

# The Evolution of Insertion Sequences Within Enteric Bacteria

Jeffrey G. Lawrence,<sup>1</sup> Howard Ochman and Daniel L. Hartl<sup>2</sup>

Department of Genetics, Washington University School of Medicine, St. Louis, Missouri 63110

Manuscript received October 10, 1991

Accepted for publication January 8, 1992

## ABSTRACT

To identify mechanisms that influence the evolution of bacterial transposons, DNA sequence variation was evaluated among homologs of insertion sequences IS1, IS3 and IS30 from natural strains of *Escherichia coli* and related enteric bacteria. The nucleotide sequences within each class of IS were highly conserved among *E. coli* strains, over 99.7% similar to a consensus sequence. When compared to the range of nucleotide divergence among chromosomal genes, these data indicate high turnover and rapid movement of the transposons among clonal lineages of *E. coli*. In addition, length polymorphism among IS appears to be far less frequent than in eukaryotic transposons, indicating that nonfunctional elements comprise a smaller fraction of bacterial transposon populations than found in eukaryotes. IS present in other species of enteric bacteria are substantially divergent from *E. coli* elements, indicating that IS are mobilized among bacterial species at a reduced rate. However, homologs of IS1 and IS3 from diverse species provide evidence that recombination events and horizontal transfer of IS among species have both played major roles in the evolution of these elements. IS3 elements from *E. coli* and *Shigella* show multiple, nested, intragenic recombinations with a distantly related transposon, and IS1 homologs from diverse taxa reveal a mosaic structure indicative of multiple recombination and horizontal transfer events.

**I**NSERTION sequences (IS) are short segments of DNA, typically 1–2 kilobasepairs (kb) in length, with the ability to translocate within and among replicons (GALAS and CHANDLER 1989). IS mediate numerous molecular and genetic phenomena, including gene activation (GLANDSDORFF, CHARLIER and ZAFARULLA 1980), repression (SAEDLER *et al.* 1974), deletion (CHOW and BROKER 1981), rearrangement (SAEDLER *et al.* 1980), recombination (LIAB 1980), and transfer (CHANDLER, CLERGET and CARO 1980), and are also of epidemiological importance, owing to their ability to form composite transposons and mobilize antibiotic resistance determinants (BERG 1977; KLECKNER *et al.* 1975). Although numerous classes of IS have been characterized from the genomes of enteric bacteria, factors that contribute to the evolution of IS within and among bacterial species have not been clearly defined. Elucidating the evolutionary forces acting upon insertion elements is necessary for understanding the genetics of mobile genetic elements in prokaryotes, the impact of IS-mediated events, and the evolution and epidemiology of composite transposons.

Bacterial reproduction confers a clonal population structure if gene exchange is infrequent. Although recent studies have shown that recombination does occur at chromosomal loci among natural isolates of

*E. coli* (DUBOSE, DYKHUIZEN and HARTL 1988; STOLTZFUS, LESLIE and MILKMAN 1988), linkage disequilibrium among enzyme electrophoretic types has indicated that large scale genetic exchange of metabolic genes among natural strains of *E. coli* is not common (SELANDER and LEVIN 1980; CAUGANT, LEVIN and SELANDER 1981; OCHMAN and SELANDER 1984a; WHITTAM, OCHMAN and SELANDER 1984). However, the contribution of genetic exchange to the evolution of transposable elements within *E. coli* has not been well characterized. Not only are IS present in variable numbers and positions within bacterial genomes (SAWYER *et al.* 1987; LAWRENCE *et al.* 1989), but a substantial portion of the *E. coli* transposon pool is plasmid borne (SAWYER *et al.* 1987; HALL *et al.* 1989), which increases opportunities for genetic transfer.

SAWYER *et al.* (1987) examined the distribution and abundance of IS1, IS2, IS3, IS4, IS5 and IS30, and HALL *et al.* (1989) the distribution of IS103, in natural isolates of *E. coli*, chosen to represent the range of phenotypic variation detected by protein electrophoresis. They determined that IS could be classified into three groups based on the apparent strength of regulation of transposition: IS1 and IS5 represented weakly regulated transposons; IS2, IS4, IS30 and IS103 represented moderately regulated transposons; and IS3 represented a class of strongly regulated transposons (see also HARTL and SAWYER 1988a,b). Although these analyses suggest different patterns of

<sup>1</sup> Current address: Department of Biology, University of Utah, Salt Lake City, Utah 84112.

<sup>2</sup> To whom correspondence should be addressed.

regulation, they did not establish which factors govern the evolution of insertion sequences. In addition, the degree to which these factors differentially influence the evolution of each class of transposons is not known.

To address these issues, isoforms of three insertion sequences, *IS1*, *IS3* and *IS30*, were studied from isolates of the ECOR reference collection of *E. coli* (OCHMAN and SELANDER 1984b). In addition, IS homologs from related species of enteric bacteria were isolated and analyzed. Nucleotide variation among classes of IS was examined to elucidate the mechanisms influencing the evolution of IS within *E. coli*, to evaluate the extent to which genetic exchange contributes to the evolution of transposon populations with *E. coli*, to determine if similar processes influence the evolution of IS among bacterial species, to assess the relative proportions of functional and non-functional elements, to identify evolutionarily conserved reading frames and sequence motifs, and to ascertain whether patterns of evolution differ among distinct classes of insertion sequences. In this manner, we may assess the evolutionary influences which differ between bacterial IS and metabolic genes, as well as those which differ between prokaryotic and eukaryotic transposons.

#### MATERIALS AND METHODS

**Strains:** Strains of the ECOR collection (OCHMAN and SELANDER 1984b), ATCC 35320-35391, were obtained from laboratory collections. The phenotypic (SELANDER, CAUGANT and WHITTAM 1987) and insertion sequence (SAWYER *et al.* 1987) profiles of these strains have been described. *Escherichia fergusonii* ATCC 35469 and ATCC 35471, *Escherichia hermannii* ATCC 33652, *Escherichia vulneris* ATCC 29943, *Shigella dysenteriae* ATCC 13313, *Shigella flexneri* ATCC 29508, *Shigella sonnei* ATCC 29930 and *Serratia odorifera* ATCC 3307 were obtained from laboratory stocks. Relationships among these taxa inferred from chromosomal gene sequences have been described (LAWRENCE, OCHMAN and HARTL 1991).

**Southern blotting:** Chromosomal DNA was isolated, digested with restriction endonucleases, size fractionated on agarose gels and transferred to nylon membranes as described previously (SAWYER *et al.* 1987). Probes were prepared as internal portions of IS amplified via the polymerase chain reaction (PCR; SAIKI *et al.* 1985, 1988). Primers utilized for amplification annealed internal to the inverted repeats: *IS1* forward: GATTTAGTGTATGATGG; *IS1* reverse: GATAGTGTTTTATGTTTC; *IS3* forward: GGA-CACGCGGCTAAGTG; *IS3* reverse: TGGACACAGGC-CTAAGCG; *IS30* forward: GCAACAGTTATGTGAAA; *IS30* reverse: AATGCAACACCCCTTTC. Amplification products were purified by the method of LAWRENCE, HARTL and OCHMAN (1991a), and labeled to high specific activity by the method of FEINBERG and VOGELSTEIN (1983). Membranes were hybridized under moderate stringency conditions as described (SAWYER *et al.* 1987).

