

## The RecA Protein as a Model Molecule for Molecular Systematic Studies of Bacteria: Comparison of Trees of RecAs and 16S rRNAs from the Same Species

Jonathan A. Eisen

Department of Biological Sciences, Stanford University, Stanford, CA 94305-5020, USA (email: jeisen@leland.stanford.edu)

Received: 1 July 1995 / Accepted: 25 July 1995

**Abstract.** The evolution of the RecA protein was analyzed using molecular phylogenetic techniques. Phylogenetic trees of all currently available complete RecA proteins were inferred using multiple maximum parsimony and distance matrix methods. Comparison and analysis of the trees reveal that the inferred relationships among these proteins are highly robust. The RecA trees show consistent subdivisions corresponding to many of the major bacterial groups found in trees of other molecules including the  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$  proteobacteria, cyanobacteria, high-GC gram-positives, and the *Deinococcus-Thermus* group. However, there are interesting differences between the RecA trees and these other trees. For example, in all the RecA trees the proteins from gram-positive species are not monophyletic. In addition, the RecAs of the cyanobacteria consistently group with those of the high-GC gram-positives. To evaluate possible causes and implications of these and other differences phylogenetic trees were generated for small-subunit rRNA sequences from the same (or closely related) species as represented in the RecA analysis. The trees of the two molecules using these equivalent species-sets are highly congruent and have similar resolving power for close, medium, and deep branches in the history of bacteria. The implications of the particular similarities and differences between the trees are discussed. Some of the features that make RecA useful for molecular systematics and for studies of protein evolution are also discussed.

**Key words:** RecA — Small-subunit rRNA — Bacteria — Molecular evolution — Molecular systematics — Congruence — Protein — Phylogeny — Gram-positive — Proteobacteria

### Introduction

Molecular systematics has become the primary way to determine evolutionary relationships among microorganisms because morphological and other phenotypic characters are either absent or change too rapidly to be useful for phylogenetic inference (Woese 1987). Not all molecules are equally useful for molecular systematic studies and the molecule of choice for most such studies of microorganisms has been the small subunit of rRNA (SS-rRNA). Comparisons of SS-rRNA sequences have revolutionized the understanding of the diversity and phylogenetic relationships of all organisms, and in particular those of microorganisms (Fox et al. 1980; Olsen 1988; Olsen et al. 1994; Pace et al. 1986; Sogin 1989; Woese 1987, 1991). Some of the reasons that SS-rRNA sequence comparisons have been so useful include: SS-rRNAs are present in, and have conserved sequence, structure, and function among, all known species of free-living organisms as well as mitochondria and chloroplasts (Pace et al. 1986; Woese 1987); genes encoding SS-rRNAs are relatively easy to clone and sequence even from uncharacterized or unculturable species (Eisen et al. 1992; Lane et al. 1985; Medlin et al. 1988; Olsen et al. 1986; Weisburg et al. 1991); the conservation of some regions of primary structure and large sections of secondary structure aids alignment of SS-rRNA sequences between species (Woese 1987); the evolutionary substitution rate between species varies greatly within the molecule, allowing for this one molecule to be used to infer relationships among both close and distant relatives (Pace et al. 1986; Woese 1987); and it is generally considered unlikely that SS-rRNA genes have undergone lateral transfers between species (Pace et al. 1986), thus the history of SS-rRNA genes should cor-

respond to the history of the species from which they come. The accumulating database of SS-rRNA sequences, which now includes over 3,000 complete or nearly complete sequences (Maidak et al. 1994), provides an extra incentive to focus on this molecule.

Despite the advantages and successes of using SS-rRNA sequences to determine microbial phylogenetic relationships, there are potential problems with relying on only SS-rRNA-based phylogenetic trees (e.g., Hasegawa and Hashimoto 1993; Rothschild et al. 1986). First, there are some characteristics of SS-rRNA genes that may lead to trees based on them being inaccurate including overestimation of the relatedness of species with similar nucleotide frequencies (such as could occur in unrelated thermophiles) (Embley et al. 1993; Vawter and Brown 1993; Viale et al. 1994; Weisburg et al. 1989b; Woese et al. 1991), nonindependence of substitution patterns at different sites (Gutell et al. 1994; Schoeniger and Von Haeseler 1994) variation in substitution rates between lineages (e.g., Wolfe et al. 1992; Bruns and Szaro 1992; Nickrent and Starr 1994) and ambiguities in alignments between distantly related taxa. Even if the trees inferred from the SS-rRNA genes accurately reflect the evolutionary history of these genes, they might not accurately reflect the history of the species as a whole. For example, lateral transfers between species might cause the genomes of some species to have mosaic evolutionary histories. Although it is unlikely that SS-rRNAs have been stably transferred between species (see above), other genes may have been. Therefore, to understand the history of entire genomes, and to better understand the extent of mosaicism within species, it is important to compare and contrast the histories of different genes from the same species. Finally, since SS-rRNA genes are present in multiple copies in many bacteria (Jinks-Robertson and Nomura 1987; Nomura et al. 1977), it is possible that the genes being compared between species are paralogous, not orthologous. This could cause the gene trees to be different from the species trees. For these and other reasons, researchers interested in microbial systematics have begun to compare and contrast the relationships of other molecules with those of the SS-rRNA. The choice of which additional molecule to use is a difficult one. Many potential candidates have arisen and each has its advantages. Examples include HSP70 (Boorstein et al. 1994, Gupta et al. 1994, Rensing and Maier 1994), GroEL (Viale et al. 1994), EF-TU (Ludwig et al. 1994; Delwiche et al. 1995), ATPase- $\beta$ -subunit (Ludwig et al. 1994), 23S rRNA (Ludwig et al. 1992), and RNA polymerases (Klenk and Zillig 1994). Another potential choice is RecA.

The RecA protein of *Escherichia coli* is a small (352 aa) yet versatile protein with roles in at least three distinct cellular processes: homologous DNA recombination, SOS induction, and DNA-damage-induced mutagenesis (Kowalczykowski et al. 1994). This diversity of genetic functions is paralleled by multiple biochemi-

cal activities including DNA binding (double- and single-stranded), pairing and exchange of homologous DNA, ATP hydrolysis, and coproteolytic cleavage of the LexA,  $\lambda$ cI, and UmuD proteins (Kowalczykowski et al. 1994). It has been 30 years since the isolation of the first *recA* mutants in *E. coli* (Clark and Margulies 1965) and 15 years since the sequencing of the corresponding *recA* gene (Sancar et al. 1980; Horii et al. 1980). In that time, studies of the wild-type and mutant RecA proteins and genes have yielded a great deal of information about the structure-function relationships of the protein, as well as about the general mechanisms of homologous recombination (e.g., Clark and Sandler 1994; Kowalczykowski 1991; Roca and Cox 1990). Such studies have been facilitated greatly by the publication of the crystal structure of the *E. coli* RecA protein alone and bound to ADP (Story and Steitz 1992; Story et al. 1992).

Genes encoding proteins with extensive amino-acid sequence similarity to the *E. coli* RecA have been cloned and sequenced from many other bacterial species. Included among these are complete open reading frames from many of the major bacterial phyla as well as an open reading frame from the nucleus of *Arabidopsis thaliana* which encodes a protein that functions in the chloroplast (Table 1). Partial open reading frames are available from many additional bacterial species. The high levels of sequence similarity, even between proteins from distantly related taxa, and the demonstration that many of the functions and activities of the *E. coli* RecA are conserved in many of these other proteins (Angov and Camerini-Otero 1994; Gutman et al. 1994; Roca and Cox 1990; Wetmur et al. 1994), suggest that these proteins are homologs of the *E. coli* RecA.

The diversity and number of species from which sequences are available make RecA a potentially useful tool for molecular systematic studies of bacteria. Previously, Lloyd and Sharp (1993) tested the utility of RecA comparisons for phylogenetic studies. They concluded that RecA comparisons were probably only useful for determining relationships among closely related bacterial species. However, they were limited by the number and diversity of RecA sequences that were available at the time. I have reanalyzed the evolution of RecA using 40 additional sequences. In this paper, analysis is presented that shows that the RecA protein is a good alternative to SS-rRNA for molecular systematic studies of all bacteria, not just closely related species. Phylogenetic trees of the 65 complete RecA protein sequences were inferred using a variety of phylogenetic methods. Statistical analysis and comparisons of trees generated by the different phylogenetic methods indicate that the RecA phylogeny is highly consistent and robust. The RecA trees are compared to trees of SS-rRNA sequences from the same or very closely related species as represented in the RecA trees. Overall, the trees of the two molecules are highly congruent. The implications of the particular similarities and differences between the RecA-based and

SS-rRNA-based trees are discussed. Some of the features of RecA that make it a potentially useful molecular chronometer are also discussed.

## Methods

**Sequences and Alignment.** All RecA sequences used in this paper were obtained from the National Center for Biotechnology Information (NCBI) databases by electronic mail (Henikoff 1993) except for those from *Methylophilus methylotrophus* (Emmerson 1995, personal communication), *Xanthomonas oryzae* (Mongkolsuk 1995, pers. comm.), *Synechococcus* sp. PCC7942 (Coleman 1995), and *Borrelia burgdorferi* (Huang 1995, pers. comm.), which were kindly provided prior to submission. Accession numbers for those in databases are given in Table 1. The amino-acid sequences of the RecA proteins were aligned both manually and with the *clustalw* multiple sequence alignment program (Thompson et al. 1994) using default parameters. The RecA alignment was used as a block and aligned with the sequence of the RadA protein from an archaea (Clark and Sandler 1994) and RecA-like proteins from eukaryotes (Ogawa et al. 1993) also using *clustalw*.

For the comparison of RecA and SS-rRNA trees, a complete or nearly complete SS-rRNA sequence was chosen to represent each species for which a complete RecA protein was available. For most of the RecA proteins, a complete SS-rRNA sequence was available from the same species. The remaining species (those for which a RecA sequence was available but for which a complete or nearly complete SS-rRNA was not) were represented by a "replacement" SS-rRNA from a different species. The choice of which replacement sequence to use was determined in one of two ways. For those RecAs for which a partial SS-rRNA was available from the same species, the complete or nearly complete SS-rRNA that was most similar to the partial sequence was used. Similarity was determined by comparisons using the Ribosomal Database Project (RDP) computer server (Maidak et al. 1994) and *blastn* searches (Altschul et al. 1990) of the NCBI databases by electronic mail (Henikoff 1993). For those RecAs for which even a partial SS-rRNA sequence was not available from the same species, a replacement was chosen from a species considered to be a close relative. A SS-rRNA was not used to represent the *Shigella flexneri* RecA because this protein was identical to the *E. coli* RecA. For the majority of the SS-rRNA phylogenetic analysis, only one SS-rRNA was used to represent the two RecAs from *Myxococcus xanthus*. For some of the analysis an additional SS-rRNA from a close relative of *M. xanthus* was also included. The SS-rRNA sequences used and the species from which they come are listed in Table 1. The SS-rRNA sequences were obtained already aligned from the RDP (Maidak et al. 1994), with the exception of those from *Corynebacterium glutamicum* and *Anabaena* sp. PCC7120 which were obtained from the NCBI and were aligned to the other sequences manually. Entry names and numbers are listed in Table 1.

