

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/7193652>

Identification of genes subject to positive selection in uropathogenic strains of Escherichia coli: A comparative genomics approach

Article in *Proceedings of the National Academy of Sciences* · May 2006

DOI: 10.1073/pnas.0600938103 · Source: PubMed

CITATIONS

365

READS

157

18 authors, including:



Jian Xu

Chinese Academy of Sciences

212 PUBLICATIONS **5,504** CITATIONS

[SEE PROFILE](#)



Christopher S Reigstad

Concordia University- Chicago

25 PUBLICATIONS **1,404** CITATIONS

[SEE PROFILE](#)



Vincent Magrini

Washington University in St. Louis

214 PUBLICATIONS **25,400** CITATIONS

[SEE PROFILE](#)



Tamberlyn Bieri

Washington University in St. Louis

18 PUBLICATIONS **4,491** CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Unique organization and unprecedented diversity of the Bacteroides (Pseudobacteroides) cellulosolvens cellulosome system [View project](#)



Metastatic Colon Cancer lncRNAs [View project](#)

found in a new habitat. Recent work has begun to surmount both of these problems (22, 23): given enough sequences from related organisms, it is possible to detect positive selection occurring only in a portion of a gene in a subset of the sequences. In this report, we have used this comparative approach to understand the adaptations that *E. coli* has made to colonize and survive in the urinary tract. Using the newly sequenced genome of UTI89, a UPEC strain isolated from a patient with an acute bladder infection (14), we have identified coding sequences common to all *E. coli* but under positive selection only in UPEC strains. These results were validated by using a panel of clinical *E. coli* strains isolated from patients with UTI. Our computational method circumvents the technical difficulties that hamper *in vivo* genetic screens of UPEC in mice and may be broadly applicable to understanding strain-specific adaptation and pathogenesis in other bacteria.

Results and Discussion

UTI89 contains a 5,065,741-bp chromosome and a 114,230-bp plasmid (pUTI89). The chromosome has a GC content of 50.6%, 5,066 predicted protein-coding genes, 88 tRNA genes, and 22 rRNA genes. These values are similar to other sequenced *E. coli* strains (Table 3, which is published as supporting information on the PNAS web site). pUTI89 has a GC content of 51.0% and 145 predicted genes. The UTI89 genome contains four large pathogenicity islands (PAIs) similar to previously characterized PAIs. Other notable features include 10 putative pilus operons and ORFs encoding several toxins (two enterotoxins, a hemolysin, and cytotoxic necrotizing factor). A more detailed description of the PAIs, the genes that produce adhesive organelles, and the plasmid can be found in *Supporting Text*, which is published as supporting information on the PNAS web site.

A functional categorization of the predicted proteome of UTI89 using Clusters of Orthologous Groups (COG) is shown in Fig. 4, which is published as supporting information on the PNAS web site, and a metabolic reconstruction based on the Kyoto Encyclopedia of Genes and Genomes (KEGG) is deposited at <http://hultgren.wustl.edu/UTI89> (for similar reconstructions of the other *E. coli* genomes, see www.genome.jp/kegg).

The number of shared and novel genes was calculated (*Supporting Text* and ref. 9) for all subsets of the seven *E. coli* genomes, and an exponential curve was fit to the data to predict the number of shared and new genes that would be contributed by additional genome sequences. The size of the *E. coli* core genome (genes shared among all *E. coli*) is predicted to be 2,865. Like the pan-genomes of Group A *Streptococcus* and Group B *Streptococcus* (9), the pan-genome of *E. coli* is open: the number of new genes contributed by each new *E. coli* genome is predicted to be 441 genes, which is substantially higher than the 27 and 33 genes predicted to be added with each new Group A and B *Streptococcus* genome (9).