**DNA sequencing:** Internal segments of each IS were amplified by the PCR utilizing either genomic DNA, for strains containing a single copy of any one IS, or DNA fragments, size fractionated by gel electrophoresis and iden-

TABLE 1  
Single nucleotide differences within *IS1* copies from *E. coli*

Position	Strain <sup>a</sup>											
	K12	5	28	48	12	32	33	51	52	53	66	60
393	A	C	-	-	C	C	C	C	C	C	C	C
396	C	T	-	-	-	-	-	-	-	-	-	-
418	C	-	-	-	-	-	-	-	A	-	-	-
486	C	-	-	-	-	-	-	G	-	-	-	-

<sup>a</sup> K12, *E. coli* K12 (OHTSUBO and OHTSUBO 1978); numbers refer to ECOR strains (OCHMAN and SELANDER 1984b). Strains 5, 28, and 48 are isoforms of *IS1R*; the remaining strains are isoforms of *IS1F*. Positions which distinguish *IS1F* from *IS1R* were not included (see text). The dash indicates a nucleotide identical to that in K12.

tified by Southern blotting, as a template. The oligonucleotide primers described above allowed amplification of a 729-bp fragment comprising 95% of *IS1*, a 1220-bp fragment encompassing 97% of *IS3*, and an 1185-bp fragment comprising 97% of *IS30*. Amplification products were purified and sequenced according to the methods of DUBOSE and HARTL (1990) and LAWRENCE, HARTL and OCHMAN (1991a). For copies of *IS3* and *IS30*, the nucleotide sequence of the entire region between the amplification primers was determined. For *IS1* resident in *E. coli*, nucleotide sequences downstream of position 400 were determined. In all cases, the nucleotide sequences of both strands of each IS were determined. For genomes containing elements that could not be amplified by the PCR, chromosomal DNA was partially digested with *Sau3A*, ligated to EMBL3 digested with *Bam*HI, and packaged in Gigapack (Stratagene). Bacteriophage DNA was prepared from appropriate clones by the method of HELMS *et al.* (1985), partially digested with *Sau3A*, ligated into M13 vectors digested with *Bam*HI, and introduced into *E. coli* JM101 by electroporation (DOWER, MILLER and RAGSDALE 1988). Appropriate clones were selected for DNA sequencing.

**Computer analysis:** Parsimony analysis was implemented by PAUP (D. SWOFFORD), and divergence calculations utilized the GCG program package (DEVEREUX, HAEBERLI and SMITHIES 1984). Phylogeny testing employed the program MATRIX2 (LAWRENCE and HARTL 1992).

#### RESULTS

**Nucleotide substitutions:** The DNA sequences of insertion elements *IS1*, *IS3* and *IS30* revealed little nucleotide variation in the form of base substitutions among strains of *E. coli* (Tables 1-3). Both types of *IS1* were detected: *IS1R*, originally isolated from plasmid R100 (OHTSUBO and OHTSUBO 1978), and *IS1F*, first isolated from *S. flexneri* (OHTSUBO *et al.* 1984). Although they differ by 10% at the nucleotide level, little variation was detected within each isoform. Similarly, *IS3* and *IS30* were also virtually monomorphic, despite their presence in bacterial strains that are quite distinct based on enzyme electrophoretic profiles (Figure 1). Figure 1 also indicates strains containing multiple copies of either *IS1F* or *IS1R*. The isoforms were distinguished by (1) restriction fragment length polymorphism analysis of PCR-amplified regions of *IS1* cleaved with restriction endonucleases that distinguish

TABLE 2

Single nucleotide differences within IS3 copies from *E. coli*

Position	Strain <sup>a</sup>											
	K12	1	14	18	23	30	43	46	50	58	69	70
38	A	-	-	-	-	-	-	-	-	T	-	-
80	A	-	-	-	-	-	G	-	-	-	-	-
200	A	-	G	-	G	G	G	G	G	G	G	-
250	G	-	-	-	-	-	A	-	-	-	-	-
483	C	-	-	-	-	-	T	-	-	-	-	-
485	C	-	-	A	-	-	-	-	-	-	-	-
716	G	-	-	-	-	-	-	-	-	T	-	-
1019	G	A	-	-	-	-	-	-	-	-	-	-
1125	C	-	-	-	-	-	-	-	T	-	-	-
1214	A	-	C	-	-	-	C	-	-	-	C	-

<sup>a</sup> K12, *E. coli* K12 (TIMMERMAN and TU 1985); numbers refer to ECOR strains (OCHMAN and SELANDER 1984b). The nucleotide sequence of ECOR 63 appears in Figure 3. The dash indicates a nucleotide identical to that in K12.

TABLE 3

Single nucleotide differences within IS30 copies from *E. coli*

Position	Strain <sup>a</sup>												
	K12	1	8	14	19	23	24	31	35	36	50	56	71
74	T	A	A	A	A	A	A	A	A	A	A	A	A
75	A	T	T	T	T	T	T	T	T	T	T	T	T
141	A	G	G	G	G	G	G	G	G	G	G	G	G
217	G	-	A	-	A	-	-	-	-	-	-	A	-
220	A	C	-	-	-	-	-	-	-	-	-	-	-
236	C	-	-	-	-	-	-	-	-	-	-	T	-
266	A	-	-	-	-	-	-	-	-	-	T	-	T
276	T	-	-	-	-	-	-	-	-	-	C	-	-
461	G	-	-	-	-	-	T	-	-	-	-	-	-
479	G	-	-	-	-	-	-	-	-	-	-	-	T
486	A	-	-	-	-	-	-	C	C	-	-	-	-
594	C	-	-	-	A	-	-	-	-	-	-	-	-
755	A	G	G	G	G	G	G	G	G	-	G	G	G
758	G	-	-	-	-	-	-	A	-	-	A	-	-
899	C	-	-	-	-	-	-	-	-	-	A	-	-
983	C	-	-	-	-	T	-	-	-	-	-	-	-
1046	C	-	-	-	T	-	-	-	-	-	-	-	-

<sup>a</sup> K12, *E. coli* K12 (DALRYMPLE *et al.* 1984); numbers refer to ECOR strains (OCHMAN and SELANDER 1984b). The dash indicates a nucleotide identical to that in K12.

IS1R from IS1F (strains 2, 24, 29, 42, 55, 57, 59 and 72), and (2) DNA sequencing of a 280-bp region in particular IS1 copies that included 29 sites distinguishing the two isoforms (strains 3, 6 and 30). The distributions of IS1F and IS1R are not confined to any one group of *E. coli*, but are widely distributed among divergent strains. For example, although strains 30, 32 and 33 are quite closely related, strain 30 contains IS1R, while strains 32 and 33 contain IS1F. We have not identified any strain containing both IS1R and IS1F. The multiple copies of IS1 in strains 2, 3, 6, 24, 29, 30, 42 and 72 are isoforms of IS1R, while those in strains 55, 57 and 59 are IS1F. However, the laboratory strain *E. coli* K12 harbors five IS1R and one IS1F (UMEDA and OHTSUBO 1991).