**Phylogenetic Trees.** Phylogenetic trees were generated from the sequence alignments using computer algorithms implemented in the PHYLIP (Felsenstein 1989), PAUP (Swofford 1991), and GDE (Smith 1991; Smith et al. 1994) computer software packages. Trees of the RecA sequences were generated using two parsimony methods (the *protpars* program in PHYLIP and the *heuristic search* algorithm of PAUP) and three distance methods (the least-squares method of De Soete [1983] as implemented in GDE, and the Fitch-Margoliash [Fitch and Margoliash 1967] and neighbor-joining methods [Saitou and Nei 1987] as implemented in PHYLIP). Trees of the SS-rRNA sequences were generated using one parsimony method (the *dnapars* algorithm of PHYLIP) and the same three distance methods as used for the RecA trees. For the trees generated by the *protpars*, *dnapars*, Fitch-Margoliash, and neighbor-joining methods, 100 bootstrap replicates were conducted by the method of Felsenstein (1985) as implemented in PHYLIP.

For the distance-based phylogenetic methods listed above, estimated evolutionary distances between each pair of sequences were calculated for input into the tree-reconstruction algorithms. Pairwise distances between RecA proteins were calculated using the *protdist* program of PHYLIP and the PAM matrix-based distance correction (Felsenstein 1989). Pairwise distances between SS-rRNA sequences were calculated in two ways: the method of Olsen (1988) (as implemented by the *count* program of GDE) was used for the trees generated by the De Soete method; and the two-parameter model of Kimura (1980) (as implemented by the *dnadist* program of PHYLIP) was used for the Fitch-Margoliash and neighbor-joining trees.

Regions of the alignments for which homology of residues could not be reasonably assumed were excluded from the phylogenetic analysis. For the SS-rRNA trees, the alignment of SS-rRNA sequences was extracted from an alignment of thousands of sequences in the RDP database (Maidak et al. 1994). This RDP alignment was generated using both primary and secondary structures as a guide to assist in the assignment of homology (Maidak et al. 1994). Therefore it was assumed that the aligned regions were likely homologous. Nevertheless, regions of high sequence variation (as determined by a 50% consensus mask) were excluded from the phylogenetic analysis since these regions are perhaps most likely to contain nonhomologous residues. The SS-rRNA alignment and a list of the 1,061 alignment positions used are available on request. For the RecA analysis, the assignment of homology was based only on similarity of primary structure as determined by the *clustalw* program. Regions of ambiguity in the alignment were considered to potentially include nonhomologous residues and were excluded from the phylogenetic analysis. Such regions were identified by comparing alignments generated by the *clustalw* program using a variety of alignment parameters. Parameters varied included scoring matrices (PAM, BLOSUM, and identity matrices were used) and gap opening and extension penalties. Alignments were compared by eye to detect differences and those regions that contained different residues in the different alignments were considered ambiguous.

**Character States and Changes.** Analysis of character states and changes over evolutionary history was done using the MacClade 3.04 program (Maddison and Maddison 1992). For each alignment position, all unambiguous substitutions as well as all unambiguous nonconservative substitutions were counted. Nonconservative substitutions were defined as amino-acid changes that were not within the following groups: (V-I-L-M), (F-W-Y), (D-E-N-Q), (K-R), (G-A), (S-T).

**Computer Programs.** GDE, PHYLIP, and *clustalw* were obtained by anonymous FTP from the archive of the Biology Department at the University of Indiana (<ftp.bio.indiana.edu>). PAUP was obtained from David Swofford. GDE, PHYLIP, and *clustalw* were run on a Sparc10 workstation and MacClade and PAUP on a Power Macintosh 7100/66. Unless otherwise mentioned, all programs were run with default settings.

## Results and Discussion

The potential of using RecA for phylogenetic studies of bacteria was first addressed by Lloyd and Sharp (1993). In a detailed analysis of the evolution of *recA* genes from 25 species of bacteria, they showed that phylogenetic trees of RecA proteins appeared to be reliable for determining relationships among closely related bacterial species. Specifically, for the proteobacteria, the branching patterns of RecA proteins were highly congruent to branching patterns of SS-rRNA genes from similar species. However, the RecA and SS-rRNA trees were not highly congruent for relationships between sequences

**Table 1.** RecA and SS-rRNA sequences

Species (by phylum)	Abbrev.	RecA Seq.	Size (aa)	SS-rRNA Seq. <sup>a,b</sup>	RecA Refs.
<b>Proteobacteria</b>					
<i>Acetobacter polyoxogenes</i>	Act.po	D13183	348	ABA.PASTER*	(Tayama et al. 1993)
<i>Acidiphilium facilis</i>	Acd.f	D16538	354	ACDP.FACI2	(Inagaki et al. 1993)
<i>Acinetobacter calcoaceticus</i>	Acn.c	L26100	349	ACN.CALCOA	(Gregg-Jolly and Ornston 1994)
<i>Agrobacterium tumefaciens</i>	Ag.t	L07902	363	AG.TUMEFAC	(Wardhan et al. 1992)
<i>Azotobacter vinelandii</i>	Az.v	S96898	349	F.LUTESCEN*	(Venkatesh and Das 1992)
<i>Bordetella pertussis</i>	Bd.p	X53457	352	BRD.PERTUS	(Favre et al. 1991; Favre and Viret 1990)
<i>Brucella abortus</i>	Br.a	L00679	360	BRU.ABORTS	(Tatum et al. 1993)
<i>Burkholderia cepacia</i> <sup>d</sup>	Bu.c	D90120	347	BUR.CEPACI	(Nakazawa et al. 1990)
<i>Campylobacter jejuni</i>	Ca.j	U03121	343	CAM.JEJUNI	(Guerry et al. 1994)
<i>Enterobacter agglomerans</i> <sup>e</sup>	En.a	P33037	354	ER.HERBICO	(Rappold and Klingmueller 1993)
<i>Erwinia carotovora</i>	Er.c	X55554	342	ER.CAROTOV	(Zhao and McEntee 1990)
<i>Escherichia coli</i>	Es.c	V00328	353	E.COLI	(Hori et al. 1980; Sancar et al. 1980)
<i>Haemophilus influenzae</i>	Ha.i	L07529	354	H.INFLUENZ	(Zulty and Barcak 1993)
<i>Helicobacter pylori</i>	He.p	Z35478	347	HLB.PYLOR3	(Haas 1994)
<i>Legionella pneumophila</i>	Le.p	X55453	348	LEG.PNEUMO	(Zhao and Dreyfus 1990)
<i>Magnetospirillum magnetotacticum</i> <sup>f</sup>	Ma.m	X17371	344	MAG.MAGNE2	(Berson et al. 1990)
<i>Methylobacillus flagellatum</i>	Mb.f	M35325	344	MBS.FLAGEL	(Gomelsky et al. 1990)
<i>Methylomonas clara</i>	Mm.c	X59514	342	MLM.METHYL*	(Ridder et al. 1991)
<i>Methylophilus methylotrophus</i>	Mp.m	unpub.	342	MLP.METHY1	(Emmerson 1995, personal communication)
<i>Myxococcus xanthus</i> 1	Mx.1	L40367	342	MYX.XANTHU	(Inouye 1995, personal communication)
<i>Myxococcus xanthus</i> 2	Mx.2	L40368	358	n/a <sup>c</sup>	(Inouye 1995, personal communication)
<i>Neisseria gonorrhoeae</i>	Ne.g	X17374	348	NIS.GONORR	(Fyfe and Davies 1990)
<i>Proteus mirabilis</i>	Pr.m	X14870	355	ARS.NASONI*	(Akaboshi et al. 1989)
<i>Proteus vulgaris</i>	Pr.v	X55555	325	P.VULGARIS	(Zhao and McEntee 1990)
<i>Pseudomonas aeruginosa</i>	Ps.a	X52261	346	PS.AERUGIN	(Sano and Kageyama 1987)
<i>Pseudomonas fluorescens</i>	Ps.f	M96558	352	PS.FLAVESC*	(De Mot et al. 1993)
<i>Pseudomonas putida</i>	Ps.p	L12684	355	PS.PUTIDA	(Luo et al. 1993)
<i>Rhizobium leguminosarum</i> bv. <i>phaseoli</i>	Rz.p	X62479	360	RHB.LEGUM6	(Michiels et al. 1991)
<i>Rhizobium leguminosarum</i> bv. <i>viciae</i>	Rz.l	X59956	351	RHB.LEGUM8	(Selbitschka et al. 1991)
<i>Rhizobium meliloti</i>	Rz.m	X59957	348	RHB.MELIL2	(Selbitschka et al. 1991)
<i>Rhodobacter capsulatus</i>	Rh.c	X82183	355	RB.CAPSUL2	(Fernandez de Henestrosa 1994)
<i>Rhodobacter sphaeroides</i>	Rh.s	X72705	343	RB.SPHAER2	(Calero et al. 1994)
<i>Rickettsia prowazekii</i>	Ri.p	U01959	340	RIC.PROWAZ	(Dunkin and Wood 1994)
<i>Serratia marcescens</i>	Se.m	M22935	354	SER.MARCES	(Ball et al. 1990)
<i>Shigella flexneri</i>	Sh.f	X55553	353	n/a	(Zhao and McEntee 1990)
<i>Thiobacillus ferrooxidans</i>	Tb.f	M26933	346	THB.CALDUS*	(Ramesar et al. 1989)
<i>Vibrio anguillarum</i>	Vi.a	M80525	348	V.ANGUILLA	(Gammie and Crosa 1991; Tolmasky et al. 1992)
<i>Vibrio cholerae</i>	Vi.c	U10162	354	V.CHOLERA	(Margraf et al. 1995; Strocher et al. 1994)
<i>Xanthomonas oryzae</i>	Xa.o	unpub.	355	XAN.ORYZAE	(Mongkolsuk 1995, personal communication)
<i>Yersinia pestis</i>	Ye.p	X75336	356	YER.PESTIS	(Kryukov et al. 1993)
<b>Gram Positives</b>					
<i>Acholeplasma laidlawii</i>	Acp.l	M81465	331	ACP.LAIDLA	(Dybvig and Woodard 1992)
<i>Bacillus subtilis</i>	Ba.s	X52132	347	B.SUBTILIS	(Stranathan et al. 1990)
<i>Corynebacterium glutamicum</i>	Co.g	X77384	376	Z46753	(Billman-Jacobe 1994; Kerins et al. 1994)
<i>Lactococcus lactis</i>	La.l	M88106	365	LCC.LACTIS	(Duwat et al. 1992a)
<i>Mycobacterium leprae</i>	Myb.l	X73822	711	MYB.LEPRAE	(Davis et al. 1994)
<i>Mycobacterium tuberculosis</i>	Myb.t	X58485	790	MYB.TUBER2	(Davis et al. 1991)
<i>Mycoplasma mycoides</i>	Myp.m	L22073	345	M.MYCOIDES	(King et al. 1994)
<i>Mycoplasma pulmonis</i>	Myp.p	L22074	339	M.PULMONIS	(King et al. 1994)
<i>Staphylococcus aureus</i>	Sta.a	L25893	347	STP.AUREUS	(Bayles et al. 1994)
<i>Streptococcus pneumoniae</i>	Stc.p	Z17307	388	STC.SALIVA*	(Martin et al. 1992)
<i>Streptomyces ambofaciens</i>	Sm.a	Z30324	372	STM.AMBOFA	(Aigle et al. 1994)
<i>Streptomyces lividans</i>	Sm.l	X76076	374	STM.LIVIDA	(Nussbaumer and Wohlleben 1994)
<i>Streptomyces violaceus</i> <sup>g</sup>	Sm.v	U04837	377	STM.COELI3*	(Yao and Vining 1994)
<b>Cyanobacteria/Chloroplasts</b>					
<i>Arabidopsis thaliana</i>	Ar_C	M98039	439	NICO.TAB_C*	(Binet et al. 1993; Cerutti et al. 1992)
<i>Anabaena variabilis</i>	An.v	M29680	358	X59559*	(Owtrim and Coleman 1989)
<i>Synechococcus</i> sp. PCC7942	Sy.79	unpub.	361	PHRM.MINUT*	(Coleman 1995, personal communication)
<i>Synechococcus</i> sp. PCC7002	Sy.70	M29495	348	SYN.6301*	(Murphy et al. 1987, 1990)