Identifying Genes Under Selection in UPEC Strains. To identify genes under positive selection in UPEC strains, we performed comparisons between all seven fully sequenced *E. coli* genomes using the scheme outlined in Fig. 1A. Identification of positive selection using the PAML program requires an aligned set of related sequences (orthologs) and knowledge of their phylogenetic relationships. Reciprocal best BLAST hits in the seven genomes were assumed to represent orthologous sequences and were aligned by using CLUSTALW. The phylogenetic relationship between sequences can consist of vertical components (direct mother-to-daughter transmission) and horizontal components (all other DNA transfers, including gene conversion, recombination, and lateral transfer). In *E. coli*, a given nucleotide difference is 10–50 times more likely to have been caused by recombination or gene conversion (horizontal relationships) than by mutation (vertical) (24, 25). Therefore, accounting for horizontal relationships is crucial for inferring an accurate phylogeny (26–28) for subsequent use in detecting selection (29,

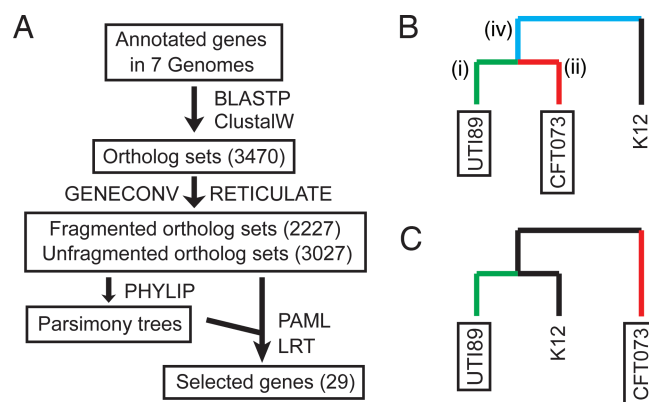


Fig. 1. Overview of the analysis and specification of foreground branches. (A) Analysis scheme. Size of each dataset (boxes) is indicated in parentheses. Programs used are indicated next to the arrows. See text for more details about how datasets were generated. (B and C) Hypothetical phylogenetic tree to indicate branch specification. UPEC strains are boxed. Evidence for positive selection was evaluated in specific lineages (termed foreground branches). B shows the sets of UPEC-specific foreground branches used: (i) UTI89 only (green), (ii) CFT073 only (red), (iii) UPEC leaves only (green + red), (iv) common UPEC branch only (cyan), (v) all UPEC (green + red + cyan). C shows that when the common UPEC branch was not present due to the tree topology, only three sets of UPEC foreground branches were used (green, red, and green + red).

30). These relationships can be assessed by using methods to detect recombination, ideally with at least two methods that differ in their underlying theory to maximize accuracy (26, 31). After accounting for recombination, vertical relationships can then be inferred with standard phylogenetic tree-building software.

We identified 3,470 ortholog sets with at least four members (the minimum number needed for subsequent analysis), representing 68.5% of all UTI89 genes. The aligned ortholog sets were tested for evidence of recombination by using the programs GENECONV and RETICULATE. Of the 3,470 ortholog sets, 443 (12.8%) showed evidence of recombination using both programs. GENECONV also predicts breakpoints where recombination has occurred. These breakpoints define subsegments (fragments) of each ortholog. Adjacent fragments, despite being in the same ortholog, have different evolutionary histories due to recombination (26). A total of 2,227 fragment sets were created from the 443 ortholog sets with evidence of recombination. The PHYLIP software package was then used to infer maximum parsimony trees from the 3,027 ortholog sets that showed no evidence of recombination and from the 2,227 fragment sets.

Positive selection was identified with the program PAML. Using a maximum likelihood algorithm, PAML assigns likelihood scores to different hypotheses (models) for selection. If a model incorporating positive selection has a higher likelihood score than a null model without positive selection, this constitutes evidence for positive selection. The null model is referred to as M1a, and the selection model as M2a. Comparison of M1a and M2a tests whether a gene is under selection in all of the sequenced *E. coli* strains. A third model (bsA) is based on the hypothesis that positive selection occurs only in certain branches/lineages: comparing the likelihoods of bsA and M1a tests whether a gene is under positive selection in a specific lineage, such as UPEC strains. In the bsA model, the branches hypothesized to have positive selection must be specified and are referred to as “foreground branches” (Fig. 1B and C).