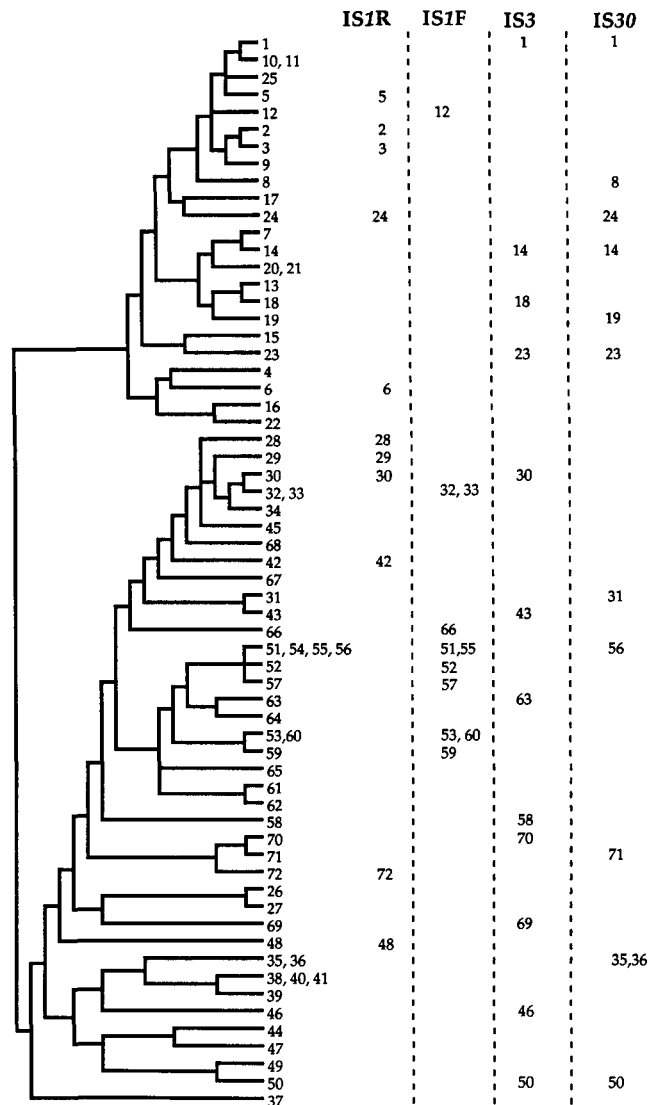


FIGURE 1.—Phylogenetic relationships of *E. coli* strains from which IS were isolated and analyzed [tree after SELANDER, CAUGANT and WHITTAM (1987)]. Strains from which copies of IS1R, IS1F, IS3 and IS30 were isolated are indicated.

Figure 2 shows the nucleotide sequences of IS1 elements resident in the genomes of related enteric bacteria. IS1-related sequences were isolated from three species of *Shigella* (*S. dysenteriae*, *S. flexneri*, and *S. sonnei*) and their nucleotide sequences agree well with IS1 sequences previously published from these taxa (OHTSUBO *et al.* 1984), differing by only single substitutions. IS1 homologs were also isolated from *E. fergusonii*, *E. hermannii* and *E. vulneris*. [Although classified as *Escherichia*, *E. hermannii* and *E. vulneris* are as divergent from *E. coli* as species of *Enterobacter* (LAWRENCE, HARTL and OCHMAN 1991b; LAWRENCE, OCHMAN and HARTL 1991).] Relationships among IS1 homologs isolated from these species are presented in Table 4. The relationships among the elements are not congruent with those inferred from chromosomal loci, notably *gap* and *ompA*, which encode glucose-3-





1 *Eco* TGTAGATCCA ATTGGTCAAC GCAACAGTTA TGTGAAAACA TGGGGTTGCG GAGGTTTTTTT  
*Ehe* -----A-----C-----TA-----A-G--

61 *Eco* GAATGAGACG AACTATTACA GCAGAGGAAA AAGCCTCTGT TTTTGAAC TAAGAAGAAG  
*Ehe* ---AT-----A-----A-----G-----

121 *Eco* GAACAGGCTT CAGTGAATA ACGAATATCC TGGGTTCAAA ACCGGGAACG ATCTTCACTA  
*Ehe* --A-----G---G-T-----G---A--G--C-----G-

181 *Eco* TGTTAAGGGA TACTGGCGGC ATAAACCCC ATGAGCGTAA GCGGGTGTGA GCTCACCTGA  
*Ehe* -----A-- -----A-A-----C-----A-----T-----

241 *Eco* CACTGTCTGA GCGCGAGGAG ATACAGCTG GTTTGTGAGC CAAAATGAGC ATTCGTGCGA  
*Ehe* -G-----A-----C-----C-----T-----C-----

301 *Eco* TAGTACTGCG GCTGAATCGC AGTCCTTCCA CGATCTCACG TGAAGTTCAG CGTAATCGGG  
*Ehe* -C-----AAT-----C-----A-----T-----

361 *Eco* GCAGACGCTA TTACAAGCTT GTTGATGCTA ATAACCGGAC CAACAGAATG GCGAAAAGGC  
*Ehe* -T--G-----C-----G-----G-----G-----TC---A-

421 *Eco* CAAAACCGTG CTTACTGGAT CAAAATTTAC CATTGGGAAA GCTTGTCTG GAAAAGCTGG  
*Ehe* -----A--G--G-----G-----G-----T-----

481 *Eco* AGATGAAATG GTCTCCAGAG CAAATATCAG GATGGTTAAG GCGAACAAAA CCACGTCAAA  
*Ehe* -----G-- -----C-----A-----G-----G-----G-----

541 *Eco* AAACGCTGCG AATATCACCT GAGACAATTT ATAAACGCT GTACTTTCGT AGCGCTGAAG  
*Ehe* -----A--A-----C-----A-----T-----T-----

601 *Eco* CGCTACACCA CCTGAATATA CAGCATCTGC GACGGTGGCA TAGCCTTCGC CATGGCAGGC  
*Ehe* -----G-- -----CG-----C-----A-----C-----C-----T-----

661 *Eco* GTCATACCCG CAAAGGGCAA AGAGGTAGCA TTAACATAGT GAACGGAAAC CCAATTCAGG  
*Ehe* -C-----T-----T-----C-----T-----C-----G-----C-----

721 *Eco* AACGTTCCCG AAATATCGAT AACAGACGCT CTCTAGGGCA TTGGGAGGCG GATTTAGTCT  
*Ehe* -----G-- -----C-----G-----A-----C-----A-----T-----

781 *Eco* CAGGTACAAA AAATCTCAT ATAGCCACAC TTGTAGACG AAAATCACGT TATACGATCA  
*Ehe* -G-----C-----C-----G-----A-----A-----T-----

841 *Eco* TCCITGACT CAGGGGCAAA GATTCTGTCT CAGTAAATCA GGCTCTTACC GACAAATTC  
*Ehe* -----G-- -----G-----A-----G-----A-----A-----T-----

901 *Eco* TGAGTTTACC GTCAGAATC AGAAAATCAC TGACATGGGA CAGAGGAATG GAANTGGCCA  
*Ehe* -----C-----AC-T-----C-GCG-----G-----

961 *Eco* GACATCTAGA ATTACTGTC AGCACCGCGG TTAAGTTTA CTCTCGGAT CCTCAGAGTC  
*Ehe* -----G-- -----A-----A-----T-----G-----C-----

1021 *Eco* CTTGGCAGCG GGGAAACAAAT GAGAACAACA ATGGGCTAAT TCGGCAGTAC TTTCTAAAA  
*Ehe* -C-----C-----A-----T-----C-----A-----C-----

1081 *Eco* AGACATGTCT TGCCCAATAT ACTCAACATG AACTAGATCT GGTGTCTGCT CAGCTAAACA  
*Ehe* -----C-----T-----G-----G-----A-----A-----A-----

1141 *Eco* ACAGACCGAG AAAGACACTG AAGTTCAAAA CACCGAAAGA GATAATTGAA AGGGGTGTTG  
*Ehe* -----G-----A-----T-----T-----A-----A-----

1201 *Eco* CATTGACAGA TTGAATCTAC A  
*Ehe* -G-----T-----

FIGURE 4.—Sequences of IS30 from enteric bacteria. *Eco*, *E. coli* K12 (DALRYMPLE, CASPERS and ARBER 1984); *Ehe*, *E. hermannii*. Numbering begins at the first base of the inverted repeat.

ships inferred from these two ORFs are significantly different ( $P < 0.001$ ; Table 5).