Table 1. Continued

Species (by phylum)	Abbrev.	RecA Seq.	Size (aa)	SS-rRNA Seq. <sup>a,b</sup>	RecA Refs.
Deinococcus-Thermus Group					
<i>Deinococcus radiodurans</i> <sup>h</sup>	De.r	U01876	363	D.RADIODUR	(Gutman et al. 1994)
<i>Thermus aquaticus</i>	Th.a	L20095	340	T.AQUATICU	(Angov and Camerini-Otero 1994; Wetmur et al. 1994)
<i>Thermus thermophilus</i>	Th.t	D13792	340	T.THMOPHL	(Kato and Kuramitsu 1993; Wetmur et al. 1994)
Chlamydia/Planctomyces					
<i>Chlamydia trachomatis</i>	Ch.t	U16739	352	CLM.TRACHO	(Larsen 1994; Zhang et al. 1994)
Spirochaetes					
<i>Borrelia burgdorferi</i>	Bo.b	unpub.	365	BOR.BURGDO	(Huang 1995, personal communication)
Bacteroides					
<i>Bacteroides fragilis</i>	Bct.f	M63029	318	BAC.FRAGIL	(Goodman and Woods 1990)
Thermophilic O <sub>2</sub> Reducers					
<i>Aquifex pyrophilus</i>	Aq.p	L23135	348	AQU.PYROPH	(Wetmur et al. 1994)
Thermotogales					
<i>Thermotoga maritima</i>	Tg.m	L23425	356	TT.MARITIM	(Wetmur et al. 1994)

<sup>a</sup> Names refer to Ribosomal Database Project entries (Maidak et al. 1994). Numbers are Genbank entries

<sup>b</sup> The SS-rRNA sequences that come from a different species than the RecA sequences are indicated by an asterisk\*. The species are ABA.PASTER (*Acetobacter pasteurianus*), F.LUTESCEN (\**Flavobacterium* lutescens), MLM.METHYL (*Methylomonas methylovora*), ARS.NASONI (*Arsenophonus nasoniae*), PS.FLAVESC (*Pseudomonas flavescens*), THB.CALDUS (*Thiobacillus caldus*), STM.COELI3 (*Streptomyces coelicolor*), STC.SALIVA (*Streptococcus salivarius*), NICO.TAB\_C (*Nicotiana tabacum* chloroplast),

X59559 (*Anabaena* sp. PCC7120), PHRM.MINUT (*Phormidium minutum*), and SYN.6301 (*Synechococcus* sp. PCC 6301)

<sup>c</sup> For most of the analyses only one SS-rRNA was used for the two *M. xanthus* RecAs. For some analyses the SS-rRNA of *Cystobacter fuscus* (CYS.FUSCUS) was also used

<sup>d</sup> Also known as *Pseudomonas cepacia*

<sup>e</sup> Also known as *Erwinia herbicola*

<sup>f</sup> Also known as *Aquaspirillum magnetotacticum*

<sup>g</sup> Also known as *Streptomyces venezuelae*

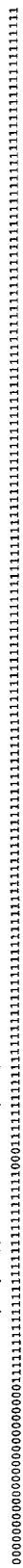
<sup>h</sup> Also known as *Micrococcus radiodurans*

from more distantly related species. Lloyd and Sharp (1993) concluded that this was due to a low resolution of the deep branches in the RecA tree. However, this low resolution of deep branches could have been due to poor representation of certain taxa in their sample set. Of the *recA* sequences available at the time, only six were from species outside the proteobacteria. The diversity as well as the number of *recA* sequences available has increased greatly since Lloyd and Sharp's study (see Table 1). Therefore, I have reanalyzed the evolution of *recA* including these additional sequences with a specific focus on determining whether *recA* comparisons can provide reasonable resolution of moderate-to-deep branches in the phylogeny of bacteria. The analysis presented here focuses on amino-acid comparisons for two reasons. First, for highly conserved proteins such as RecA, it is likely that amino-acid trees will be less biased by multiple substitutions at particular sites and base-composition variation between species than trees of the corresponding nucleotide sequences (Hasegawa and Hashimoto 1993; Viale et al. 1994). In addition, Lloyd and Sharp (1993) presented evidence suggesting that DNA-level comparisons of the *recA* genes between distantly related taxa might be misleading.

#### Alignment of RecA Sequences

An alignment of the sequences of the complete RecA proteins is shown in Fig. 1. Aligning sequences is an

integral part of any molecular systematic study because each aligned position is assumed to include only homologous residues from the different molecules. Assignment of homology, as represented by the sequence alignments, can be highly controversial, and differences in alignments can cause significant differences in phylogenetic conclusions (e.g., Gatesy et al. 1993; Lake 1991). To limit such problems, regions for which homology of residues cannot be unambiguously assigned should be excluded from phylogenetic analysis. Thus for a molecule to be useful for molecular systematic studies, alignments between species should be relatively free of ambiguities. This is one of the main advantages of using SS-rRNA sequences for phylogenetic analysis. Assignment of homology for SS-rRNA sequences can be aided by alignment of both primary and secondary structures (Woese 1987). In addition, regions of high primary structural conservation that are interspersed throughout the molecule help align less conserved regions. Since RecA is a highly conserved protein, it has the potential to be useful for phylogenetics because the assignment of homology should be relatively unambiguous (Lloyd and Sharp 1993). Regions of ambiguity in the RecA alignment shown in Fig. 1 were determined by comparing this alignment to those generated using different alignment parameters (see Methods). Regions of the alignment were considered to be ambiguous if they contained different residues in the different alignments, as suggested by Gatesy et al. (1993). Overall, the majority of the alignment was determined to be free of ambiguities and



**Fig. 1.** Alignment of complete RecA sequences. The alignment was generated using the *clustalw* multiple sequence alignment program. *Dashes* (-) represent alignment gaps. Three insertions that are present in only one sequence each (Myb.1, Myb.1, and Tg.m) and the first 80 aa of the *A. thaliana* protein are left out for space reasons and are indicated by a ·. Conservation of alignment positions as determined by the *clustalw* program is indicated by \* (identical aa in all) and · (similar aa in all). The alignment positions used in phylogenetic analysis are indicated by the sequence mask (1 = used, 0 = not used). Sequence abbreviations are described in Table 1.

[illegible]

Fig. 1. Continued.



thus can be used with confidence for the phylogenetic analysis. The four regions of ambiguity (the C- and N-termini [corresponding to *E. coli* amino acids 1–7 and 320–352] and two short regions [corresponding to *E. coli* amino acids 36–37 and 231–236]) were excluded from the phylogenetic analysis. The 313 alignment positions used are indicated by the sequence mask shown in Fig. 1.

Another potential source of error in phylogenetic reconstruction from sequences lies in assigning a weight to give insertion or deletion differences (indels) between species. Other than in the C- and N-terminal regions, there are few indels in the RecA alignment (see Fig. 1). Most of the indels are in regions of ambiguous alignment as identified above, and thus were not included in the phylogenetic analysis. The phylogenetic results were not affected whether or not the regions with the few remaining indels were included (data not shown). Of the indels in regions of unambiguous alignment most are isolated (in only one species) and only one amino acid in length. There are two very large indels—one in each of the *Mycobacterium* RecAs. These are protein introns that are removed by post-translational processes (Davis et al. 1991; Davis et al. 1994). There is a four aa indel in the *Thermotoga maritima* RecA (see Fig. 1). The only indels that have obvious phylogenetic relevance are the one aa gaps in the cyanobacterial and the *A. thaliana* RecAs—all at the same position—*E. coli* position 53 (discussed below).

Another aspect of the RecA alignment that is relevant to molecular systematics is the degree of conservation of different alignment positions. I have used the RecA phylogeny and parsimony character-state analysis to characterize the patterns of amino-acid substitutions at different sites of the molecule (see Methods). The number of substitutions varies a great deal across the molecule. The number of total substitutions ranges from 0 (at 58 positions) to 38 (at one position) with a mean of 9.4. The number of nonconservative substitutions varies from 0 (at 111 positions) to 27 (at one position) with a mean of 4.8. The variation in the substitution patterns suggests that RecA comparisons may have phylogenetically useful information at multiple evolutionary distances.

### Generation of Phylogenetic Trees

To examine the utility of the RecA comparisons for molecular systematics, the RecA trees were compared to trees of the same species based on studies of other molecules. Such a comparison is useful for a few reasons. First, congruence among trees of different molecules indicates both that the genomes of the species are not completely mosaic and that the molecular systematic techniques being used are reliable (Miyamoto and Fitch 1995). Differences in the branching patterns between

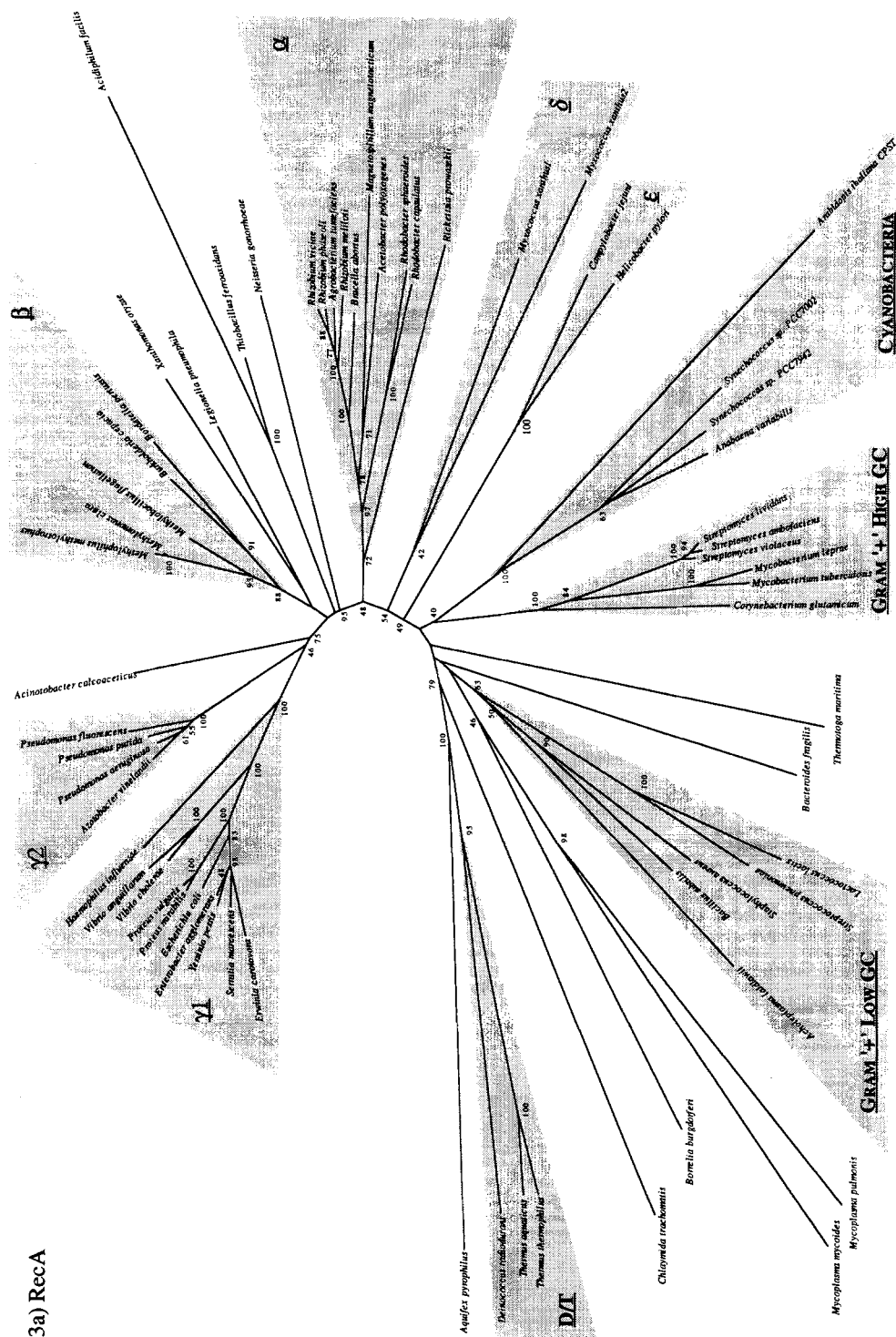
trees of different molecules can help identify genetic mosaicism, unusual evolutionary processes, or inaccuracies in one or both of the trees. Differences in resolution and significance of particular branches can help identify which molecules are useful for specific types of phylogenetic comparisons. Since differences in species sampled have profound effects on tree generation (e.g., Lecointre et al. 1993), to best compare the phylogenetic resolution of trees of different molecules the analysis should include sequences from the same species. Fortunately, SS-rRNA sequences were available for most of the species represented in the RecA data set. Therefore it was possible to generate SS-rRNA trees for essentially the same species set as represented in the RecA trees. For those species for which RecA sequences were available but SS-rRNA sequences were not, SS-rRNA sequences were used from close relatives (see Methods). A list of the sequences used is in Table 1.