Using the aligned ortholog and fragment sets and their corresponding phylogenetic trees, likelihood scores were assigned for the M1a, M2a, and bsA models. Ortholog and fragment sets that showed evidence for positive selection in UPEC branches (bsA versus M1a) but not in all *E. coli* (M2a versus M1a) were then inspected individually. Table 1 lists the resulting 29 genes desig-

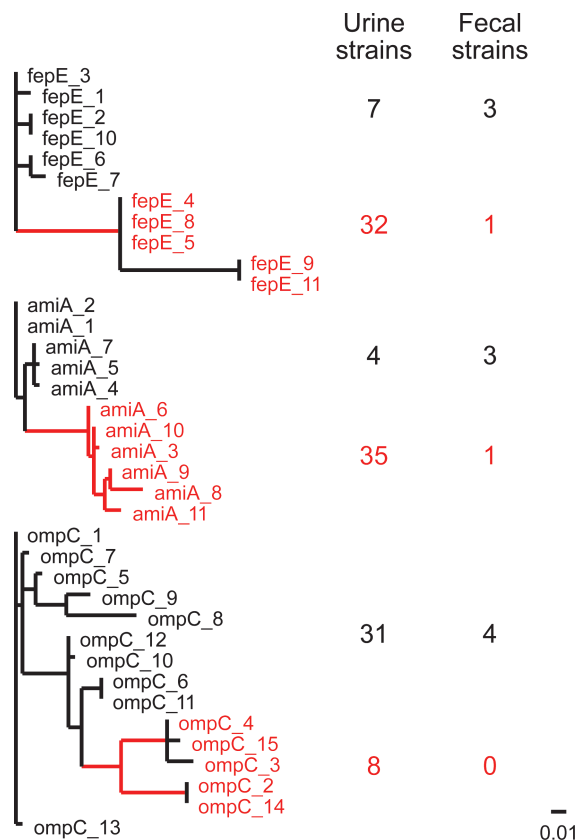


Fig. 2. *amiA*, *fepE*, and *ompC* are under selection in clinical UPEC isolates. Phylogenetic trees of unique sequences for each gene are shown. Red branches/labels indicate foreground branches that show evidence for positive selection (see Table 2 and text). Red numbers to the right of the tree indicate the number of urine and fecal isolates represented by the red (foreground) labels in the phylogenetic tree. Black numbers to the right of the tree indicate the number of urine and fecal isolates represented by black labels in the phylogenetic tree. The sites of isolation of strains represented by sequence labels are summarized in Table 4. Scale bar for phylogenetic trees is shown at the bottom right.

all of the strains (see below). This directly mirrors the analysis of *ompC* in the seven genomes, where the gene showed evidence of positive selection in one of the two UPEC strains and weaker evidence of selection in all of the strains. No evidence for positive selection was found in any branch for the adenylate kinase (*adk*), malate dehydrogenase (*mdh*), and isocitrate dehydrogenase (*icd*) genes, further confirming the results we obtained using the seven fully sequenced *E. coli* genomes.

Functional Analysis. To gain additional insights into the evolutionary pressures acting on UPEC strains, we analyzed the functional

Table 2. Analysis of resequenced genes

Gene	P value	d_n/d_s	Base pairs sequenced
<i>fepE</i> (UTI89_C0589)	3.94E-02	86.96	375
<i>ompC</i> (UTI89_C2497)	1.29E-07	999.00	756
<i>amiA</i> (UTI89_C2768)	5.48E-03	7.95	279
<i>adk</i> (UTI89_C0502)	1.00E+00	NA	501
<i>mdh</i> (UTI89_C3667)	1.00E+00	NA	435
<i>icd</i> (UTI89_C1266)	4.46E-01	NA	666

P value and d_n/d_s values from the analysis of resequenced clinical strains are shown. Number of base pairs sequenced and analyzed is shown in the last column. A d_n/d_s value of N/A means that no evidence for positive selection was found.