The single IS30 isolated from *E. hermannii* was 89.5% identical to the *E. coli* K12 element (Figure 4). Southern blot analysis revealed the presence of sequences related to IS30 in the genomes of *Escherichia blattae* and *E. fergusonii*; however, the relative hybridization intensity indicated that either the IS30 sequence was much more divergent than that resident in the *E. hermannii* genome, or only a small fragment of the element remained in that genome (data not shown).

**Length variation:** In addition to the nucleotide substitutions and recombinations described above, several insertions and deletions were observed within various IS (Table 6). The IS30 from ECOR 24 carried an insertion of the unrelated transposon IS3411. In addition, the IS3 from ECOR 18 carried a 4-bp duplication which may have arisen as a target duplication from transposon insertion and subsequent excision. In total, five of the 35 elements examined from *E. coli*

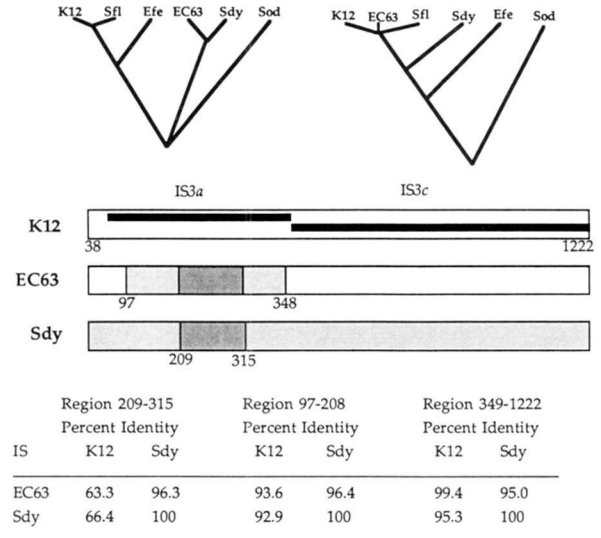


FIGURE 5.—Representation of intragenic recombination within IS3. K12, *E. coli* K12; EC63, *E. coli* ECOR 63; Efe, *E. fergusonii*; Sdy, *S. dysenteriae*; Sfl, *S. flexneri*; Sod, *S. odorifera*. The location of the two putative ORFs, IS3a and IS3c, are indicated with heavy lines. Dendrograms were constructed using parsimony methods from sequences of either the IS3a or IS3c reading frames. Consistency indices for the trees are 0.905 and 0.848, respectively.

TABLE 6

Length differences within insertion sequences from *E. coli*

IS	Strain	Position	Event
IS1	ECOR32	544-545	Deletion of AT
IS1	ECOR33	544-545	Deletion of AT
IS3	ECOR63	71	Deletion of A
IS3	ECOR63	98-99	Deletion of CC
IS3	ECOR18	493/494	Insertion of CAGT
IS3	ECOR50	1166-1182	Deletion of 17 bp
IS30	ECOR24	1071-1072	Deletion of TT
IS30	ECOR24	1070/1071	Insertion of IS3411

exhibited variation in length which would disrupt a major ORF.

The IS30 and IS1 sequences from other enteric species were the same length as the *E. coli* element, but several IS3 elements differed in size (Figure 3). The IS3 copy in *S. dysenteriae* contains a 3-bp insertion conserving the IS3c reading frame. The *S. odorifera* element exhibits 12- and 2-bp deletions, and a 2-bp insertion, which together conserve the reading frame of IS3c. The 3 bp deleted from IS3a in ECOR 63, although conserving the reading frame in this element, probably arose as two separate events (Table 6). Since the putative initiation codon has been mutated to AAG in this isoform, it is unclear whether the second deletion was favored by selection to restore the proper reading frame.

**Distribution of IS3411:** The IS30 from ECOR 24 exhibits a 1.2-kb insertion which proved to be IS3411. Although these sequences were plasmid borne, IS30::IS3411 transposons were not detected in other

**TABLE 7**  
Distribution of *IS3411* within the ECOR strains

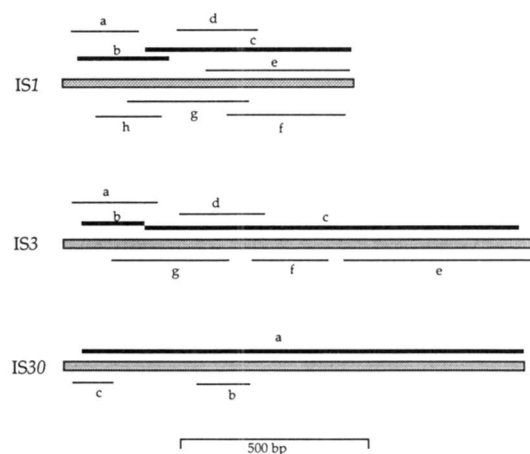
Strain <sup>a</sup>	Chromosomal	Plasmid	Strain	Chromosomal	Plasmid
2	6	1	39	13	2
3	6	1	40	2	0
4	1	1	41	0	1
7	0	1	43	0	2
8	3	0	44	0	2
9	6	1	49	0	1
10	6	1	50	8	1
11	11	1	51	1	1
13	0	1	52	1	0
15	1	0	53	0	1
18	1	0	54	0	1
19	5	0	55	1	0
20	2	0	56	1	0
21	1	1	57	1	0
23	0	1	59	0	1
24	7	6	60	1	0
25	5	0	61	1	0
30	1	1	62	6	1
34	4	1	63	1	0
35	0	7	64	8	1
36	9	1	70	2	0
37	0	2	71	2	1
38	0	2	72	2	1

<sup>a</sup> Number refers to ECOR strain (OCHMAN and SELANDER 1984b). Only strains harboring *IS3411* are listed.

ECOR strains. The distribution of *IS3411* within ECOR strains is presented in Table 7. A total of 174 copies of *IS3411*, 127 chromosomal and 47 plasmid, were detected among 46 strains. The copy number of 3.7 per infected strain is typical of other IS (SAWYER *et al.* 1987). Analysis of the distribution of *IS3411* among ECOR strains by the method of SAWYER *et al.* (1987) suggested a pattern of regulation similar to that of *IS1* and *IS5*. Models with no explicit regulation of transposition (that is, a linear increase of transposition rate with copy number), and a moderate decrease in fitness with increasing copy number, fit the data reasonably well ( $P \sim 0.25$ ). Models with constant or decreasing rate of transposition with increasing copy numbers were strongly rejected.

**ORF analysis:** GALAS and CHANDLER (1989) proposed 8, 7 and 3 potential ORF for *IS1*, *IS3* and *IS30*, respectively (Figure 6). To determine which ORFs were evolutionarily conserved, and therefore the most likely to encode proteins, nucleotide substitutions in each ORF within divergent copies of each IS were classified as nonsynonymous substitutions, which altered the amino acid sequence, or synonymous substitutions. Conserved ORFs were identified as those with an excess of synonymous substitutions relative to nonsynonymous substitutions, an absence of chain terminating substitutions, and preservation of initiation and termination codons. Data for certain comparisons are presented in Tables 8–10.

*IS1* contains two evolutionarily conserved ORFs,



**FIGURE 6.**—Genetic organization of bacterial insertion sequences (after GALAS and CHANDLER 1989). Reading frames transcribed from the positive strand are placed above each element; those transcribed from the complementary strand are placed below. Heavy lines indicate conserved ORFs.