Phylogenetic trees of the RecAs and SS-rRNAs were generated from the sequence alignments using multiple phylogenetic techniques (see Methods). The trees were generated without an outgroup and thus can be considered unrooted. However, since rooting of trees is helpful for a variety of reasons, a root was determined for both RecA and SS-rRNA trees. In both cases the root was determined to be the sequence from *A. pyrophilus*. For the SS-rRNA trees, this rooting was chosen because analyses of sequences from all three kingdoms of organisms indicate that the deepest-branching bacterial SS-rRNA is that of *A. pyrophilus* (Burggraf et al. 1992; Pitulle et al. 1994). Although it seems reasonable to assume that the deepest branching bacterial RecA would also be that of *A. pyrophilus*, if there have been lateral transfers or other unusual evolutionary processes, the RecA trees could be rooted differently from the SS-rRNA trees. Therefore, the rooting of the RecA sequences was tested by constructing trees using likely RecA homologs from Archaea and eukaryotes as outgroups (see Methods). In both neighbor-joining and *prot-pars* trees, the deepest-branching bacterial protein was that of *A. pyrophilus* (not shown). However, the alignments of the RecAs with the Archaeal and eukaryotic RecA-like proteins include many regions of ambiguity. Therefore only 140 alignment positions were used in this analysis and the trees showed little resolution within the bacteria. In addition, the bootstrap values for the deep branching of the *A. pyrophilus* RecA were low (less than 30 in all cases). Thus although the rooting of the RecA trees to the *A. pyrophilus* sequence is reasonable it should be considered tentative. The rooting will likely be better resolved as more sequences become available of the Archaeal and eukaryote RecA-like proteins.

The analysis and comparison of the phylogenetic trees focused on a few specific areas. First, bootstrap values were used to get an estimate of the degree to which the inferred branching patterns reflect the characteristics of the entire molecule. In addition, since phylogenetic



A quick glance at the trees in Figs. 2 and 3 shows that the patterns for each molecule are highly robust (there is high resolution in the consensus trees) and that the pat-



**Fig. 3.** Fitch-Margoliash Trees for RecA (A) and SS-rRNA (B). Trees were generated from the multiple sequence alignments by the method of Fitch and Margoliash. Regions of ambiguous alignment and indels were excluded from the analysis (see Methods). For the RecA tree, distances were calculated using the *protdist* program of PHYLIP with a PAM-matrix based distance correction. For the SS-rRNA tree, distances were calculated using the *dnadist* program of PHYLIP and the Kimura-2-parameter distance correction. Consensus clades representing groups found in all phylogenetic methods are *highlighted*. Branch lengths and *scale bars* correspond to estimated evolutionary distance. Bootstrap values when over 40 are indicated.

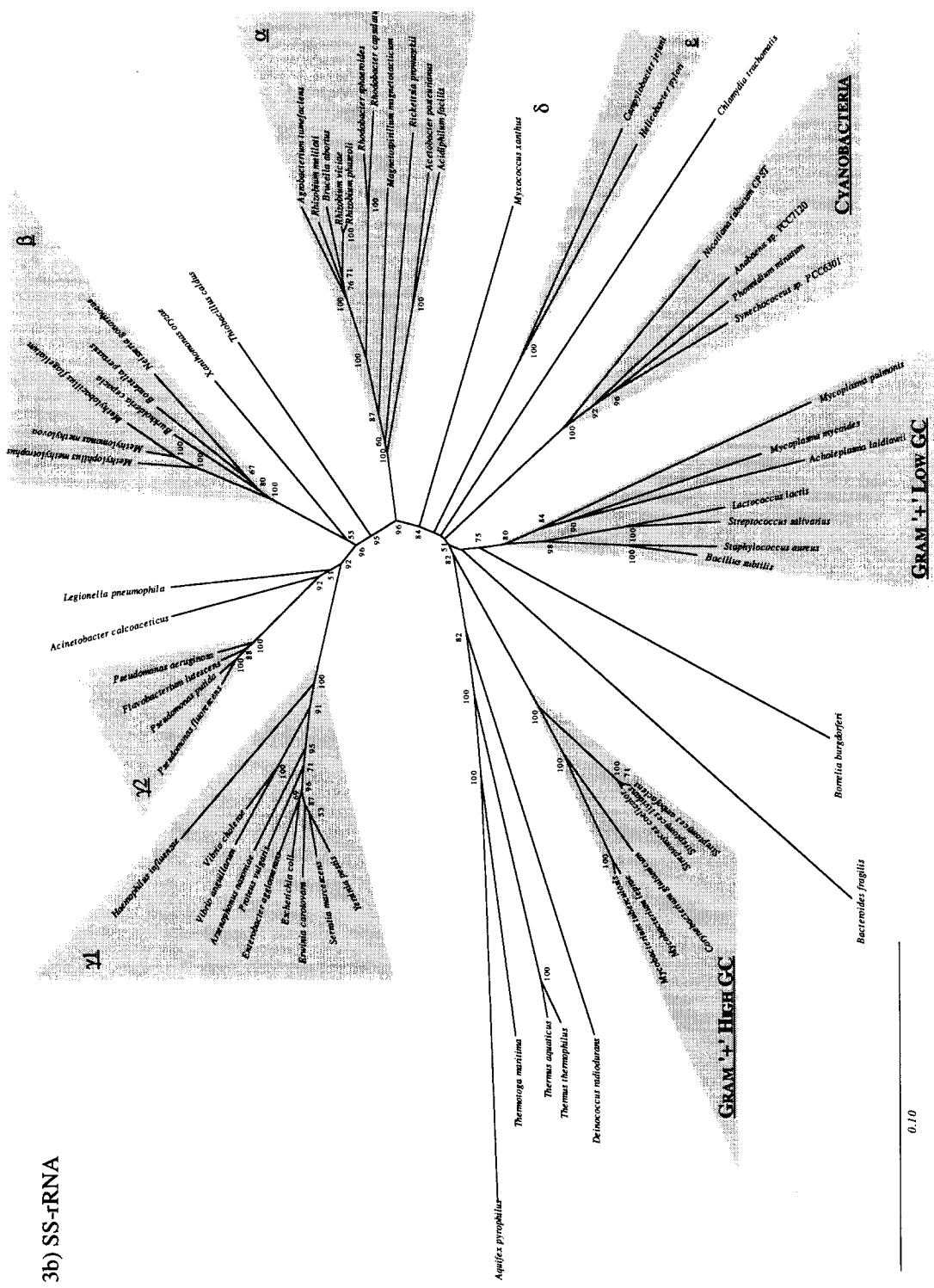


Fig. 3. Continued.

**Table 2.** Consensus phylogenetic groups

Clade	Species in RecA consensus clade <sup>2</sup>	Comparable consensus clade for SS-rRNA? <sup>3,4</sup>	Bootstrap % RecA			Bootstrap % ssRNA <sup>1</sup>		
			PPars	NJ	FM	DPars	NJ	FM
Proteobacteria-γ1	<i>Escherichia coli</i> , <i>Shigella flexneri</i> , <i>Yersinia pestis</i> , <i>Erwinia carotovora</i> , <i>Serratia marcescens</i> , <i>Enterobacter agglomerans</i> , <i>Proteus vulgaris</i> , <i>P. mirabilis</i> , <i>Vibrio cholerae</i> , <i>V. anguillarum</i> , <i>Haemophilus influenzae</i>	YES	78	91	100	100	100	100
Proteobacteria-γ2	<i>Azotobacter vinelandii</i> , <i>Pseudomonas aeruginosa</i> , <i>P. putida</i> , <i>P. fluorescens</i>	YES	100	100	100	100	100	100
Proteobacteria-γ	γ1, γ2, <i>Acinetobacter calcoaceticus</i>	YES (+ <i>L. pneumophila</i> )	33	63	75	48	85	92
Proteobacteria-β1	<i>Methylobacillus flagellatum</i> , <i>Methylomonas clara</i> , <i>Methylophilus methylotrophus</i> , <i>Burkholderia cepacia</i> , <i>Bordetella pertussis</i>	YES (+ <i>N. gonorrhoeae</i> )	74	84	88	100	100	100
Proteobacteria-β2	<i>Thiobacillus ferrooxidans</i> , <i>Acidiphilium facilis</i>	NO	100	100	100	N/A	N/A	N/A
Proteobacteria-βγ	γ, β1, β2, <i>Xanthomonas oryzae</i> , <i>Neisseria gonorrhoeae</i> , <i>Legionella pneumophila</i>	YES (– <i>A. facilis</i> )	53	86	95	90	94	95
Proteobacteria-α	<i>Rhodobacter capsulatus</i> , <i>Rh. sphaeroides</i> , <i>Rhizobium meliloti</i> , <i>R. viciae</i> , <i>R. phaseoli</i> , <i>Acetobacter polyoxogenes</i> , <i>Magnetospirillum magnetotacticum</i> , <i>Brucella abortus</i> , <i>Agrobacterium tumefaciens</i> , <i>Rickettsia prowazekii</i>	YES (+ <i>A. facilis</i> )	14	68	72	100	100	100
Proteobacteria-αβγ	α, β, γ	YES	10	57	58	93	96	96
Proteobacteria-δ	<i>Myxococcus xanthus</i> 1, <i>M. xanthus</i> 2	YES	43	71	42	N/A <sup>5</sup>	N/A	N/A
Proteobacteria-ε	<i>Campylobacter jejuni</i> , <i>Helicobacter pylori</i>	YES	100	100	100	100	100	100
Proteobacteria	γ, β, α, δ, ε	NO	14	38	49	N/A	N/A	36
Gram “+” High GC	<i>Corynebacterium glutamicum</i> , <i>Streptomyces ambofaciens</i> , <i>S. violaceus</i> , <i>S. lividans</i> , <i>Mycobacterium tuberculosis</i> , <i>M. leprae</i>	YES	97	100	100	100	100	100
Gram “+” Low GC	<i>Bacillus subtilis</i> , <i>Lactococcus lactis</i> , <i>Streptococcus pneumoniae</i> , <i>Staphylococcus aureus</i> , <i>Acholeplasma laidlawii</i>	YES (+ <i>Mycoplasmas</i> )	27	59	63	50	56	80
Mycoplasmas	<i>Mycoplasma mycoides</i> , <i>M. pulmonis</i>	YES (+ <i>A. laidlawii</i> )	88	100	98	71	88	84
Cyanobacteria	<i>Arabidopsis thaliana</i> , <i>Anabaena variabilis</i> , <i>Synechococcus</i> sp. PCC7942, <i>Syn. sp.</i> PCC7002	YES	100	96	91	100	100	100
Deinococcus-Thermus	<i>Deinococcus radiodurans</i> , <i>Thermus aquaticus</i> , <i>T. thermophilus</i>	NO	95	96	95	N/A	N/A	N/A