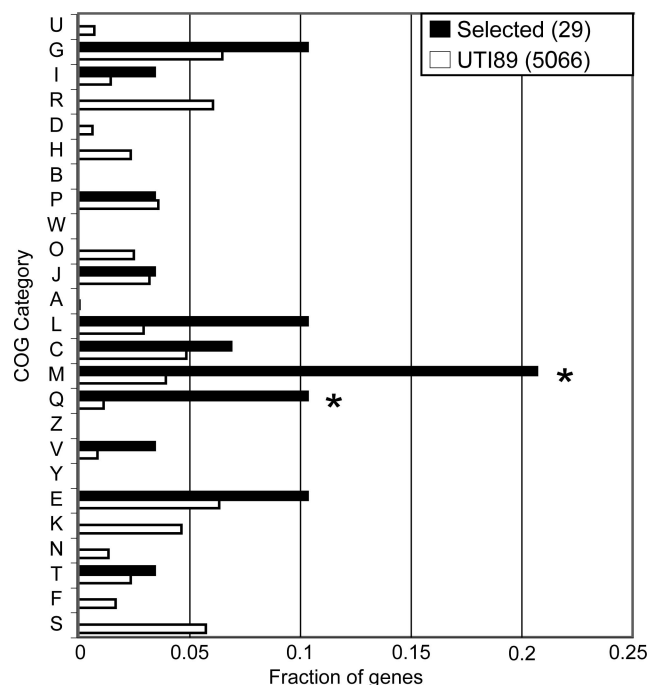


Fig. 3. UPEC-selected genes are enriched for genes in two COG functional categories. COG category codes are indicated on the y axis. The fraction of genes in each COG category is shown on the x axis. Black bars indicate genes under positive selection in UPEC strains (n = 29). Open bars are for all genes annotated in UTI89 (n = 5,066). COG categories that are significantly enriched (P < 0.05, binomial test) in the set of UPEC-selected genes relative to all UTI89 genes are indicated by an asterisk. COG category codes are as follows: U, intracellular trafficking and secretion; G, carbohydrate transport and metabolism; I, lipid transport and metabolism; R, general function prediction only; D, cell cycle control, mitosis and meiosis; H, coenzyme transport and metabolism; B, chromatin structure and dynamics; P, inorganic ion transport and metabolism; W, extracellular structures; O, posttranslational modification, protein turnover, chaperones; J, translation; A, RNA processing and modification; L, replication, recombination and repair; C, energy production and conversion; M, cell wall/membrane biogenesis; Q, secondary metabolites biosynthesis, transport and catabolism; Z, cytoskeleton; V, defense mechanisms; Y, nuclear structure; E, amino acid transport and metabolism; K, transcription; N, cell motility; T, signal transduction mechanisms; F, nucleotide transport and metabolism; and S, function unknown.

features of the 29 USG. Two COG categories were significantly enriched (P < 0.05, binomial test) among the USG compared with the entire set of genes in the UTI89 genome: (i) M, cell wall/membrane biogenesis; and (ii) Q, secondary metabolites biosynthesis, transport, and catabolism (Fig. 3). A third COG category (L, replication, recombination, and repair) was also enriched but with borderline significance (P = 0.0527).

UPEC often have high mutation rates compared with other *E. coli* strains (up to 5×10^{-7} per nucleotide per generation, a value that is 100- to 1,000-fold higher than for wild-type strains) (24, 35). Notably, a generally high mutation rate would conservatively bias our analysis (see *Supporting Text* for further details). Elevated mutation rate confers a fitness advantage to UPEC strains in a mouse UTI model: mutator strains persist longer in bladder and kidney than wild-type strains, and serial passage increases the virulence of mutator strains over that of wild-type (36). The strains tested were *mutS* mutants that are defective in initial recognition of mispaired DNA bases (36). Little is known about the mechanistic basis of elevated mutation rates in UPEC and its significance in the context of the drastic changes in population size (P. C. Seed, K. J. Wright, G. G. Anderson, and S.J.H., unpublished observations) and growth rate (12) that occur during IBC formation. Intriguingly,

Detection of Recombination. Both GENECONV and RETICULATE were run on the aligned nucleotide sequences: GENECONV was run by using the “/r” (silent sites only) option; RETICULATE was run with the “treat sites with more than two characters as sites with more than two characters” option, and *P* values were calculated from 10,000 randomizations of the data. If the *P* values reported by both programs were <0.05 , the aligned sequences were fragmented at recombination breakpoints identified by GENECONV (using the endpoints of all reported fragments). These fragments were then treated independently to infer phylogenetic trees and detect selection.