**TABLE 8**  
Comparison of *IS1* ORFs from *E. coli* and *S. flexneri*

ORF <sup>a</sup>	Begin <sup>b</sup>	End	Stop <sup>c</sup>	Percent Divergence <sup>d</sup>		Protein similarity
				Replacement	Silent	
a	27 (GTG)	239 (TGA)		13.8	2.8	75.4
b	56 (GTG)	331 (TAA)		4.0	32.1	95.2
c	250 (ATG)	753 (TAA)		5.1	33.0	92.0
d	353 (GTG)	511 (TGA)	1	16.6	1.9	71.6
e	376 (ATG)	753 (TAA)		4.2	35.4	93.5
f	719 (ATG)	429 (TAA)	2	7.5	30.7	89.0
g	468 (GTG)	199 (TGA)	2	7.5	11.1	88.4
h	304 (GTG)	89 (GGA)	2	10.4	8.1	83.3
bc	56 (GTG)	753 (TAA)		3.8	36.8	94.6

<sup>a</sup> ORF notation of GALAS and CHANDLER (1989).

<sup>b</sup> Begin, putative initiation codon; End, putative termination codon. Codons in parentheses are present in *S. flexneri*.

<sup>c</sup> Stop = number of nonsense substitutions in the putative ORF.

<sup>d</sup> Divergence values calculated by the method of PERLER *et al.* (1980).

*IS1b* and *IS1c* (Table 8). SEKINE and OHTSUBO (1989) proposed ribosomal frameshifting between these two reading frames as a requisite for the proper production of the *IS1* transposase. This hypothesis is supported by the present data, since both ORFs are conserved among divergent copies. The putative full length protein is described by the *IS1bc* ORF (Table 8), created by a  $-1$  frameshift between these two reading frames. As expected, the 5' portion of *IS1c*, which would be translated in the frame of *IS1b*, contains an excess of nonsynonymous substitutions in the *IS1c* frame. The remaining potential reading frames in *IS1* are not evolutionarily conserved, bearing either nonsense mutations in the divergent homolog, loss of initiation or termination codons, or an excess of nonsynonymous substitutions which would yield a substantially altered protein product.

TABLE 9

Comparison of IS3 ORFs from *E. coli* and *S. odorifera*

ORF <sup>a</sup>	Begin <sup>b</sup>	End	Stop <sup>c</sup>	Percent divergence <sup>d</sup>		Protein similarity
				Replacement	Silent	
a	57 (GTG)	362 (TGA)		8.9	109.7	85.4
b	157 (GTG)	357 (TGA)		33.0	6.3	52.9
c	362 (ATG)	1225 (?)		8.5	116.4	88.5
d	415 (GTG)	573 (TGA)	3	44.3	34.4	55.6
e	1240 (?)	968 (TGA)	2	20.3	19.5	74.4
f	691 (ATG)	542 (TGC)	(1)	39.0	26.3	53.2
g	505 (GTC)	233 (TGT)	(1)	18.3	45.0	77.2

<sup>a</sup> ORF notation of GALAS and CHANDLER (1989).<sup>b</sup> Begin, putative initiation codon; End, putative termination codon. Codons in parentheses are present in *S. odorifera*.<sup>c</sup> Stop = number of nonsense substitutions in the putative ORF.<sup>d</sup> Divergence values calculated by the method of PERLER *et al.* (1980).

IS3 also exhibits two evolutionarily conserved ORFs, IS3a and IS3c (Table 9). Both ORFs present the excess of synonymous over nonsynonymous substitutions typical of protein coding regions. IS30 bears a single evolutionarily conserved ORF, IS30a, which covers 94% of its length (Table 10). Although IS30c appears somewhat conserved, it is encoded on the complementary DNA strand as IS30a (Figure 6). It is unlikely that IS30c encodes a peptide since it lacks conventional promoters and a Shine-Dalgarno sequence.

## DISCUSSION

**Variation within *E. coli*:** Eleven independent copies of IS1 representing both major types, 12 copies of IS3, and 12 copies of IS30 were isolated and examined from natural strains of *E. coli* and related enteric bacteria. The data presented in Tables 1–3 indicate that IS present in otherwise divergent isolates of *E. coli* are nearly monomorphic. The copies showed four or fewer substitutions relative to a consensus sequence (0–0.3% different), and these differences were typically restricted to a particular element. In contrast, sequence difference for chromosomal genes in representative strains of *E. coli* ranges from 0 to 3.4% at the *trp* locus (MILKMAN and CRAWFORD 1981), 2–4% at *phoA* (DUBOSE, DYKHUIZEN and HARTL 1988), and 4–16% at *gnd* (DYKHUIZEN and GREEN 1986). The sequence homogeneity of IS copies is also in strong contrast to the diversity in the number and chromosomal position of IS observed among *E. coli* strains (SAWYER *et al.* 1987; LAWRENCE *et al.* 1989).

The high rate of transposition of IS in *E. coli* genomes and their frequent occurrence on plasmids (SAWYER *et al.* 1987) suggests rapid turnover and frequent transfer of IS among strains. This model of population structure is supported by the sequence homogeneity among *E. coli* insertion sequences. Within *E. coli*, each family of IS turns over so rapidly

TABLE 10

Comparison of IS30 ORFs from *E. coli* and *E. hermannii*

ORF <sup>a</sup>	Begin <sup>b</sup>	End	Percent divergence <sup>c</sup>		Protein similarity
			Replacement	Silent	
a	63 (ATG)	1214 (TGA)	2.5	46.1	96.0
b	557 (GTG)	388 (TCG)	15.4	0.0	72.3
c	202 (ATG)	8 (TGA)	11.4	25.8	83.9

<sup>a</sup> ORF notation of GALAS and CHANDLER (1989).<sup>b</sup> Begin, putative initiation codon; End, putative termination codon. Codons in parentheses are present in *E. hermannii*.<sup>c</sup> Divergence values calculated by the method of PERLER *et al.* (1980).

that there is insufficient time to accumulate substantial genetic divergence. High rates of interstrain transfer are also supported by the statistically significant co-occurrence of unrelated IS in *E. coli* genomes, which can be explained quantitatively by simultaneous plasmid-mediated transfer (HARTL and SAWYER 1988a,b). Rapid dissemination and turnover of IS copies in *E. coli* contrasts with the much smaller rate of exchange among conventional chromosomal genes (DUBOSE, DYKHUIZEN and HARTL 1988; STOLTZFUS, LESLIE and MILKMAN 1988; SELANDER and LEVIN 1980; WHITTAM, OCHMAN and SELANDER 1984).

Alternatively, insertion sequences may have been conserved among *E. coli* isolates due to intense selection. Variation at the *gap* locus, encoding the glycolytic enzyme glyceraldehyde-3-phosphate dehydrogenase, is also quite low (NELSON, WHITTAM and SELANDER 1991), as selection on both function and codon usage constrain the evolution of this gene even among distantly related enteric bacteria (LAWRENCE, HARTL and OCHMAN 1991b). However, nucleotide variation among IS homologs from related enteric species do not support this model. Rather, the lack of nucleotide variation among *E. coli* IS supports a model of high turnover within *E. coli* genomes with significant mobilization between isolates.

**Variation among enteric bacteria:** In contrast to the high rate of transfer of IS among *E. coli* genomes, horizontal transfer of IS among bacterial species is less frequent. For example, the genome of the closely related bacterium *Salmonella typhimurium* contains none of the IS present in *E. coli* K12 (GALAS and CHANDLER 1989), and the only IS isolated from *S. typhimurium* (IS200) has not been detected in *E. coli* (LAM and ROTH 1983a,b). The difference in IS pools between *E. coli* and *Salmonella* implies either that IS are very recent invaders of enteric genomes, that transfer of IS among bacterial species is not common, or that *E. coli* IS cannot proliferate in *Salmonella* genomes.