<sup>1</sup> Bootstrap values are shown for comparable clade<sup>2</sup> Groups found in trees generated by neighbor-joining, Fitch-Margoliash, De Soete, *protpars*, and PAUP<sup>3</sup> For those groups which have 1 or 2 additional species in the SS-rRNA tree, the extra species are listed<sup>4</sup> Groups found in trees generated by neighbor-joining, Fitch-Margoliash, De Soete and *dnaps*<sup>5</sup> Bootstraps were only calculated for trees with the one δ sequence (see Methods)

terns are similar between the two molecules. To aid comparison of the trees of the two molecules, sequences have been grouped into consensus clades based on the patterns found in the consensus trees (Fig. 2, Table 2). Clades of RecA sequences were chosen to represent previously characterized bacterial groups as well as possible. Comparable clades were determined for the SS-rRNA sequences (Table 2). The clades are named after the rRNA-based classification of most of the members of the clade

(Maidak et al. 1994). These clades are highlighted in the trees in Figs. 2 and 3. Sequences from the same or similar species are aligned in the middle in Fig. 2 to ease comparison of the two consensus trees. Besides being found in trees generated by all the phylogenetic methods used, the consensus clades have high bootstrap values for the methods in which bootstrapping was performed (Table 2). Thus the clades appear to be consistent and reliable groupings of the RecA and SS-rRNA sequences.

In the following sections, some of the implications of the similarities and differences within and between the RecA and SS-rRNA trees are discussed. The discussion has been organized by phylogenetic groups.

### Proteobacteria

The proteobacteria phylum includes most but not all the traditional gram-negative bacterial species (Stackebrandt et al. 1988). This phylum has been divided into five phylogenetically distinct groups ( $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\epsilon$ ) mostly based on SS-rRNA comparisons (Olsen et al. 1994; Rainey et al. 1993; Stackebrandt et al. 1988; Woese 1987). The available RecA sequences are heavily biased toward the proteobacteria (Table 1) and thus much of the discussion will focus on this phylum. With the species represented in this analysis, the proteobacterial RecA sequences form a monophyletic clade in all phylogenetic methods (Fig. 2). In contrast, with essentially the same species set, the proteobacterial SS-rRNA sequences do not consistently form a clade (Fig. 2, positions of *Campylobacter jejuni*, *Helicobacter pylori*, and *M. xanthus*), although they do in some of the phylogenetic methods (e.g., Fig. 3). This was surprising since the proteobacterial group was defined based on SS-rRNA comparisons (Stackebrandt et al. 1988). When additional SS-rRNA sequences are included in phylogenetic analysis, *M. xanthus*, *C. jejuni*, and *H. pylori* consistently branch with the other proteobacteria (Maidak et al. 1994; Olsen et al. 1994). The lack of resolution of the position of these species in the SS-rRNA trees versus the RecA trees was not due to using only one SS-rRNA sequence to represent the two *M. xanthus* RecAs—the same pattern was seen in trees generated including the SS-rRNA sequence from another  $\delta$  species in addition to the one from *M. xanthus* (not shown). Thus in this case the RecA trees can be considered to have higher resolution than the SS-rRNA trees since the RecA trees show a relationship between species that is only consistently detected in SS-rRNA trees with more sequences.

Subdivisions corresponding to the  $\alpha$ ,  $\beta$ ,  $\gamma$ ,  $\delta$ , and  $\epsilon$  groups are detected in both the RecA and SS-rRNA trees and the placement of species into these subdivisions is nearly the same for the two molecules (Fig. 2, Table 2). Thus the RecA comparisons support the division of the proteobacteria into these groups as well as the classification of particular species into the groups here. There are other phylogenetic patterns that are the same in the RecA and SS-rRNA trees here. Examples include the separation of the *Pseudomonas*-*Azotobacter*  $\gamma$ s ( $\gamma$ 2 here) from the *Haemophilus*, *Proteus*, and enteric  $\gamma$ s ( $\gamma$ 1 here); the monophyly of the enteric bacteria (represented here by *E. coli*, *S. flexneri*, *Erwinia carotovora*, *Enterobacter agglomerans*, and *Yersinia pestis*); the relatedness of the *Rhizobium* species, *Agrobacterium tumefaciens*, and *Brucella abortus*; the placement of *Acinetobacter cal-*

*coaceticus* into the  $\gamma$  supergroup; an affiliation between the  $\gamma$ s and the  $\beta$ s; and the grouping of *Legionella pneumophila*, *Neisseria gonorrhoeae*, *Xanthomonas oryzae*, and the *Thiobacillus* species somewhere in the  $\gamma$ - $\beta$  supergroup. In all these cases, the relationships have been suggested by other studies of SS-rRNA sequences (see Maidak et al. 1994; Olsen et al. 1994; Woese 1987). Finding the same patterns in the RecA trees serves to confirm the previous suggestions of the phylogenetic associations indicated between these species. Thus even though the RecA trees are based on analysis of highly conserved protein sequences, they do appear to have resolution for even close relatives as suggested by Lloyd and Sharp (1993).

Most of the differences between the RecA and SS-rRNA trees for the Proteobacteria are in areas of low resolution (differences among the trees generated by the different methods) or low bootstrap values for one or both of the molecules and thus are probably not biologically significant. For example, the differences in the grouping of the  $\delta$  and  $\epsilon$  clades within the Proteobacteria discussed above appear to be due to a lack of resolution of the SS-rRNA trees with the species represented here. In addition, the branching order between *Haemophilus influenzae*, the *Proteus* species, the *Vibrios*, and the enteric species is ambiguous in the SS-rRNA trees yet consistent in the RecA trees. In other cases, the SS-rRNA trees appear to have more resolution than the RecA trees. For example, the specific position of the RecA from *L. pneumophila* is ambiguous (Fig. 2a) yet the SS-rRNA of this species consistently groups with the  $\gamma$ 1 and  $\gamma$ 2 groups, and thus can be considered part of the  $\gamma$  clade (Fig. 2b). Analysis of other SS-rRNA sequences suggests that the position of the Legionellaceae in the  $\gamma$  subgroup is robust (Fry et al. 1991; Weisburg et al. 1989a). Similarly, the exact position of the *N. gonorrhoeae* RecA is ambiguous, yet its SS-rRNA groups consistently with the  $\beta$  species.

There are branching patterns within the proteobacteria that have high resolution and robustness for each molecule but are different between the two. The most striking example of this is the phylogenetic position of the sequences from *Acidiphilium facilis*. The *A. facilis* RecA branches with the *Thiobacillus ferrooxidans* RecA in the  $\beta$ - $\gamma$  supergroup in all trees (Fig. 2) and the node joining these two species has very high bootstrap values (Table 2). However, the corresponding *A. facilis* SS-rRNA consistently branches with species in the  $\alpha$  clade also with high bootstrap values. Thus either the SS-rRNA and RecA genes have different phylogenetic histories, or one of the trees is inaccurate. The grouping of *Acidiphilium* species within the  $\alpha$  subgroup appears to be a reliable representation of the SS-rRNA relationships (Sievers et al. 1994), so it is unlikely that the SS-rRNA trees here are biased by species sampling. It has been suggested that the *A. facilis* RecA sequence contains many sequencing errors and it is currently being resequenced

(Roca 1995, pers. comm.). Errors in the sequence would explain the unusual amino acids found in the *A. facilis* RecA in otherwise highly conserved regions (Fig. 1) and the extremely long branch length for this sequence in all phylogenetic methods (Fig. 3). Thus the position of the *A. facilis* RecA in the trees here may not represent the actual evolutionary history of this gene.

*M. xanthus*, the only  $\delta$  Proteobacteria represented in this analysis, is the only species known to encode two RecA proteins. There are at least two plausible explanations for this: lateral transfer from another species or gene duplication. The phylogenetic analysis of the two proteins helps limit the possibilities for when and how a duplication or lateral transfer could have occurred. In all the RecA trees, the two *M. xanthus* proteins branch together, showing that they are more related to each other than to any other known RecAs. However, the node joining them is quite deep, indicating that the degree of evolutionary separation between them is quite high. Thus if a duplication event was what led to these two genes in the same species, it apparently happened reasonably early in the history of the  $\delta$  clade. If one of these genes was obtained by a lateral transfer from another species, most likely the donor was another  $\delta$  species. It is interesting that the bootstrap values for the node joining the two *M. xanthus* RecAs are relatively low in all methods (Table 2). This indicates that the branching together is not very stable and is affected by the choice of alignment positions used in the phylogenetic analysis. Perhaps there was a gene conversion event after a lateral transfer or duplication and only certain regions of the *recA* genes underwent the conversion. Alternatively, the low bootstraps could also be explained if a duplication occurred right at or near the time of separation of the  $\delta$  clade from the other proteobacterial groups. The specific history of these two genes will probably be best resolved by studies of RecAs in other  $\delta$  species.

### Gram-Positives

Previous studies have shown that gram-positive species are divided into multiple phylogenetically distinct groups (Woese 1987). Whether these distinct groups are monophyletic has been the subject of a great deal of research and debate (e.g., Gupta et al. 1994; Van De Peer et al. 1994; Weisburg et al. 1989c; Woese 1987). For example, studies of HSP70 genes (Viale et al. 1994) and some studies of rRNA genes (Woese 1987) suggest the gram-positives are monophyletic while studies of EF-TUs (Ludwig et al. 1994), ATPase $\beta$  (Ludwig et al. 1994), and different studies of rRNA genes (Van De Peer et al. 1994) suggest they are not.

Species from two of the gram-positive groups, the low-GCs and the high-GCs, are represented in the analysis here (Table 1). In all the RecA and SS-rRNA trees inferred in this study, the sequences from the high-GC

species cluster together (Fig. 2). In addition these species have the same branching patterns within this group in all trees of both molecules. Thus the RecA data support the phylogenetic coherence of as well as the branching topology within the high-GC clade. In contrast, the RecA and SS-rRNA trees are not congruent for the relationships among sequences from low-GC gram-positive species. In all the SS-rRNA trees, the sequences from species considered to be low-GC gram-positives are monophyletic, as might be expected, since the classification of these species was based on SS-rRNA comparisons. However, in all the RecA trees the sequences from the low-GCs are not monophyletic (e.g., Fig. 3). This may be due to a combination of poor species sampling and unusual evolutionary patterns. In four of the five RecA trees only one RecA, that of the spirochaete *Borrelia burgdorferi*, prevents the low-GCs as a whole from being monophyletic (e.g., Fig. 3). The bootstrap values for the position of the *B. burgdorferi* RecA are relatively low in all of these trees, and since this is the only sample from the spirochaetes, its position may be unreliable. In addition, in three of the four SS-rRNA trees, the *B. burgdorferi* sequence is an outgroup to the low-GCs. Thus with the species sampled here the *B. burgdorferi* sequences tend to group with the sequences from low-GCs. Yet another factor that could contribute to a biased placement of the *B. burgdorferi* RecA is the apparent high rate of sequence change in the mycoplasmal RecAs, which can be seen by their long branch lengths (Fig. 3a). A rapid rate of mycoplasmal protein evolution has been thought to complicate trees of other proteins (e.g., Ludwig et al. 1994). The inclusion of additional sequences from the spirochetes and other low-GC gram-positives may help resolve whether this difference between the RecA and SS-rRNA trees is biologically significant.