Detection of Selection. The CODEML program from the PAML package (Version 3.14b) was used for all calculations. The following

models were run for each set of genes or gene fragments (if recombination was detected): site models M1a and M2a and branch-site model A for each set of foreground branches. Options set in the control file followed those in the lysozyme example directory of the PAML distribution package.

A likelihood ratio test was used to compare model M2a with M1a, and branch-site model A with model M1a. The significance cutoff was set at 1/3,470 (the reciprocal of the number of genes tested).

We thank Justin Fay for suggestions and comments on the text. This work was supported by National Institutes of Health Grants P50-ARO49475, T32HG000045, DK64540, and DK51406.

- McFeters, G. A., Barry, J. P. & Howington, J. P. (1993) *Water Res.* **27**, 645–650.
- Keseler, I. M., Collado-Vides, J., Gama-Castro, S., Ingraham, J., Paley, S., Paulsen, I. T., Peralta-Gil, M. & Karp, P. D. (2005) *Nucleic Acids Res.* **33**, D334–D337.
- Blattner, F. R., Plunkett, G., 3rd, Bloch, C. A., Perna, N. T., Burland, V., Riley, M., Collado-Vides, J., Glasner, J. D., Rode, C. K., Mayhew, G. F., *et al.* (1997) *Science* **277**, 1453–1474.
- Welch, R. A., Burland, V., Plunkett, G., 3rd, Redford, P., Roesch, P., Rasko, D., Buckles, E. L., Liou, S. R., Boutin, A., Hackett, J., *et al.* (2002) *Proc. Natl. Acad. Sci. USA* **99**, 17020–17024.
- Perna, N. T., Plunkett, G., 3rd, Burland, V., Mau, B., Glasner, J. D., Rose, D. J., Mayhew, G. F., Evans, P. S., Gregor, J., Kirkpatrick, H. A., *et al.* (2001) *Nature* **409**, 529–533.
- Hayashi, T., Makino, K., Ohnishi, M., Kurokawa, K., Ishii, K., Yokoyama, K., Han, C. G., Ohtsubo, E., Nakayama, K., Murata, T., *et al.* (2001) *DNA Res.* **8**, 11–22.
- Wei, J., Goldberg, M. B., Burland, V., Venkatesan, M. M., Deng, W., Fournier, G., Mayhew, G. F., Plunkett, G., 3rd, Rose, D. J., Darling, A., *et al.* (2003) *Infect. Immun.* **71**, 2775–2786.
- Jin, Q., Yuan, Z., Xu, J., Wang, Y., Shen, Y., Lu, W., Wang, J., Liu, H., Yang, J., Yang, F., *et al.* (2002) *Nucleic Acids Res.* **30**, 4432–4441.
- Tettelin, H., Maignani, V., Cieslewicz, M. J., Donati, C., Medini, D., Ward, N. L., Angiuoli, S. V., Crabtree, J., Jones, A. L., Durkin, A. S., *et al.* (2005) *Proc. Natl. Acad. Sci. USA* **102**, 13950–13955.
- Mulvey, M. A., Lopez-Boado, Y. S., Wilson, C. L., Roth, R., Parks, W. C., Heuser, J. & Hultgren, S. J. (1998) *Science* **282**, 1494–1497.
- Martinez, J. J., Mulvey, M. A., Schilling, J. D., Pinkner, J. S. & Hultgren, S. J. (2000) *EMBO J.* **19**, 2803–2812.