The inferred relationships of IS3 homologs from *E. coli*, *S. dysenteriae*, *E. fergusonii* and *S. odorifera* are consistent with the phylogenetic relationships inferred



from chromosomal genes (LAWRENCE, HARTL and OCHMAN 1991b; LAWRENCE, OCHMAN and HARTL 1991). The IS $\beta$  copies from *E. fergusonii* are virtually identical to each other (Figure 3) and 1.9% and 15.6% divergent from the *E. coli* element at nonsynonymous and synonymous sites, respectively. This level of divergence is consistent with that observed for the *ompA* gene (LAWRENCE, HARTL and OCHMAN 1991b). The IS $\beta$  homolog from *S. odorifera* is substantially different from the *E. coli* sequence. It exhibits 9.0% and 119% divergence (corrected for multiple substitutions) at nonsynonymous and synonymous sites, respectively; these values are also consistent with the *ompA* data (LAWRENCE, HARTL and OCHMAN 1991b). The IS $\beta$  copies from *Shigella* are more closely related to those from *E. coli* than are those from *E. fergusonii*, which is consistent with the taxonomic placement of these species (EDWARDS and EWING 1962). These data suggest that IS $\beta$  elements have been resident in these species since their separation and have been evolving at approximately the same rate as chromosomal genes.

In contrast, the relationships inferred from IS1 sequences are not congruent with those inferred from chromosomal genes (Table 4). Sequence divergences among IS1 copies from *E. hermannii* and *E. vulneris* are substantially smaller than those observed for *gap* and *ompA*, both of which show strong conservation of both synonymous and nonsynonymous sites (LAWRENCE, OCHMAN and HARTL 1991). In addition, while the *ompA* data indicate that *E. fergusonii* is more closely related to *E. coli* (0.8% divergent at nonsynonymous sites) than to either *E. hermannii* or *E. vulneris* (5.0% and 9.8% divergent at nonsynonymous sites, respectively), the IS1 data indicate just the opposite (3.9% vs. 3.0% and 2.1% divergence at nonsynonymous sites for the IS1bc ORF. Comparisons with ISIR and ISIF yield similar results). Although these comparisons are potentially confounded by recombination events among these elements (see below), it is probable that IS1 was introduced into these genomes by horizontal transfer since their separation. In addition, the divergence of IS $\beta$  in *E. hermannii* is substantially less than that of *ompA* or *gap*, which suggests more recent invasion of this species by IS $\beta$ . Yet the absence of IS which closely resemble *E. coli* elements in other species indicates that the rate of horizontal transfer of IS among bacterial species must be considerably lower than the rate of intraspecific transfer.

**Intragenic recombination:** Both IS1 and IS $\beta$  provide evidence that intragenic recombination plays a role in the evolution of insertion sequences. Two nested intragenic recombinations involve IS $\beta$  homologs resident in *E. coli* (ECOR 63) and *S. dysenteriae*, as well as a distantly related element (Figure 5). The taxonomic relationships implicit in sequences within the recombined region (209–315 bp) are significantly

different from those implicit in the remainder of the element ( $P = 0.000$ , Table 5), supporting a transfer of genetic material between the *Shigella* IS $\beta$  and a distantly related element. Moreover, this region, as well as flanking sequences (97–348 bp), were transferred a second time between the *Shigella* IS $\beta$  and element now present in *E. coli* (Figure 5). As another indication of recombination, IS1 homologs in *E. vulneris* and *E. hermannii* are also chimeras of distantly related elements (Tables 4 and 5). At least three separate regions exhibiting statistically distinct evolutionary histories can be identified among IS1 sequences present in the genomes of enteric bacteria. Therefore, gene conversion, as well as rapid turnover, may play an important role in maintaining sequence homogeneity among IS within bacterial species. Recombination within bacterial genes has also been reported in *E. coli* (DUBOSE, DYKHUIZEN and HARTL 1988), *Streptococcus pneumoniae* (DOWSON *et al.* 1989), and *Neisseria gonorrhoeae* (SPRATT 1989).

**Length polymorphisms:** IS copies exhibit variation in length. Six of 35 elements (17%) examined from *E. coli* show length differences (Table 6), all but one being unique, and all the elements appear to be non-functional. In one case, IS $\beta$ 411 was inserted into an unrelated insertion sequence (IS $\beta$ 30 in ECOR 24). The insertion site reveals not only a 3-bp target duplication, but a deletion of 2 bp downstream of the site of insertion. In addition, the 4-bp insertion in the ECOR 18 copy of IS $\beta$  (cag/CAGT/tt) may be a remnant of the target site duplication of an excised transposon. In contrast, in certain natural populations of *Drosophila melanogaster*, the KP element, which contains a deletion of over 1700 bp, accounts for over 50% of the P element homologs (BLACK *et al.* 1987). In *Drosophila teissieri*, at least 75% of the *mariner* elements are deleted for more than half of their length (MARUYAMA and HARTL 1991). The few, unique, defective transposons detected among *E. coli* isolates indicate that substantially different factors contribute to the evolution of transposon populations within prokaryotic and eukaryotic genomes.

**Ribosomal frameshifting:** Analysis of nucleotide sequence variation revealed two evolutionarily conserved ORFs in divergent isolates of IS1 and IS $\beta$  and a single conserved ORF in IS $\beta$ 30. It is not surprising that IS1 maintains two conserved ORFs. SEKINE and OHTSUBO (1989) have demonstrated that ribosomal frameshifting between IS1b and IS1c is required for the production of the IS1 transposase (see also LÜTHI *et al.* 1990), and our data support this hypothesis. The single conserved ORF of IS $\beta$ 30 probably encodes that element's transposase. The finding of two conserved ORFs within IS $\beta$  was unexpected. However, a comparison between the IS1 and IS $\beta$  ORFs reveals striking similarity between the genetic organization of these

TABLE 11

Organization of ORFs and potential frameshifting signals in bacterial insertion sequences

Element	5' ORF <sup>a</sup>	3' ORF	Signal <sup>b</sup>	Location <sup>c</sup>
IS1	56-331 (b)	250-750 (c)	AAAAAAC	307
IS3	57-365 (a)	362-1225 (c)	AAAAGG	326
IS10	4-168 (a)	108-1313 (b)	AAAAAT	84
IS21	102-1274 (a)	1271-2068 (d)	??	
IS51	1261-938 (d)	941-42 (f)	AAAAAAC	1258
IS150	48-569 (a)	566-1414 (c)	AAAAAG	561
IS3411	55-381 (a)	378-1097 (b)	AAAAAAT	373

<sup>a</sup> Coordinates of 5' and 3' open reading frames in putative frameshifting events. Letters indicate ORF notation of Galas and Chandler (1989).

<sup>b</sup> Potential frameshifting signal.

<sup>c</sup> Coordinates of frameshifting signal.

unrelated elements. Both show two conserved open reading frames, the shorter located at the 5' end of the element and oriented in the -1 reading frame relative to the longer 3' ORF. Both exhibit the polyadenine frameshifting signal as well as a potential hairpin structure to facilitate ribosome slippage. In addition to IS1, these motifs have been implicated as the source of ribosomal frameshifting in the bovine leukemia virus (YOSHINAKA *et al.* 1986), mouse mammary tumor virus (HIZI *et al.* 1987; JACKS *et al.* 1987), and human T cell leukemia virus (SHIMOTOHNO *et al.* 1985) genomes. The 3' ORF in IS3, like that in IS1, lacks a conventional promoter region for transcription initiation and a Shine-Dalgarno sequence for translation initiation. These similarities suggest that ribosomal frameshifting may be also responsible for the production of the IS3 transposase.