With the species represented here, the branching between the high- and low-GCs is unresolved in both the RecA and SS-rRNA trees. Interestingly, in all the RecA trees, the high-GCs form a group with the cyanobacteria. Analysis of other genes has suggested that the cyanobacteria and gram-positives are sister groups (e.g., Van De Peer et al. 1994; Viale et al. 1994; Woese 1987). The RecAs are one of if not the only case in which cyanobacterial sequences consistently group with the high-GCs to the exclusion of the low-GCs. However, the bootstrap values for the node linking these two groups are moderate (31–40), indicating that this association is a good, but not great, representation of the relationships of RecA sequences.

### Cyanobacteria

The RecA and SS-rRNA trees both show the cyanobacteria forming a coherent clade. The nuclear encoded chloroplast RecA from *A. thaliana* groups consistently with the cyanobacterial RecAs. This suggests that the *A.*

*thaliana recA* gene is derived from the *recA* gene of a cyanobacterial-like ancestor to the *A. thaliana* chloroplast and that, as has been demonstrated for many other genes, it was transferred to the nucleus after endosymbiosis. Given the high degree of sequence conservation in RecAs, it is possible that studies of chloroplast evolution might be aided by sequencing of additional nuclear encoded chloroplast RecAs. In addition, all the RecAs from this group (including the *A. thaliana* RecA) contain an alignment gap not found in any other RecAs (see Fig. 1). This could serve as a sequence signature for cyanobacterial and chloroplast RecAs and further serves to demonstrate the relatedness among chloroplasts and cyanobacteria. As discussed above, the cyanobacterial RecAs group with those of the high-GC gram-positives in all trees.

#### *Deinococcus/Thermus* Group

The RecAs of *D. radiodurans* and the two *Thermus* species form a clade with high bootstrap values in all the trees (see Table 2, Fig. 2). Analysis of other data suggests that these species are part of a clade (Ludwig et al. 1994; Weisburg et al. 1989b). However, these sequences do not consistently form a clade in the SS-rRNA trees here (although they do in the *dnaps* tree (not shown)). Inclusion of additional SS-rRNA sequences allows for better resolution of this clade, probably because of GC content variation among the species (Embley et al. 1993). Thus with the species used here, the RecA sequences show resolution of the *Deinococcus-Thermus* group while the SS-rRNA sequences do not. This may be due to less of a GC bias in the RecA sequences, which was documented by Lloyd and Sharp (1993). The RecA analysis also supports previous assertions that this group is one of the deeper-branching bacterial phyla (Weisburg et al. 1989b) and shows that RecA has resolution even for deep branches. However this conclusion relies on the rooting of the RecA trees to the *A. pyrophilus* sequence, which has low support (see above).

#### Other Taxa

There is little resolution in the RecA trees regarding the position of the *Thermotoga maritima*, *Chlamydia trachomatis*, and *Bacteroides fragilis* proteins. These RecA proteins do not show consistent affiliations with any individual sequences or groups (Figs. 2, 3) and the bootstrap values for their positions in the individual trees are low (Fig. 3). I believe that this is due to these sequences being the only representatives from large phylogenetic groups (Thermotogales, Chlamydia, and Bacteroides, respectively). Using the same sets of sequences as in the RecA trees, the SS-rRNA trees show a similar lack of resolution for sequences that are individual representatives of large groups (in this case, *C. trachomatis*, *B. fragilis*, and *B. burgdorferi*). It would be useful to have

more RecA genes from these phylogenetic groups to better determine if the RecA- and SS-rRNA-based trees are congruent for these bacterial groups. It is interesting that although the specific position of the *T. maritima* RecA is ambiguous, it never branches below the *Deinococcus-Thermus* sequences as the *T. maritima* SS-rRNA does in all the SS-rRNA trees. Thus even if the rooting of the RecA tree with *A. pyrophilus* is incorrect, the *A. pyrophilus* and *T. maritima* RecAs never branch immediately near each other as they do in the SS-rRNA trees. Since the RecA trees appear to be less biased by GC content variation than the SS-rRNA trees (Lloyd and Sharp 1993) it seems plausible that the close branching of the *T. maritima* and *A. pyrophilus* SS-rRNAs may be caused by GC content convergence.

#### Conclusions

Comparison of phylogenetic results for particular taxa using different genes can help determine what genes are useful for evolutionary studies as well as whether different genes have different histories (as could be caused by lateral transfers). However, in order to make direct comparisons it is important to remove as many variables in the studies of the different genes. For example, many researchers studying bacterial systematics compare phylogenetic trees of particular genes to standard trees of SS-rRNA sequences. Yet when these trees have differences with the SS-rRNA trees it is not always clear whether the differences are due to use of different techniques (SS-rRNA trees frequently are constructed with maximum likelihood methods while such methods are still difficult to apply to large numbers of protein sequences), the inclusion of different sets of species (there are some 3000 SS-rRNA sequences that can be used), or true differences in branching or resolution power of different molecules. In the analysis presented here I have compared phylogenetic trees of RecA and SS-rRNA sequences using similar techniques with essentially the same sets of species. Overall, the branching patterns and powers of resolution of the two molecules are highly similar. The similar branching patterns lends support to the general pattern of bacterial systematics inferred from SS-rRNA sequences. This indicates either that the potential problems with SS-rRNA trees have little effect on phylogenetic analysis or that the RecA trees are biased in the same ways by these problems. In some cases, the RecA trees have resolution while the SS-rRNA trees do not (e.g., for the monophyly of the Proteobacteria and the grouping of *D. radiodurans* and the *Thermus* species) and in other cases the SS-rRNA trees have resolution (e.g., the position of *T. maritima*; the placement of *L. pneumophila* within the  $\gamma$ -Proteobacteria and the monophyly of the low-GC gram-positives). The lack of resolution of some of the deep branches in the RecA trees appears to be a problem of the species sampled—a sim-



ilar lack of resolution is seen in the SS-rRNA trees when using the same species-set. RecA comparisons may be particularly useful in cases of high GC bias. However the presence of two highly diverged *recA* genes in one species means that some caution should be used when using RecA studies for phylogenetic analysis, especially for the  $\delta$  proteobacteria. It remains to be seen whether some of the unusual patterns in the RecA trees (such as the grouping of the cyanobacteria with the high-GC gram-positives and the branching of *T. maritima* above the *Deinococci-Thermus* group) are supported by future studies.

In conclusion I would like to emphasize some of the features of RecA that make it a good choice for molecular systematic studies. Among protein encoding genes RecA is relatively easy to clone from new species—either by degenerate PCR (e.g., Duwat et al. 1992a,b; Dybvig et al. 1992; Dybvig and Woodard 1992; Quivey and Faustoferri 1992) or functional complementation of the radiation sensitivity of *recA* mutants from other species (Calero et al. 1994; De Mot et al. 1993; Favre et al. 1991; Gomelsky et al. 1990; Tatum et al. 1993). RecA protein function appears to be conserved in all bacteria and there are similar proteins in eukaryotes and archaea (Clark and Sandler 1994), although whether these can be used reliably for phylogenetics of all three kingdoms remains to be seen. As with SS-rRNAs, some regions of RecA are virtually completely conserved between species and other regions are variable even between close relatives. This allows for resolution of relationships among both close and distant relatives. The extensive size and sequence conservation among RecAs makes alignments virtually unambiguous, limiting complications due to incorrectly assigning homology. In addition, since RecA sequences can be compared at the protein and the DNA level it may be possible to limit problems due to nucleotide composition convergence between species. However, perhaps most importantly, I have shown here that phylogenetic trees of RecA sequences have similar topologies and similar resolution, even for deep branches, to trees of SS-rRNA sequences from the same species. This not only demonstrates that the genomes of these species are not completely mosaic but also that molecular systematics of bacteria is reliable, and that RecA comparisons are useful for such molecular systematic studies.

Finally, I would like to suggest two additional reasons why researchers might want to choose RecA for molecular systematic studies. First, the cloning and sequencing of *recA* genes from new species facilitate the creation of *recA* mutants which are useful to have for bacterial species. Also, with the availability of the crystal structure of the *E. coli* protein and with information about the phenotypes of hundreds of *recA* mutants, I believe RecA can become a model for studies of protein evolution.

**Acknowledgments.** I would like to thank P.C. Hanawalt for support and encouragement; M.B. Eisen for help with computer analysis; S.

Smith and J. Felsenstein for making their computer programs freely available; M. Huang, J. Coleman, A.J. Clark, S. Sandler, W. Finch, and S. Mongkolsuk for making sequences available prior to publication; M. Feldman for the use of computer equipment; and D. Pollock, M.B. Eisen, M.-I. Benito, H. Hamilton, D. Distel, and J.D. Palmer for helpful discussion. Finally, I want to thank A.I. Roca for help in all aspects of this project including and especially the identification of errors in *recA* sequences in databases and aid in getting access to unpublished sequences. During the course of this research I have received support from a predoctoral fellowship from the N.S.F., an N.I.H. grant to P.C. Hanawalt, and a tuition grant from the Woods Hole Institution to attend the Summer Workshop on Molecular Evolution, 1994.

## References

- Aigle B, Schneider D, Decaris B (1994) Genbank entry Z30324
- Akaboshi E, Yip ML, Howard-Flanders P (1989) Nucleotide sequence of the *recA* gene of *Proteus mirabilis*. Nucleic Acids Res 17:4390
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. J Mol Biol 215:403–410
- Angov E, Camerini-Otero RD (1994) The *recA* gene from the thermophile *Thermus aquaticus* YT-1: cloning, expression, and characterization. J Bacteriol 176:1405–1412
- Ball TK, Wasmuth CR, Braunagel SC, Benedik MJ (1990) Expression of *Serratia marcescens* extracellular proteins requires *recA*. J Bacteriol 172:342–349
- Bayles KW, Brunskill EW, Iandolo JJ, Hruska LL, Huang S, Pattee PA, Smiley BK, Yasbin RE (1994) A genetic and molecular characterization of the *recA* gene from *Staphylococcus aureus*. Gene 147: 13–20
- Berson AE, Peters MR, Waleh NS (1990) Nucleotide sequence of *recA* gene of *Aquaspirillum magnetotacticum*. Nucleic Acids Res 18:675
- Billman-Jacobe H (1994) Genbank entry X77384
- Binet M-N, Osman M, Jagendorf AT (1993) Genomic nucleotide sequence of a gene from *Arabidopsis thaliana* encoding a protein homolog of *Escherichia coli* RecA. Plant Physiol 103:673–674
- Boorstein WR, Ziegelhoffer T, Craig EA (1994) Molecular evolution of the HSP70 multigene family. J Mol Evol 38:1–17
- Bruns TD and Szaro TM (1992) Rate and mode differences between nuclear and mitochondrial small-subunit rRNA genes in mushrooms. Mol Biol Evol 9:836–855.
- Burggraf S, Olsen GJ, Stetter KO, Woese CR (1992) A phylogenetic analysis of *Aquifex pyrophilus*. Syst Appl Microbiol 15:352–356
- Calero S, Fernandez de Henestrosa AR, Barbe J (1994) Molecular cloning, sequence and regulation of expression of the *recA* gene of the phototrophic bacterium *Rhodobacter sphaeroides*. Mol Gen Genet 242:116–120
- Cerutti H, Osman M, Grandoni P, Jagendorf AT (1992) A homolog of *Escherichia coli* RecA protein in plastids of higher plants. Proc Natl Acad Sci USA 89:8068–8072
- Clark AJ, Margulies AD (1965) Isolation and characterization of recombination-deficient mutants of *Escherichia coli* K-12. Proc Natl Acad Sci USA 53:451–459
- Clark AJ, Sandler SJ (1994) Homologous genetic recombination: the pieces begin to fall into place. Crit Rev Microbiol 20:125–142
- Davis EO, Sedgwick SG, Colston MJ (1991) Novel structure of the *recA* locus of *Mycobacterium tuberculosis* implies processing of the gene product. J Bacteriol 173:5653–5662
- Davis EO, Thangaraj HS, Brooks PC, Colston MJ (1994) Evidence of selection for protein introns in the *recAs* of pathogenic mycobacteria. Embo J 13:699–703
- Delwiche CF, Kuschel M, Palmer JD (1995) Phylogenetic analysis of *tufA* sequences indicates a cyanobacterial origin of all plastids. Mol Phylogen Evol 4:110–128
- De Mot R, Laeremans T, Schoofs G, Vanderleyden J (1993) Charac-

