huansheng Cao, Kevin Dougherty, Evelyn Fetridge, Catherine Ruggiero, Chad Thompson, and several anonymous reviewers for helpful comments on this manuscript, and Elizabeth Coffey for early contributions to this project. This work was supported by the National Institutes of Health (grant number 1R15GM081862-01A1). This is contribution number 246 of the Louis Calder Center—Biological Field Station, Fordham University.

## Literature Cited

- Agresti A. 1996. An introduction to categorical data analysis. New York: John Wiley and Sons.
- Campbell A. 2002. Eubacterial genomes. In: Craig NL, Craigie R, Gellert M, Lambowitz A, editors. Mobile DNA II. Washington (DC): ASM Press. pp. 1024–1039.
- Chandler M, Mahillon J. 2002. Insertion sequences revisited. In: Craig NL, Craigie R, Gellert M, Lambowitz A, editors. Mobile DNA II. Washington (DC): ASM Press. p. 305–366.
- Cochran WG. 1954. Some methods for strengthening the common  $\chi^2$  test. *Biometrics*. 10:417–451.
- Cooper VS, Schneider D, Blot M, Lenski RE. 2001. Mechanisms causing rapid and parallel losses of ribose catabolism in evolving populations of *Escherichia coli* B. *J Bacteriol*. 183:2834–2841.
- Doolittle WF, Sapienza C. 1980. Selfish genes, the phenotype paradigm and genome evolution. *Nature*. 284:601–603.
- Doran T, et al. 1997. IS900 targets translation initiation signals in *Mycobacterium avium* subsp. *paratuberculosis* to facilitate expression of its *hed* gene. *Microbiology*. 143:547–552.
- Inglis NF, Stevenson K, Heaslip DG, Sharp JM. 2003. Characterisation of IS901 integration sites in the *Mycobacterium avium* genome. *FEMS Microbiol Lett*. 221:39–47.
- Liu LF, Wang JC. 1987. Supercoiling of the DNA template during transcription. *Proc Natl Acad Sci U S A*. 84:7024–7027.
- Lodge JK, Berg DE. 1990. Mutations that affect Tn5 insertion into pBR322: importance of local DNA supercoiling. *J Bacteriol*. 172:5956–5960.
- Lynch M. 2006. Streamlining and simplification of microbial genome architecture. *Annu Rev Microbiol*. 60:327–349.
- Moran MD. 2003. Arguments for rejecting the sequential Bonferroni in ecological studies. *Oikos*. 100:403–405.
- Nuzhdin SV. 1999. Sure facts, speculations, and open questions about the evolution of transposable element copy number. *Genetica*. 107:129–137.
- Orgel LE, Crick FHC. 1980. Selfish DNA: the ultimate parasite. *Nature*. 284:604–607.
- Rice WR. 1989. Analyzing tables of statistical tests. *Evolution*. 43:223–225.
- Safi H, et al. 2004. IS6110 functions as a mobile, monocyte-activated promoter in *Mycobacterium tuberculosis*. *Mol Microbiol*. 52:999–1012.

- Siguier P, Perochon J, Lestrade L, Mahillon J, Chandler M. 2006. ISfinder: the reference centre for bacterial insertion sequences. *Nucleic Acids Res.* 34:D32–D36.
- Tetu SG, Holmes AJ. 2008. A family of insertion sequences that impacts integrons by specific targeting of gene cassette recombination sites, the IS1111-attC group. *J Bacteriol.* 190:4959–4970.
- Touchon M, Rocha EPC. 2007. Causes of insertion sequences abundance in prokaryotic genomes. *Mol Biol Evol.* 24:969–981.
- Wagner A. 2006. Periodic extinctions of transposable elements in bacterial lineages: evidence from intragenomic variation in multiple genomes. *Mol Biol Evol.* 23:723–733.
- Zaghloul L, et al. 2007. The distribution of insertion sequences in the genome of *Shigella flexneri* strain 2457T. *FEMS Microbiol Lett.* 277:197–204.

**Associate editor:** Laurence Hurst