

The genome sequence of *Rickettsia prowazekii* and the origin of mitochondria

Siv G. E. Andersson*, Alireza Zomorodipour*, Jan O. Andersson*, Thomas Sicheritz-Pontén*, U. Cecilia M. Alsmark*, Raf M. Podowski*, A. Kristina Näslund*, Ann-Sofie Eriksson*, Herbert H. Winkler† & Charles G. Kurland*

* Department of Molecular Biology, University of Uppsala, Uppsala S-75124, Sweden

† Department of Microbiology and Immunology, University of South Alabama, Mobile, Alabama 36688, USA

We describe here the complete genome sequence (1,111,523 base pairs) of the obligate intracellular parasite *Rickettsia prowazekii*, the causative agent of epidemic typhus. This genome contains 834 protein-coding genes. The functional profiles of these genes show similarities to those of mitochondrial genes: no genes required for anaerobic glycolysis are found in either *R. prowazekii* or mitochondrial genomes, but a complete set of genes encoding components of the tricarboxylic acid cycle and the respiratory-chain complex is found in *R. prowazekii*. In effect, ATP production in *Rickettsia* is the same as that in mitochondria. Many genes involved in the biosynthesis and regulation of biosynthesis of amino acids and nucleosides in free-living bacteria are absent from *R. prowazekii* and mitochondria. Such genes seem to have been replaced by homologues in the nuclear (host) genome. The *R. prowazekii* genome contains the highest proportion of non-coding DNA (24%) detected so far in a microbial genome. Such non-coding sequences may be degraded remnants of 'neutralized' genes that await elimination from the genome. Phylogenetic analyses indicate that *R. prowazekii* is more closely related to mitochondria than is any other microbe studied so far.

The *Rickettsia* are α -proteobacteria that multiply in eukaryotic cells only. *R. prowazekii* is the agent of epidemic, louse-borne typhus in humans. Three features of this endocellular parasite deserve our attention. First, *R. prowazekii* is estimated to have infected 20–30 million humans in the wake of the First World War and killed another few million following the Second World War (ref. 1). Because it is the descendent of free-living organisms^{2–4}, its genome provides insight into adaptations to the obligate intracellular lifestyle, with probable practical value. Second, phylogenetic analyses based on sequences of ribosomal RNA and heat-shock proteins indicate that mitochondria may be derived from the α -proteobacteria^{5,6}. Indeed, the closest extant relatives of the ancestor to mitochondria seem to be the *Rickettsia*^{7–10}. That modern *Rickettsia* favour an intracellular lifestyle identifies these bacteria as the sort of organism that might have initiated the endosymbiotic scenario leading to modern mitochondria¹¹. Finally, the genome of *R. prowazekii* is a small one, containing only 1,111,523 base pairs (bp). Its phylogenetic placement and many other characteristics identify it as a descendant of bacteria with substantially larger genomes^{2–4}. Thus *Rickettsia*, like mitochondria, are good examples of highly derived genomes, the products of several types of reductive evolution.

The genome sequence of *R. prowazekii* indicates that these three features may be related. For example, prokaryotic genomes evolving within a cell dominated by a much larger, eukaryote genome and constrained by bottle-necked population dynamics will tend to lose genetic information^{12,13}. Predictable sets of expendable genes will tend to disappear from the prokaryotic genome when they are made redundant by the activities of nuclear genes. Likewise, non-essential sequences and otherwise highly conserved gene clusters may be obliterated by deleterious mutations that are fixed in clonal parasite or organelle populations because they cannot be eliminated by selection. This process is ongoing in the *Rickettsia* genomes, as shown by the identification of sequences that have recently become pseudogenes. Also, a large fraction (~25%) of non-coding sequences in this genome may be gene remnants that have been

degraded by mutation and have not yet been removed from the genome. Finally, transfer of genes from a mitochondrial ancestor to the nucleus of the host would both reduce the mitochondrial genome size and stabilize the symbiotic relationship. Phylogenetic reconstructions that identify genes in the *Rickettsia* genome as sister clades to eukaryotic homologues found in the nucleus or the organelle support this interpretation. *Rickettsia* and mitochondria probably share an α -proteobacterial ancestor and a similar evolutionary history.