- terization of the *recA* gene from *Pseudomonas fluorescens* OE 28.3 and construction of a *recA* mutant. *J Gen Microbiol* 139:49–57
- De Soete G (1983) A least squares algorithm for fitting additive trees to proximity data. *Psychometrika* 48:621–626
- Dunkin SM, Wood DO (1994) Isolation and characterization of the *Rickettsia prowazekii recA* gene. *J Bacteriol* 176:1777–1781
- Duwat P, Ehrlich SD, Gruss A (1992a) A general method for cloning *recA* genes of gram-positive bacteria by polymerase chain reaction. *J Bacteriol* 174:5171–5175
- Duwat P, Ehrlich SD, Gruss A (1992b) Use of degenerate primers for polymerase chain reaction cloning and sequencing of the *Lactococcus lactis* subsp. *lactis recA* gene. *Appl Environ Microbiol* 58:2674–2678
- Dybvig K, Hollingshead SK, Heath DG, Clewell DB, Sun F, Woodard A (1992) Degenerate oligonucleotide primers for enzymatic amplification of *recA* sequences from gram-positive bacteria and mycoplasmas. *J Bacteriol* 174:2729–2732
- Dybvig K, Woodard A (1992) Cloning and DNA sequence of a mycoplasmal *recA* gene. *J Bacteriol* 174:778–784
- Eisen JA, Smith SW, Cavanaugh CM (1992) Phylogenetic relationships of chemoautotrophic bacterial symbionts of *Solemya velum* Say (Mollusca: Bivalvia) determined by 16S rRNA sequence analysis. *J Bacteriol* 174:3416–3421
- Embley TM, Thomas RH, Williams RAD (1993) Reduced thermophilic bias in the 16S rDNA sequence from *Thermus ruber* provides further support for a relationship between *Thermus* and *Deinococcus*. *Syst Appl Microbiol* 16:25–29
- Favre D, Cryz SJ Jr, Viret JF (1991) Cloning of the *recA* gene of *Bordetella pertussis* and characterization of its product. *Biochimie* 73:235–44
- Favre D, Viret JF (1990) Nucleotide sequence of the *recA* gene of *Bordetella pertussis*. *Nucleic Acids Res* 18:4243
- Felsenstein J (1985) Confidence limits on phylogenies: an approach using the bootstrap. *Evolution* 39:783–791
- Felsenstein J (1989) PHYLIP—phylogeny inference package, version 3.2. *Cladistics* 5:164–166
- Fernandez de Henestrosa AR (1994) Genbank entry X82183
- Fitch WM, Margoliash E (1967) Construction of phylogenetic trees. *Science* 155:279–284
- Fox GE, Stackebrandt E, Hespell RB, Gibson J, Maniloff J, Dyer TA, Wolfe RS, Balch WE, Tanner RS, Magrum LJ, Zablen LB, Blake-more R, Gupta R, Bonen L, Lewis BJ, Stahl DA, Leuhrs KH, Chen KN, Woese CR (1980) The phylogeny of prokaryotes. *Science* 209:457–463
- Fry NK, Warwick S, Saunders NA, Embley TM (1991) The use of 16S ribosomal RNA analyses to investigate the phylogeny of the family Legionellaceae. *J Gen Microbiol* 137:1215–1222
- Fyfe JA, Davies JK (1990) Nucleotide sequence and expression in *Escherichia coli* of the *recA* gene of *Neisseria gonorrhoeae*. *Gene* 93:151–156
- Gammie AE, Crosa JH (1991) Co-operative autoregulation of a replication protein gene. *Mol Microbiol* 5:3015–3023
- Gatesy J, Desalle R, Wheeler W (1993) Alignment-ambiguous nucleotide sites and the exclusion of systematic data. *Mol Phylog Evol* 2:152–157
- Gomelsky M, Gak E, Chistoserdov A, Bolotin A, Tsygankov YD (1990) Cloning, sequence and expression in *Escherichia coli* of the *Methylobacillus flagellatum recA* gene. *Gene* 94:69–75
- Goodman HJ, Woods DR (1990) Molecular analysis of the *Bacteroides fragilis recA* gene. *Gene* 94:77–82
- Gregg-Jolly LA, Ornston LN (1994) Genbank entry L26100
- Guerry P, Pope PM, Burr DH, Leifer J, Joseph SW, Bourgeois AL (1994) Development and characterization of *recA* mutants of *Campylobacter jejuni* for inclusion in attenuated vaccines. *Infect Immunol* 62:426–432
- Gupta RS, Golding GB, Singh B (1994) Hsp70 phylogeny and the relationship between archaeobacteria, eubacteria, and eukaryotes. *J Mol Evol* 39:537–540
- Gutell RR, Larsen N, Woese CR (1994) Lessons from an evolving rRNA: 16S and 23S rRNA structures from a comparative perspective. *Microbiol Rev* 58:10–26
- Gutman PD, Carroll JD, Masters CI, Minton KW (1994) Sequencing, targeted mutagenesis and expression of a *recA* gene required for the extreme radioresistance of *Deinococcus radiodurans*. *Gene* 141:31–37
- Haas R (1994) Genbank entry Z35478
- Hasegawa M, Hashimoto T (1993) Ribosomal RNA tree misleading? *Nature* 361:23
- Henikoff S (1993) Sequence analysis by electronic mail server. *Trends Biochem Sci* 18:267–268
- Hillis DM (1995) Approaches for assessing phylogenetic accuracy. *Syst Biol* 44:3–16
- Horii T, Ogawa T, Ogawa H (1980) Organization of the *recA* gene of *Escherichia coli*. *Proc Natl Acad Sci USA* 77:313–317
- Inagaki K, Tomono J, Kishimoto N, Tano T, Tanaka H (1993) Cloning and sequence of the *recA* gene of *Acidiphilium facilis*. *Nucleic Acids Res* 21:4149
- Jinks-Robertson S, Nomura M (1987) Ribosomes and tRNA, In: Neidhardt FC (ed) *Escherichia coli* and *Salmonella typhimurium* cellular and molecular biology. American Society for Microbiology, Washington, DC, pp 1358–1385
- Kato R, Kuramitsu S (1993) RecA protein from an extremely thermophilic bacterium, *Thermus thermophilus* HB8. *J Biochem* 114:926–929
- Kerins SM, Fitzpatrick R, O'Donohue M, Dunican L (1994) Genbank entry X75085
- Kimura M (1980) A simple method for estimating evolutionary rate of base substitutions through comparative studies of nucleotide sequences. *J Mol Evol* 16:111–120
- King KW, Woodard A, Dybvig K (1994) Cloning and characterization of the *recA* genes from *Mycoplasma pulmonis* and *M. mycoides* subsp. *mycoides*. *Gene* 139:111–115
- Klenk H-P, Zillig W (1994) DNA-dependent RNA polymerase subunit b as a tool for phylogenetic reconstructions: branching topology of the archaeal domain. *J Mol Evol* 38:420–432
- Kowalczykowski SC, Dixon DA, Eggleston AK, Lauder SS, Rehauer WM (1994) Biochemistry of homologous recombination in *Escherichia coli*. *Microbiol Rev* 58:401–465
- Kowalczykowski SC (1991) Biochemical and biological function of *Escherichia coli* RecA protein: behavior of mutant RecA proteins. *Biochimie* 73:289–304
- Kryukov VM, Suchkov IY, Sazykin IS, Mishankin BN (1993) Genbank entry X75336
- Lake JA (1991) The order of sequence alignment can bias the selection of tree topology. *Mol Biol Evol* 8:378–385
- Lane DJ, Harrison AP, Stahl DA, Pace B, Giovannoni SJ, Olsen GJ, Pace NR (1992) Evolutionary relationships among sulfur- and iron-oxidizing eubacteria. *J Bacteriol* 174:269–278.
- Lane DJ, Pace B, Olsen GJ, Stahl DA, Sogin ML, Pace NR (1985) Rapid determination of 16S rRNA sequences for phylogenetic analysis. *Proc Natl Acad Sci USA* 82:6955–6959
- Larsen SH (1994) Genbank entry U16739
- Lecointre G, Philippe H, Van Le HL, Le Guyader H (1993) Species sampling has a major impact on phylogenetic inference. *Mol Phylog Evol* 2:205–224
- Lloyd AT, Sharp PM (1993) Evolution of the *recA* gene and the molecular phylogeny of bacteria. *J Mol Evol* 37:399–407
- Ludwig W, Kirchhof G, Klugbauer N, Weizenegger M, Betzl D, Ehrmann M, Hertel C, Jilg S, Tatzel R, Zitzelsberger H, Liebl S, Hochberger M, Shah J, Lane D, Wallnöfer PR, Schleifer KH (1992) Complete 23S ribosomal RNA sequences of gram-positive bacteria with a low DNA G plus C content. *Syst Appl Microbiol* 15:487–501
- Ludwig W, Neumaier J, Klugbauer N, Brockmann E, Roller C, Jilg S, Reetz K, Schachtner I, Ludvigsen A, Bachleitner M, Fischer U, Schleifer KH (1994) Phylogenetic relationships of bacteria based