- Justice, S. S., Hung, C., Theriot, J. A., Fletcher, D. A., Anderson, G. G., Footer, M. J. & Hultgren, S. J. (2004) *Proc. Natl. Acad. Sci. USA* **101**, 1333–1338.
- Anderson, G. G., Palermo, J. J., Schilling, J. D., Roth, R., Heuser, J. & Hultgren, S. J. (2003) *Science* **301**, 105–107.
- Mulvey, M. A., Schilling, J. D. & Hultgren, S. J. (2001) *Infect. Immun.* **69**, 4572–4579.
- Zhang, L., Foxman, B., Manning, S. D., Tallman, P. & Marrs, C. F. (2000) *Infect. Immun.* **68**, 2009–2015.
- Guyer, D. M., Kao, J. S. & Mobley, H. L. (1998) *Infect. Immun.* **66**, 4411–4417.
- Sannes, M. R., Kuskowski, M. A., Owens, K., Gajewski, A. & Johnson, J. R. (2004) *J. Infect. Dis.* **190**, 2121–2128.
- Srinivasan, U., Foxman, B. & Marrs, C. F. (2003) *J. Clin. Microbiol.* **41**, 285–289.
- Sokurenko, E. V., Chesnokova, V., Dykhuizen, D. E., Ofek, I., Wu, X. R., Krogfelt, K. A., Struve, C., Schembri, M. A. & Hasty, D. L. (1998) *Proc. Natl. Acad. Sci. USA* **95**, 8922–8926.
- Kreitman, M. (2000) *Annu. Rev. Genom. Hum. Genet.* **1**, 539–559.
- Fay, J. C. & Wu, C. I. (2003) *Annu. Rev. Genom. Hum. Genet.* **4**, 213–235.
- Yang, Z. & Nielsen, R. (2002) *Mol. Biol. Evol.* **19**, 908–917.
- Suzuki, Y. & Nei, M. (2001) *Mol. Biol. Evol.* **18**, 2179–2185.
- Guttman, D. S. & Dykhuizen, D. E. (1994) *Science* **266**, 1380–1383.
- Feil, E. J., Maiden, M. C., Achtman, M. & Spratt, B. G. (1999) *Mol. Biol. Evol.* **16**, 1496–1502.
- Wiuf, C., Christensen, T. & Hein, J. (2001) *Mol. Biol. Evol.* **18**, 1929–1939.
- Jakobsen, I. B. & Easteal, S. (1996) *Comput. Appl. Biosci.* **12**, 291–295.
- Feil, E. J. & Spratt, B. G. (2001) *Annu. Rev. Microbiol.* **55**, 561–590.
- Anisimova, M., Nielsen, R. & Yang, Z. (2003) *Genetics* **164**, 1229–1236.
- Suzuki, Y. & Nei, M. (2004) *Mol. Biol. Evol.* **21**, 914–921.
- Posada, D. (2002) *Mol. Biol. Evol.* **19**, 708–717.
- Wong, W. S., Yang, Z., Goldman, N. & Nielsen, R. (2004) *Genetics* **168**, 1041–1051.
- Anisimova, M., Bielawski, J. P. & Yang, Z. (2002) *Mol. Biol. Evol.* **19**, 950–958.
- Anisimova, M., Bielawski, J. P. & Yang, Z. (2001) *Mol. Biol. Evol.* **18**, 1585–1592.
- Denamur, E., Bonacorsi, S., Giraud, A., Duriez, P., Hilali, F., Amorin, C., Bingen, E., Andreumont, A., Picard, B., Taddei, F., *et al.* (2002) *J. Bacteriol.* **184**, 605–609.
- Labat, F., Pradillon, O., Garry, L., Peuchmaur, M., Fantin, B. & Denamur, E. (2005) *FEMS Immunol. Med. Microbiol.* **44**, 317–321.
- Viswanathan, M., Burdett, V., Baitinger, C., Modrich, P. & Lovett, S. T. (2001) *J. Biol. Chem.* **276**, 31053–31058.
- Burdett, V., Baitinger, C., Viswanathan, M., Lovett, S. T. & Modrich, P. (2001) *Proc. Natl. Acad. Sci. USA* **98**, 6765–6770.