General features of the genome

The circular chromosome of *R. prowazekii* strain Madrid E has 1,111,523 bp and an average G+C content of 29.1% (Figs 1, 2). The genome contains 834 complete open reading frames with an average length of 1,005 bp. Protein-coding genes represent 75.4% of the genome and 0.6% of the genome encodes stable RNA. We have assigned biological roles to 62.7% of the identified genes and pseudogenes; 12.5% of the identified genes match hypothetical coding sequences of unknown function and the remaining 24.8% represent unusual genes with no similarities to genes in other organisms (Table 1). Multivariate statistical analysis has shown that there is no major variation in codon-usage patterns among genes that are expressed in different amounts, indicating that codon-usage patterns in *R. prowazekii* may be dominated mainly by mutational forces¹⁴. G+C-content values at the three codon positions average 40.4, 31.2 and 18.6%, and these values are similar at different positions in the genome. We classified the open reading frames with significant sequence-similarity scores to gene sequences in the public databases into functional categories (Table 1) that allow comparisons with the metabolic profiles of other bacterial genomes^{15–23}.

Non-coding DNA. The coding content of previously sequenced bacterial genomes is, on average, 91%, ranging from 87% in *Haemophilus influenzae* to 94% in *Aquifex aeolicum*. In comparison, a large fraction of the *R. prowazekii* genome, 24%, represents non-coding DNA (Fig. 3). A small fraction of this corresponds to