The frameshifting sequence motif may also be identified in other bacterial insertion sequences. The 19 IS described by GALAS and CHANDLER (1989) may be separated into three classes: (1) IS bearing one, long, uninterrupted open reading frame, such as IS4, IS5, IS6, IS30, IS50, IS186, IS701 and IS903; (2) IS bearing numerous scattered ORFs, such as IS66, IS136 and IS600; and (3) IS bearing two closely situated open reading frames, the smaller 5' ORF being in the -1 reading frame relative to the 3' ORF. These elements are listed in Table 11. Although the genetic organization of certain elements resembles that of IS1, clearly not every IS is regulated by ribosomal frameshifting. IS21 exhibits no obvious frameshifting signal between its closely positioned ORFs, while the poly(A) signal in IS10 would lead to translation termination prior to reaching the second ORF. These cases aside, IS3, IS150 and IS3411 all exhibit the sets of sequence motifs responsible for frameshifting in IS1. Aside from IS3, however, it is not clear which ORFs are evolutionarily conserved, and there is no empirical evidence for frameshifting. Nevertheless, the repetition of this motif among bacterial IS suggests that

ribosomal frameshifting may be a common mechanism for the regulation of IS transposition.

In summary, DNA sequence variation among bacterial IS within *E. coli* and related species indicates that (1) insertion sequences are rapidly mobilized among strains within a species; (2) horizontal transfer of IS occurs, but is less frequent than intraspecific transfer; (3) few defective transposons are present in the bacterial genome, relative to nonautonomous elements within eukaryotic genomes; (4) recombination occurs between IS at significant rates; and (5) functionally important ORFs can be recognized by their evolutionary conservation.

We thank S. SAWYER for aid with the IS3411 analysis, M. STURMOSKI for technical assistance, and D. DYKHUIZEN and an anonymous reviewer for helpful comments on the manuscript. This work was supported by grants GM 40322 (D.L.H.) and GM 40995 (H.O.) from the National Institutes of Health. The DNA sequences reported in this paper have been submitted to GenBank and have been assigned accession numbers Z11603 through Z11609.

#### LITERATURE CITED

- BERG, D. E., 1977 Insertion and excision of the transposable kanamycin resistance determinant Tn5, pp. 205-212 in *DNA Insertion Elements, Plasmids, and Episomes*, edited by A. BUKHARI, J. A. SHAPIRO and S. ADHYA. Cold Spring Harbor Laboratory, Cold Spring Harbor, N.Y.
- BLACK, D. M., M. S. JACKSON, M. G. KIDWELL and A. DOVER, 1987 KP elements repress P-induced hybrid dysgenesis in *Drosophila melanogaster* using a novel and general method. *Cell* **25**: 693-704.
- CAUGANT, D. A., B. R. LEVIN and R. K. SELANDER, 1981 Genetic diversity and temporal variation in the *E. coli* population of a human host. *Genetics* **98**: 467-490.
- CHANDLER, M., M. CLERGET and L. CARO, 1980 IS1-promoted events associated with drug resistance plasmids. Cold Spring Harbor Symp. Quant. Biol. **45**: 157-165.
- CHOW, L. T., and T. R. BROKER, 1981 Adjacent insertion sequences IS2 and IS5 in bacteriophage Mu mutants and IS5 in lambda *darg* bacteriophage. *J. Bacteriol.* **133**: 1427-1436.
- DALRYMPLE, B., P. CASPERS and W. ARBER, 1984 Nucleotide sequence of the prokaryotic mobile genetic element IS30. *EMBO J.* **3**: 2145-2149.
- DEVEREUX, J., P. HAEBERLI and O. SMITHIES, 1984 A comprehensive set of sequence analysis programs for the VAX. *Nucleic Acids Res.* **12**: 387-395.
- DOWER, W. J., J. F. MILLER and C. RAGSDALE, 1988 High efficiency transformation of *E. coli* by high voltage electroporation. *Nucleic Acids Res.* **16**: 6127-6145.
- DOWSON, C. G., A. HUTCHINSON, J. A. BRANNIGAN, R. C. GEORGE, D. HANSMAN, J. LINARES, A. TOMASZ, J. MAYNARD SMITH and B. G. SPRATT, 1989 Horizontal transfer of penicillin-binding genes in penicillin-resistant clinical isolates of *Streptococcus pneumoniae*. *Proc. Natl. Acad. Sci. USA* **86**: 8842-8846.
- DUBOSE, R. F., D. E. DYKHUIZEN and D. L. HARTL, 1988 Genetic exchange among natural isolates of bacteria: recombination within the *phoA* locus of *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **85**: 7036-7040.
- DUBOSE, R. F., and D. L. HARTL, 1990 Rapid purification of PCR products for DNA sequencing using Sepharose CL-6B spin columns. *Biotechniques* **8**: 271-273.
- DYKHUIZEN, D. E., and L. GREEN, 1986 DNA sequence variation, DNA phylogeny, and recombination in *E. coli*. *Genetics* **113** (Suppl.): s71.