- on comparative sequence analysis of elongation factor TU and ATP-synthase beta-subunit genes. *Antonie Van Leeuwenhoek* 64: 285–305
- Luo J, Burns G, Sokatch JR (1993) Construction of chromosomal *recA* mutants of *Pseudomonas putida* PpG2. *Gene* 136:263–266
- Maddison WP, Maddison DR (1992) MacClade version 3. Sinauer Associates, Sunderland, MA
- Maidak BL, Larsen N, McCaughey MJ, Overbeek R, Olsen GJ, Fogel K, Blandy J, Woese CR (1994) The ribosomal database project. *Nucleic Acids Res* 22:3485–3487
- Margraf RL, Roca AI, Cox MM (1995) The deduced *Vibrio cholerae* RecA amino acid sequence. *Gene* 152:135–136
- Martin B, Ruellan JM, Angulo JF, Devoret R, Claverys JP (1992) Identification of the *recA* gene of *Streptococcus pneumoniae*. *Nucleic Acids Res* 20:6412
- Medlin L, Elwood HJ, Stickel S, Sogin ML (1988) The characterization of enzymatically amplified eukaryotic 16S-like ribosomal RNA-coding regions. *Gene* 71:491–500
- Michiels J, Vande Broek A, Vanderleyden J (1991) Molecular cloning and nucleotide sequence of the *Rhizobium phaseoli* *recA* gene. *Mol Gen Genet* 228:486–490
- Miyamoto MM, Fitch WM (1995) Testing species phylogenies and phylogenetic methods with congruence. *Syst Biol* 44:64–76
- Murphy RC, Bryant DA, Porter RD, de Marsac NT (1987) Molecular cloning and characterization of the *recA* gene from the cyanobacterium *Synechococcus* sp. strain PCC 7002. *J Bacteriol* 169:2739–2747
- Murphy RC, Gasparich GE, Bryant DA, Porter RD (1990) Nucleotide sequence and further characterization of the *Synechococcus* sp. strain PCC 7002 *recA* gene: complementation of a cyanobacterial *recA* mutation by the *Escherichia coli* *recA* gene. *J Bacteriol* 172: 967–976
- Nakazawa T, Kimoto M, Abe M (1990) Cloning, sequencing, and transcriptional analysis of the *recA* gene of *Pseudomonas cepacia*. *Gene* 94:83–88
- Nomura M, Morgan EA, Jaskunas S (1977) Genetics of bacterial ribosomes. *Annu Rev Genet* 11:297–347
- Nussbaumer B, Wohleben W (1994) Identification, isolation and sequencing of the *recA* gene of *Streptomyces lividans* TK24. *FEMS Microbiol Lett* 118:57–63
- Ogawa T, Yu X, Shinohara A, Egelman EH (1993) Similarity of the yeast RAD51 filament to the bacterial RecA filament. *Science* 259: 1896–1899
- Olsen GJ (1988) Phylogenetic analysis using ribosomal RNA. *Methods Enzymol* 164:793–812
- Olsen GJ, Lane DJ, Giovannoni SJ, Pace NR, Stahl DA (1986) Microbial ecology and evolution: a rRNA approach. *Annu Rev Microbiol* 40:337–365
- Olsen GJ, Woese CR, Overbeek R (1994) The winds of (evolutionary) change: breathing new life into microbiology. *J Bacteriol* 176:1–6
- Owtrim GW, Coleman JR (1989) Regulation of expression and nucleotide sequence of the *Anabaena variabilis* *recA* gene. *J Bacteriol* 171:5713–5719
- Pace NR, Olsen GJ, Woese CR (1986) Ribosomal RNA phylogeny and the primary lines of evolutionary descent. *Cell* 45:325–326
- Pitulle C, Yang Y, Marchiani M, Moore ERB, Siefert JL, Aragno M, Jurtshuk PJ, Fox GE (1994) Phylogenetic position of the genus *Hydrogenobacter*. *Int J Syst Bacteriol* 44:620–626
- Quivey RG Jr, Faustoferri RC (1992) In vivo inactivation of the *Streptococcus mutans* *recA* gene mediated by PCR amplification and cloning of a *recA* DNA fragment. *Gene* 116:35–42
- Rainey FA, Toalster R, Stackebrandt E (1993) *Desulfurella acetivorans* a thermophilic, acetate-oxidizing, and sulfur-reducing organism, represents a distinct lineage within the Proteobacteria. *Syst Appl Microbiol* 16:373–379.
- Ramesar RS, Abratt V, Woods DR, Rawlings DE (1989) Nucleotide sequence and expression of a cloned *Thiobacillus ferrooxidans* *recA* gene in *Escherichia coli*. *Gene* 78:1–8
- Rappold CSJ, Klingmueller W (1993) Genbank entry P33037
- Rensing SA, Maier UG (1994) Phylogenetic analysis of the stress-70 protein family. *J Mol Evol* 39:80–86
- Ridder R, Marquardt R, Esser K (1991) Molecular cloning and characterization of the *recA* gene of *Methylobacterium clara* and construction of *recA* deficient mutant. *Appl Microbiol Biotechnol* 35:23–31
- Roca AI, Cox MM (1990) The RecA protein: structure and function. *Crit Rev Biochem Mol Biol* 25:415–456
- Rothschild LJ, Ragan MA, Coleman AW, Heywood P, Gerbi SA (1986) Are rRNA sequences the Rosetta stone of phylogenetics. *Cell* 47:640
- Saitou N, Nei M (1987) The neighbor joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol* 4:406–425
- Sancar A, Stachek C, Konigsberg W, Rupp WD (1980) Sequences of the *recA* gene and protein. *Proc Natl Acad Sci USA* 77:2611–2615
- Sano Y, Kageyama M (1987) The sequence and function of the *recA* gene and its protein in *Pseudomonas aeruginosa* PAO. *Mol Gen Genet* 208:412–419
- Schoeniger M, Von Haeseler A (1994) A stochastic model for the evolution of autocorrelated DNA sequences. *Mol Phylog Evol* 3: 240–247
- Selbitschka W, Arnold W, Priefer UB, Rottschäfer T, Schmidt M, Simon R, Puhler A (1991) Characterization of *recA* genes and *recA* mutants of *Rhizobium meliloti* and *Rhizobium leguminosarum* biovar viciae. *Mol Gen Genet* 229:86–95
- Sievers M, Ludwig W, Teuber M (1994) Phylogenetic positioning of *Acetobacter*, *Gluconobacter*, *Rhodopila* and *Acidiphilium* species as a branch of acidophilic bacteria in the alpha-subclass of proteobacteria based on 16S ribosomal DNA sequences. *Syst Appl Microbiol* 17:189–196
- Smith SW (1991) Genetic data environment, version 2.2a. Harvard Genome Laboratory, Cambridge, MA
- Smith SW, Overbeek R, Woese CR, Gilbert W, Gillevet PM (1994) The genetic data environment: an expandable GUI for multiple sequence analysis. *CABIOS* 10:671–675.
- Sogin ML (1989) Evolution of eukaryotic microorganisms and their small subunit ribosomal RNA. *Am Zool* 29:487–500
- Stackebrandt E, Murray RGE, Trüper HG (1988) *Proteobacteria* classis nov., a name for the phylogenetic taxon that includes 'purple bacteria and their relatives.' *Int J Syst Bacteriol* 38:321–325
- Story RM, Bishop DK, Kleckner N, Steitz TA (1993) Structural relationship of bacterial RecA proteins to recombination proteins from bacteriophage T4 and yeast. *Science* 259:1892–1896
- Story RM, Steitz TA (1992) Structure of the RecA protein-ADP complex. *Nature* 355:374–376
- Story RM, Weber IT, Steitz TA (1992) The structure of the *E. coli* RecA protein monomer and polymer. *Nature* 355:318–325
- Stranathan MC, Bayles KW, Yasbin RE (1990) The nucleotide sequence of the *recE+* gene of *Bacillus subtilis*. *Nucleic Acids Res* 18:4249
- Stroehrer UH, Lech AJ, Manning PA (1994) Gene sequence of *recA+* and construction of *recA* mutants of *Vibrio cholerae*. *Mol Gen Genet* 244:295–302
- Swofford D (1991) Phylogenetic analysis using parsimony (PAUP), version 3.0d. Illinois Natural History Survey, Champaign, IL
- Tatum FM, Morfitt DC, Halling SM (1993) Construction of a *Brucella abortus* *recA* mutant and its survival in mice. *Microbiol Pathogen* 14:177–185
- Tayama K, Fukaya M, Takemura H, Okumura H, Kawamura Y, Horinouchi S, Beppu T (1993) Cloning and sequencing the *recA+* genes of *Acetobacter polyoxogenes* and *Acetobacter aceti*: construction of *recA-* mutants of by transformation-mediated gene replacement. *Gene* 127:47–52
- Thompson JD, Higgins DG, Gibson TJ (1994) Clustal w: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22:4673–4680

- Tolmasky ME, Gammie AE, Crosa JH (1992) Characterization of the *recA* gene of *Vibrio anguillarum*. *Gene* 110:41–48
- Van De Peer Y, Neefs JM, De Rijk P, De Vos P, De Wachter R (1994) About the order of divergence of the major bacterial taxa during evolution. *Syst Appl Microbiol* 17:32–38
- Vawter L, Brown WM (1993) Rates and patterns of base change in the small subunit ribosomal RNA gene. *Genetics* 134:597–608
- Venkatesh TV, Das HK (1992) The *Azotobacter vinelandii recA* gene: sequence analysis and regulation of expression. *Gene* 113:47–53
- Viale AM, Arakaki AK, Soncini FC, Ferreyra RG (1994) Evolutionary relationships among eubacterial groups as inferred from GroEL (chaperonin) sequence comparisons. *Int J Syst Bacteriol* 44:527–533
- Wardhan H, McPherson MJ, Harris CA, Sharma E, Sastry GR (1992) Molecular analysis of the *recA* gene of *Agrobacterium tumefaciens* C58. *Gene* 121:133–136
- Weisburg WG, Barns SM, Pelletier DA, Lane DJ (1991) 16S ribosomal DNA amplification for phylogenetic study. *J Bacteriol* 173:697–703
- Weisburg WG, Dobson ME, Samuel JE, Dasch GA, Mallavia LP, Baca O, Mendelco L, Sechrest JE, Weiss E, Woese CR (1989a) Phylogenetic diversity of the Rickettsiae. *J Bacteriol* 171:4202–4206
- Weisburg WG, Giovannoni SG, Woese CR (1989b) The *Deinococcus-Thermus* phylum and the effect of rRNA composition on phylogenetic tree construction. *Syst Appl Microbiol* 11:128–134
- Weisburg WG, Tully JG, Rose DL, Petzel JP, Oyaizu H, Yang D, Mandelco L, Sechrest J, Lawrence TG (1989c) A phylogenetic analysis of the mycoplasmas: basis for their classification. *J Bacteriol* 171:6455–6467
- Wetmur JG, Wong DM, Ortiz B, Tong J, Reichert F, Gelfand DH (1994) Cloning, sequencing, and expression of RecA proteins from three distantly related thermophilic eubacteria. *J Biol Chem* 269:25928–25935
- Woese CR (1991) The use of ribosomal RNA in reconstructing relationships among bacteria. In: Selander RK, Clark AG, Whittam TS (eds) *Evolution at the molecular level*. Sinauer Associates, Sunderland, MA, pp 1–24
- Woese CR (1987) Bacterial evolution. *Microbiol Rev* 51:221–271
- Woese CR, Achenbach L, Rouviere P, Mandelco L (1991) Archaeal phylogeny: reexamination of the phylogenetic position of *Archaeoglobus fulgidus* in the light of certain composition-induced artifacts. *Syst Appl Microbiol* 14:364–371
- Yao W, Vining LC (1994) Cloning and sequence analysis of a *recA*-like gene from *Streptomyces venezuelae* ISP5230. *FEMS Microbiol Lett* 118:51–56
- Zhang D, Fan H, McClarty G, Brunham RC (1994) Genbank entry U15281
- Zhao X, Dreyfus LA (1990) Expression and nucleotide sequence analysis of the *Legionella pneumophila recA* gene. *FEMS Microbiol Lett* 70:227–232
- Zhao XJ, McEntee K (1990) DNA sequence analysis of the *recA* genes from *Proteus vulgaris*, *Erwinia carotovora*, *Shigella flexneri* and *Escherichia coli* B/r. *Mol Gen Genet* 222:369–376
- Zulty JJ, Barcak GJ (1993) Structural organization, nucleotide sequence, and regulation of the *Haemophilus influenzae rec-1+* gene. *J Bacteriol* 175:7269–7281