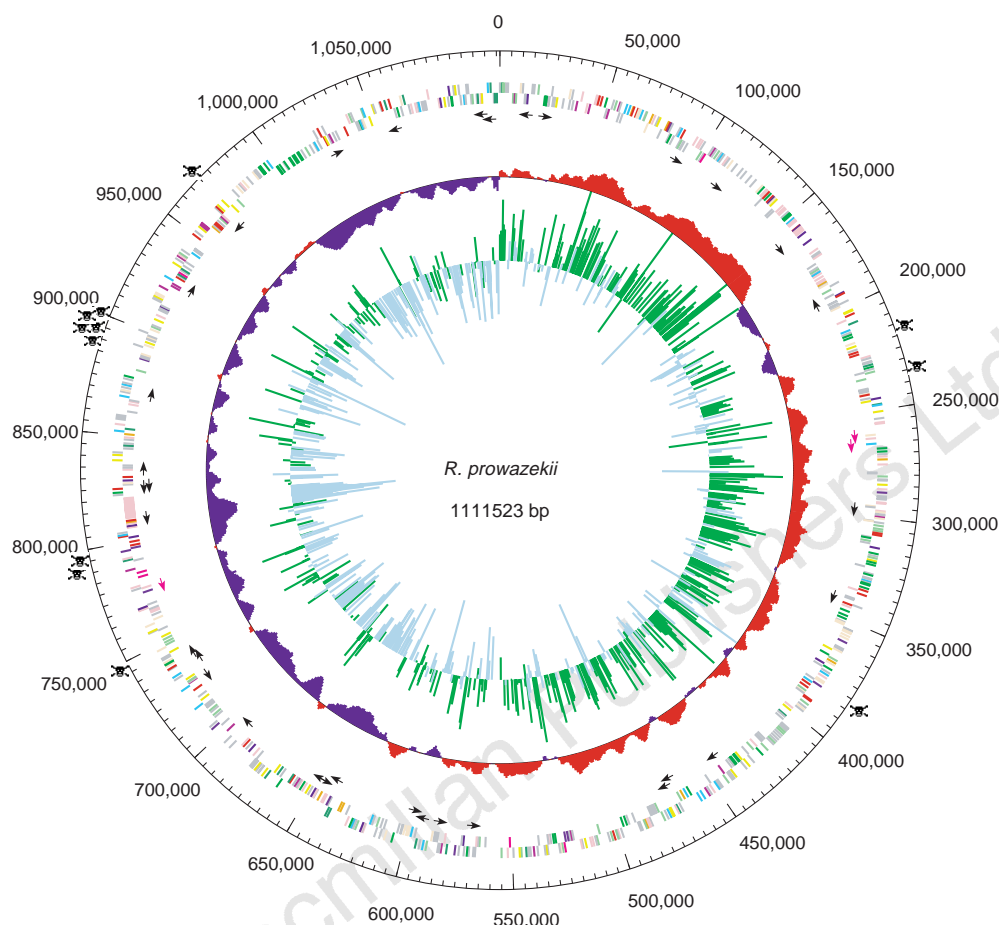


Figure 1 Overall structure of the *R. prowazekii* genome. The putative origin of replication is at 0 kb. The outer scale indicates the coordinates (in base pairs). The positions of pseudogenes are highlighted with death's heads. The distribution of genes is shown on the first two rings within the scale. The location and direction of transcription of rRNA are shown by pink arrows and of tRNA genes by black arrows. The next circle in shows GC-skew values measured over all bases in the genome. Red and purple colours denote positive and negative signs, respectively.

pseudogenes (0.9% of the genome) and less than 0.2% of the genome is accounted for by non-coding repeats. The remaining 22.9% contains no open reading frames of significant length and it has the low G+C content (mean 23.7%) that is characteristic of spacer sequences in the *R. prowazekii* genome¹⁴. A region of 30 kilobases (kb) located at position 886–916 kb contains as much as 41.6% non-coding DNA and 11.5% pseudogenes. The non-coding DNA in this region has a small, but significantly higher, G+C content (mean 27.3%) than non-coding DNA in other areas of the genome (mean 23.7%) ($P < 0.001$), indicating that it may correspond to inactivated genes that are being degraded by mutation (Fig. 3).

Origin of replication. The origin of replication has not been experimentally identified in the *R. prowazekii* genome, but we identified *dnaA* at ~750 kb. However, the genes flanking the *dnaA* gene differ from the conserved motifs found in *Escherichia coli* and *Bacillus subtilis* (*rnpA*–*rpmH*–*dnaA*–*dnaN*–*recF*–*gyrB*). In *R. prowazekii*, the genes *rnpA* and *rpmH* are located in the vicinity of *dnaA*, but in the reverse orientation compared to the consensus motif, and *dnaN*, *recF* and *gyrB* are located elsewhere.

The origin and end replication in microbial genomes are often associated with transitions in GC skew ($G - C/G + C$) values²⁴. In *R. prowazekii* we observe transitions in the GC skew values at

around 0 and 500–600 kb (Fig. 1). There is a weak asymmetry in the distribution of genes in the two strands, such that the first half of the genome has a 1.6-fold higher gene density on one strand and the second half of the genome has a 1.6-fold higher gene density on the other strand. The shift in coding-strand bias correlates with the shift in GC-skew values. As most genes are transcribed in the direction of replication in microbial genomes, the origin of replication may correspond to the shift in GC-skew values at the position that we have chosen as the start point for numbering. Indeed, several short sequence stretches that are characteristic of *dnaA*-binding motifs are found in the intergenic region of genes *RP001* and *RP885* at 0 kb, supporting this interpretation.

Stable RNA sequences and repeat elements. We identified 33 genes encoding transfer RNA, corresponding to 32 different isoacceptor-tRNA species. There is a single copy of each of the rRNA genes, with *rrs* located more than 500 kb away from the *rrl*–*rrf* gene cluster

around 0 and 500–600 kb (Fig. 1). There is a weak asymmetry in the distribution of genes in the two strands, such that the first half of the genome has a 1.6-fold higher gene density on one strand and the second half of the genome has a 1.6-fold higher gene density on the other strand. The shift in coding-strand bias correlates with the shift in GC-skew values. As most genes are transcribed in the direction of replication in microbial genomes, the origin of replication may correspond to the shift in GC-skew values at the position that we have chosen as the start point for numbering. Indeed, several short sequence stretches that are characteristic of *dnaA*-binding motifs are found in the intergenic region of genes *RP001* and *RP885* at 0 kb, supporting this interpretation.