- Dermic, D., Halupecki, E., Zahradka, D. & Petranovic, M. (2005) *Res. Microbiol.* **156**, 304–311.
- Fitzpatrick, D. A., Creevey, C. J. & McInerney, J. O. (2005) *J. Mol. Evol.* **61**, 90–98.
- Finlay, B. B. & Falkow, S. (1997) *Microbiol. Mol. Biol. Rev.* **61**, 136–169.
- Murray, G. L., Attridge, S. R. & Morona, R. (2003) *Mol. Microbiol.* **47**, 1395–1406.
- West, N. P., Sansonetti, P., Mounier, J., Exley, R. M., Parsot, C., Guadagnini, S., Prevost, M. C., Prochnicka-Chalufour, A., Delepiere, M., Tanguy, M., *et al.* (2005) *Science* **307**, 1313–1317.
- Svanborg, C., Godaly, G. & Hedlund, M. (1999) *Curr. Opin. Microbiol.* **2**, 99–105.
- Pece, S., Giuliani, G., Di Leo, A., Fumarola, D., Antonaci, S. & Jirillo, E. (1997) *Recent Prog. Med.* **88**, 237–241.
- Roux, A., Beloin, C. & Ghigo, J. M. (2005) *J. Bacteriol.* **187**, 1001–1013.
- Nikaido, H. (1996) in *Escherichia coli and Salmonella typhimurium: Cellular and Molecular Biology*, eds. Neidhardt, F. C., Curtiss, R., Ingraham, J. L., Lin, E. C. C., Low, K. B., Magasanik, B., Reznikoff, W. S., Riley, M., Schaecter, M. & Umberger, H. E. (Am. Soc. Microbiol., Washington, DC), pp. 29–47.
- Snyder, J. A., Haugen, B. J., Buckles, E. L., Lockatell, C. V., Johnson, D. E., Donnenberg, M. S., Welch, R. A. & Mobley, H. L. (2004) *Infect. Immun.* **72**, 6373–6381.
- Batchelor, E., Walther, D., Kenney, L. J. & Goulian, M. (2005) *J. Bacteriol.* **187**, 5723–5731.
- Bernardini, M. L., Sanna, M. G., Fontaine, A. & Sansonetti, P. J. (1993) *Infect. Immun.* **61**, 3625–3635.
- Heidrich, C., Templin, M. F., Ursinus, A., Merdanovic, M., Berger, J., Schwarz, H., de Pedro, M. A. & Holtje, J. V. (2001) *Mol. Microbiol.* **41**, 167–178.
- Katayama, T., Takata, M. & Sekimizu, K. (1997) *Mol. Microbiol.* **26**, 687–697.
- Nurse, P., Levine, C., Hassing, H. & Mariani, K. J. (2003) *J. Biol. Chem.* **278**, 8653–8660.
- Weinberg, E. D. (1978) *Microbiol. Rev.* **42**, 45–66.
- Janke, B., Dobrindt, U., Hacker, J. & Blum-Oehler, G. (2001) *FEMS Microbiol. Lett.* **199**, 61–66.
- Woodrow, G. C., Young, I. G. & Gibson, F. (1979) *Biochim. Biophys. Acta* **582**, 145–153.
- Zhang, Z. & Gerstein, M. (2003) *J. Biol.* **2**, 11.1–11.4.
- Wagner, G. P., Fried, C., Prohaska, S. J. & Stadler, P. F. (2004) *Mol. Biol. Evol.* **21**, 2116–2121.
- Lee, W. & Chen, S. L. (2002) *BioTechniques* **33**, 1334–1341.
- Thompson, J. D., Higgins, D. G. & Gibson, T. J. (1994) *Nucleic Acids Res.* **22**, 4673–4680.
- Sawyer, S. (1989) *Mol. Biol. Evol.* **6**, 526–538.
- Felsenstein, J. (1989) *Cladistics* **5**, 164–166.
- Yang, Z. (1997) *Comput. Appl. Biosci.* **13**, 555–556.