- EDWARDS, P. R., and W. H. EWING, 1962 *Identification of Enterobacteriaceae*. Burgess, Minneapolis.
- FEINBERG, A. P., and B. VOGELSTEIN, 1983 A technique for radiolabeling DNA restriction endonuclease fragments to high specific activity. *Anal. Biochem.* **132**: 6–13.
- GALAS, D. J., and M. CHANDLER, 1989 Bacterial insertion sequences, pp. 109–162 in *Mobile DNA*, edited by D. E. BERG and M. M. HOWE. American Society for Microbiology, Washington, D.C.
- GLANDSDORFF, N., D. CHARLIER and M. ZAFARULLAH, 1980 Activation of gene expression by IS2 and IS3. Cold Spring Harbor Symp. Quant. Biol. **45**: 153–156.
- HALL, B. G., L. L. PARKER, P. W. BETTS, R. F. DUBOSE, S. A. SAWYER and D. L. HARTL, 1989 IS103: a new insertion element in *Escherichia coli*. Characterization and distribution in natural populations. *Genetics* **121**: 423–431.
- HARTL, D. L., and S. A. SAWYER, 1988a Multiple correlations among insertion sequences in the genome of natural isolates of *Escherichia coli*, pp. 91–106 in *Transposition*, edited by A. J. KINGSMAN, S. M. KINGSMAN and K. F. CHATER. Cambridge University Press, Cambridge.
- HARTL, D. L., and S. A. SAWYER, 1988b Why do unrelated insertion sequences occur together in the genome of *Escherichia coli*? *Genetics* **118**: 537–541.
- HELMS, C., M. Y. GRAHAM, J. E. DUTCHIK and M. V. OLSON, 1985 A new method for purifying lambda DNA from phage lysates. *DNA* **4**: 39–49.
- HIZI, A., L. E. HENDERSON, T. D. COPELAND, R. C. SOWDER, C. V. HIXSON and S. OROSZLAN, 1987 Characterization of mouse mammary tumor virus *gag-pro* gene products and the ribosomal frameshift site by protein sequencing. *Proc. Natl. Acad. Sci. USA* **84**: 7041–7045.
- JACKS, T. K. TOWNSLEY, H. E. VARMUS and J. MAJORS, 1987 Two efficient ribosomal frameshifting events are required for synthesis of mouse mammary tumor virus *gag*-related polyproteins. *Proc. Natl. Acad. Sci. USA* **84**: 4298–4302.
- KLECKNER, N., R. K. CHAN, B. K. TYE and D. BOTSTEIN, 1975 Mutagenesis by insertion of a drug resistance element carrying an inverted repetition. *J. Mol. Biol.* **97**: 561–575.
- LAM, S., and J. R. ROTH, 1983a Genetic mapping of IS200 copies in *Salmonella typhimurium* LT2. *Genetics* **105**: 801–811.
- LAM, S., and J. R. ROTH, 1983b IS200: a *Salmonella*-specific insertion sequence. *Cell* **34**: 951–960.
- LAWRENCE, J. G., and D. L. HARTL, 1992 Inference of horizontal genetic transfer from molecular data: an approach using the bootstrap. *Genetics* (in press).
- LAWRENCE, J. G., D. L. HARTL and H. OCHMAN, 1991a Sequencing products of the polymerase chain reaction. *Methods Enzymol.* (in press).
- LAWRENCE, J. G., D. L. HARTL, and H. OCHMAN, 1991b Molecular considerations in the evolution of bacterial genes. *J. Mol. Evol.* **33**: 241–250.
- LAWRENCE, J. G., H. OCHMAN and D. L. HARTL, 1991 Molecular and evolutionary relationships among enteric bacteria. *J. Gen. Microbiol.* **137**: 1911–1921.
- LAWRENCE, J. G., D. E. DYKHUIZEN, R. F. DUBOSE and D. L. HARTL, 1989 Phylogenetic analysis of *Escherichia coli* based on insertion sequence fingerprinting. *Mol. Biol. Evol.* **6**: 1–14.
- LIAB, M., 1980 IS5 increases recombination in adjacent regions as shown for the repressor gene of coliphage lambda. *Gene* **12**: 277–280.
- LÜTHI, K., M. MOSER, J. RYSER and H. WEBER, 1990 Evidence for a role of translational frameshifting in the expression of transposition activity of the bacterial insertion element IS1. *Gene* **88**: 15–20.
- MARUYAMA, K., and D. L. HARTL, 1991 Evolution of the transposable element *mariner* in *Drosophila* species. *Genetics* **128**: 319–329.
- MILKMAN, R., and I. P. CRAWFORD, 1981 Clustered third base substitutions among wild strains of *Escherichia coli*. *Science* **221**: 378–379.
- NELSON, K., T. S. WHITTAM and R. K. SELANDER, 1991 Nucleotide polymorphism and evolution in the glyceraldehyde-3-phosphate gene (*gapA*) in natural populations of *Salmonella* and *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **88**: 6667–6671.
- OCHMAN, H., and R. K. SELANDER, 1984a Evidence for clonal population structure in *Escherichia coli*. *Proc. Natl. Acad. Sci. USA* **81**: 198–201.
- OCHMAN, H., and R. K. SELANDER, 1984b Standard reference collection of *Escherichia coli* from natural populations. *J. Bacteriol.* **157**: 517–524.
- OHTSUBO, H., and E. OHTSUBO, 1978 Nucleotide sequence of an insertion element, IS1. *Proc. Natl. Acad. Sci. USA* **75**: 615–619.
- OHTSUBO, E., H. OHTSUBO, W. DOROSZKIEWICZ, K. NYMAN, D. ALLEN and D. DAVIDSON, 1984 An evolutionary analysis of iso-IS1 elements from *Escherichia coli* and *Shigella* strains. *J. Gen. Appl. Microbiol.* **30**: 359–376.
- PERLER, F., A. EFSTRATIADIS, P. LOMEDICO, W. GILBERT, R. KOLODNER and J. DODGSON, 1980 The evolution of genes: the chicken preproinsulin gene. *Cell* **20**: 555–566.
- SAEDLER, H., H. J. REIF, S. HU and N. DAVIDSON, 1974 IS2: a genetic element for turn-off and turn-on of gene activity in *E. coli*. *Mol. Gen. Genet.* **132**: 265–289.
- SAEDLER, H., G. CORNELIS, J. CULLUM, B. SCHUMACHER and H. SOMMER, 1980 IS1 mediated DNA rearrangements. Cold Spring Harbor Symp. Quant. Biol. **45**: 93–98.
- SAIKI, R. K., S. SCHARF, F. FALOONA, K. B. MULLIS, G. T., HORN, H. A. ERLICH and N. A. ARNHEIM, 1985 Enzymatic amplification of beta-globin genomic sequences and restriction site analysis for diagnosis of sickle cell anemia. *Science* **230**: 1350–1354.
- SAIKI, R. K., D. H. GELFAND, S. STOFFEL, S. J. SCHARF, R. HIGUCHI, G. T. HORN, K. B. MULLIS and H. A. ERLICH, 1988 Primer-directed enzymatic amplification of DNA with a thermostable DNA polymerase. *Science* **239**: 487–491.
- SAWYER, S., D. DYKHUIZEN, R. DUBOSE, L. GREEN, T. MUTANGADURA-MHLANGA, D. WOLCZYK and D. HARTL, 1987 Distribution and abundance of insertion sequences among natural isolates of *Escherichia coli*. *Genetics* **115**: 51–63.
- SEKINE, Y., and E. OHTSUBO, 1989 Frameshifting is required for production of the transposase encoded by insertion sequence I. *Proc. Natl. Acad. Sci. USA* **86**: 4609–4613.
- SELANDER, R. K., and B. R. LEVIN, 1980 Genetic diversity and structure in *Escherichia coli* populations. *Science* **210**: 545–547.
- SELANDER, R. K., D. A. CAUGANT and T. S. WHITTAM, 1987 Genetic structure and variation in natural populations of *Escherichia coli*, pp. 1625–1648 in *Escherichia coli and Salmonella typhimurium: Cellular and Molecular Biology*, edited by F. C. NEIDHARDT. American Society for Microbiology, Washington, D.C.
- SHIMOTOHNO, K., Y. TAKAHASHI, T. GOJOBORI, D. W. GOLDE, I. S. Y. CHEN, M. MIWA and T. SUGIMURA, 1985 Complete nucleotide sequence of an infectious clone of human T-cell leukemia virus type II: an open reading frame for the protease gene. *Proc. Natl. Acad. Sci. USA* **82**: 3101–3105.
- SPRATT, B. G., 1989 Hybrid penicillin-binding proteins in penicillin-resistant strains of *Neisseria gonorrhoeae*. *Nature* **332**: 173–176.
- STOLTZFUS, A., J. F. LESLIE and R. MILKMAN, 1988 Molecular evolution of the *Escherichia coli* chromosome. I. Analysis of structure and natural variation in a previously uncharacterized region between *trp* and *tonB*. *Genetics* **120**: 345–358.
- TIMMERMAN, K., and C. P. D. TU, 1985 Complete sequence of IS3. *Nuc Acids Res* **13**: 2127–2139.

WHITTAM, T. S., H. OCHMAN and R. K. SELANDER, 1984 Geographical components of linkage disequilibrium in natural populations of *Escherichia coli*. *Mol. Biol. Evol.* **1**: 67-83.

UMEDA, M., and E. OHTSUBO, 1991 Four types of *IS1* with differences in nucleotide sequence reside in the *Escherichia coli* K-12 chromosome. *Gene* **98**: 1-5.

YOSHINAKA, Y., I. KATOH, T. D. COPELAND, G. W. SMYTHERS and S. OROSZLAN, 1986 Bovine leukemia virus protease: purification, chemical analysis, and in-vitro processing of *gag* precursor polyproteins. *J. Virol.* **57**: 826-832.

Communicating editor: W.-H. LI