Stable RNA sequences and repeat elements. We identified 33 genes encoding transfer RNA, corresponding to 32 different isoacceptor-tRNA species. There is a single copy of each of the rRNA genes, with *rrs* located more than 500 kb away from the *rrl*–*rrf* gene cluster

Figure 2 Linear map of the *R. prowazekii* chromosome. The position and orientation of known genes are indicated by arrows. Coding regions are colour-coded according to their functional roles. The positions of tRNA genes are indicated (inverted triangle on stalk). For additional information, see <http://evolution.bmc.uu.se/~siv/gnomics/Rickettsia.html>.

(Fig. 1). Comparison of the sequences from ten different *Rickettsia* species indicates that the disruption of the rRNA gene operon preceded the divergence of the typhus group and spotted fever group *Rickettsia* (S.G.E.A. *et al.*, unpublished observations). In addition, the genome contains a short sequence with similarity to a 213-nucleotide RNA molecule in *Bradyrhizobium japonicum* that may regulate transcription²⁵.

There are unusually few repeat sequences in this genome. We identified four different types of repeat sequence: all of these are located in intergenic regions. There is a sequence of 80 bp that is repeated seven times downstream of *rpmH* and *rnpA* in the *dnaA* region. A repetitive sequence of 325 bp is found at two intergenic regions that are more than 80 kb apart, downstream of the genes *ksgA* and *rnh*, respectively. A 440-bp-long repetitive sequence has been identified at two intergenic sites, 140 kb apart; one of these sites is downstream of *rff* and the others downstream of *pdhA* and

pdhB. Finally, two similar sequences of 730 bp are located immediately next to each other at 850 kb.

Paralogous families. We have identified 54 paralogous gene families comprising 147 gene products. Of these, 125 have an assigned function. Most paralogues encode proteins with transport functions, such as the ABC transporters, the proline/betaine transporters and the ATP/ADP transporters. Five paralogous genes located next to each other at 115 kb encode putative integral membrane proteins with unknown functions.

Biosynthetic pathways

A striking feature of the *R. prowazekii* genome is the small proportion of biosynthetic genes compared with free-living proteobacterial relatives (such as *Haemophilus influenzae*, *Helicobacter pylori* and *E. coli*)^{15,19,20}. This scarcity of biosynthetic functions is also seen in diverse endocellular and epicellular parasites^{16–18,23}. This scarcity of biosynthetic functions is also seen in diverse endocellular and epicellular parasites^{16–18,23}.

Amino-acid metabolism. As many as 43 and 69 genes required for amino-acid biosynthesis are found in *Helicobacter pylori* and *Haemophilus influenzae*, respectively. In contrast, *Mycoplasma genitalium* and *Borrelia burgdorferi* contain only *glyA*, which encodes serine hydroxymethyltransferase. This gene is also found in *R. prowazekii* (Table 1). Serine hydroxymethyltransferase catalyses the conversion of serine and tetrahydrofolate into glycine and methylenetetrahydrofolate, respectively. A role in tetrahydrofolate metabolism may account for the ubiquity of *glyA* in bacteria.

Seven genes normally associated with lysine biosynthesis (*lysC*, *asd*, *dapA*, *dapB*, *dapD*, *dapE* and *dapF*) are also present in *R. prowazekii*. The biosynthetic pathways leading to lysine, methionine and threonine share the first two of these (*lysC* and *asd*). However, none of the downstream genes for threonine biosynthesis are found in *R. prowazekii*. Likewise, the lysine pathway is incomplete, and *lysA*, which encodes the enzyme that converts meso-diaminopimelate to lysine, is missing. The likely role of the upstream genes of this pathway in *R. prowazekii* is the biosynthesis of diaminopimelate, an essential envelope component. We have therefore classified these genes as 'cell-envelope' genes (Table 1).

We have identified other genes that are superficially involved in the metabolism of amino acids, but which apparently function in deamination pathways that divert amino acids into the tricarboxylic acid (TCA) cycle. For example, there is *aatA*, encoding aspartate aminotransferase, which catalyses the degradation of aspartate to oxaloacetate and glutamate. *tdcB* encodes threonine deaminase, which converts threonine into α -ketobutyrate. Another gene (*ilvE*) encodes branched-chain-amino-acid aminotransferase, which converts leucine, isoleucine or valine into glutamate. *pccA* and *pccB* encode propionyl-CoA carboxylase, which converts propionyl-CoA, an intermediate in the breakdown of methionine, valine and isoleucine, into succinyl-CoA. The *pccA* and *pccB* gene products show greatest similarity to the eukaryotic proteins that are located in the mitochondrial matrix.

Nucleotide biosynthesis. No genes required for the *de novo* syntheses of nucleosides have been found in the *R. prowazekii* genome. However, four genes required for the conversion of nucleoside monophosphates into nucleoside diphosphates (*adk*, *gmk*, *cmk* and *pyrH*) are present. There are also two genes encoding ribonucleotide reductase, which converts ribonucleoside diphosphates into deoxyribonucleoside diphosphates. Nucleoside diphosphate kinase (encoded by *ndk*), which converts NDPs and dNDPs to NTPs and dNTPs, is also present in *R. prowazekii*. Finally, there is a complete set of genes for the conversion of dCTP and dUTP into TTP, including *thyA*, which codes for thymidylate synthase. Thus, the *R. prowazekii* genome encodes all of the enzymes required for the interconversion of nucleoside monophosphates into all of the other required nucleotides. The nucleoside monophosphates are probably imported from the eukaryotic host.

Table 1 Asterisks indicate putative pseudogenes. Abbreviations of species names

are: Bacteria: *Acinetobacter calcoaceticus* (B-Aca), *Actinobacillus actinomycetem-comitans* (B-Aac), *Acyrtosiphon condii* (B-Aco), *Agrobacterium tumefaciens* (B-Atu), *Alcaligenes eutrophus* (B-Aeu), *Anabena sp.* PCC7120 (B-Asp), *Anabena variabilis* (B-Ava), *Anacystis nidulans* (B-Ani), *Azorhizobium caulinodans* (B-Aca), *Azospirillum brasilense* (B-Abr), *Azotobacter vinelandii* (B-Avi), *Bacillus caldotenax* (B-Bca), *Bacillus stercorophilus* (B-Bst), *Bacillus subtilis* (B-Bsu), *Bartonella bacilliformis* (B-Bba), *Bartonella henselae* (B-Bhe), *Bordetella pertussis* (B-Bpe), *Borrelia burgdorferi* (B-Bbu), *Bradyrhizobium japonicum* (B-Bja), *Brucella abortus* (B-Bab), *Brucella ovis* (B-Bov), *Caulobacter crescentus* (B-Ccr), *Chlamydia trachomatis* (B-Ctr), *Chloroflexus aurantiacus* (B-Cau), *Chromatium visium* (B-Cvt), citrus-greening-disease-associated bacterium (B-Cgr), *Clostridium acetobutylicum* (B-Cac), *Clostridium pasteurianum* (B-Cpa), *Clostridium thermosaccharolyticum* (B-Cts), *Coxiella burnetii* (B-Cbu), *Erwinia chrysanthemi* (B-Ech), *Escherichia coli* (B-Eco), *Haemophilus influenzae* (B-Hin), *Helicobacter pylori* (B-Hpy), *Klebsiella pneumoniae* (B-Kpn), *Legionella pneumophila* (B-Lpn), *Leucothrix mucor* (B-Lmu), *Liberobacter africanum* (B-Laf), *Methylobacterium extorquens* (B-Mex), *Micrococcus luteus* (Mlu), *Moraxella catarrhalis* (Mca), *Mycobacterium leprae* (Mle), *Mycobacterium smegmatis* (B-Msm), *Mycobacterium tuberculosis* (B-Mtu), *Mycoplasma capricolum* (B-Mca), *Mycoplasma genitalium* (B-Mge), *Mycoplasma pneumoniae* (B-Mpn), *Paracoccus denitrificans* (B-Pde), *Pasteurella haemolytica* (B-Pha), *Plectonema boryanum* (B-Pbo), *Proteus mirabilis* (B-Pmi), *Proteus vulgaris* (B-Pvu), *Pseudomonas aeruginosa* (B-Pae), *Pseudomonas fluorescens* (B-Pfl), *Pseudomonas putida* (B-Ppu), *Pseudomonas syringae* (B-Psy), *Rhizobium melliloti* (B-Rme), *Rhizobium sp.* NGR234 (B-Rsp), *Rhodobacter capsulatus* (B-Rca), *Rhodobacter sphaeroides* (B-Rsp), *Rhodobacter sulfidophilus* (B-Rsu), *Rhodospseudomonas blastica* (B-Rbl), *Rhodospirillum rubrum* (B-Rru), *Rickettsia japonicum* (B-Rja), *Rickettsia rickettsii* (B-Rri), *Rickettsia typhi* (B-Rty), *Salmonella typhi* (B-Sti), *Salmonella typhimurium* (B-Sty), *Shigella flexneri* (B-Sfl), *Spiroplasma citri* (B-Sci), *Staphylococcus aureus* (B-Sau), *Staphylococcus carnosus* (B-Scs), *Streptococcus pneumoniae* (B-Spn), *Streptomyces clavuligerus* (B-Scl), *Streptomyces coelicolor* (B-Sco), *Synechocystis* PCC 6803 (B-Syn), *Thermus aquaticus* (B-Taq), *Thermus thermophilus* (B-Tth), *Thiobacillus cuprinus* (B-Tcu), *Treponema hyodysenteriae* (B-Thy), *Vibrio alginolyticus* (B-Val), *Vibrio cholera* (B-Vch), *Vibrio parahaemolyticus* (B-Vpa), *Vibrio proteolyticus* (B-Vpr), *Wolbachia sp.* (B-Wsp), *Yersinia enterocolitica* (B-Yen), *Zooglea ramigera* (B-Zra), *Zymomonas mobilis* (B-Zmo). Archaea: *Methanococcus jannaschii* (A-Mja), *Sulfolobus acidocaldarius* (A-Sac). Eukaryotes: *Apis mellifera* (E-Ame), *Arabidopsis thaliana* (E-Ath), *Atratyloides japonica* (E-Aja), *Bos taurus* (E-Bta), *Candida albicans* (E-Cal), *Caenorhabditis elegans* (E-Cel), *Dicytostellium discoideum* (E-Ddi), *Flaveria trinervia* (E-Ftr), *Giardia theta* (E-Gth), *Glycine max* (E-Gma), *Haematobia irritans* (E-Hir), *Homo sapiens* (E-Hsa), *Marchantia polymorpha* (E-Mpa), *Mus musculus* (E-Mmu), *Prototheca wickerhamii* (E-Pwi), *Petunia hybrida* (E-Phy), *Pisum sativum* (E-Psa), *Porphyra purpurea* (E-Ppu), *Odontella sinensis* (E-Osi), *Reclinomonas americana* (E-Ram), *Rattus norvegicus* (E-Rno), *Rhizopus oryzae* (E-Ror), *Saccharomyces cerevisiae* (E-Sce), *Schizosaccharomyces pombe* (E-Spo), *Solanum tuberosum* (E-Stu), *Spinacia oleracea* (E-Sol).

Energy metabolism

Early in its infectious cycle, *R. prowazekii* uses the ATP of the host with the help of membrane-bound ATP/ADP translocases. However, *R. prowazekii* is also capable of generating ATP, which may compensate for the gradual depletion of cytosolic ATP later in the infection. *R. prowazekii*'s repertoire of genes involved in ATP production and transport include determinants for the TCA cycle, the respiratory-chain complexes, the ATP-synthase complexes and the ATP/ADP translocases (Table 1). Genes to support anaerobic glycolysis are absent.

Pyruvate dehydrogenase. Pyruvate is imported into mitochondria directly from the cytoplasm and converted into acetyl-CoA by pyruvate dehydrogenase. The genes encoding three components (E1–E3) of the pyruvate dehydrogenase complex are found in *R. prowazekii*, indicating that it too uses cytosolic pyruvate. Pyruvate dehydrogenase (E1) consists of two subunits (α and β) in *R. prowazekii*, mitochondria and Gram-positive bacteria; the corresponding genes are clustered in the genome. In contrast, proteobacteria such as *E. coli*, *Haemophilus influenzae* and *Helicobacter pylori* have a single subunit for the E1 component and these have little similarity to the α and β subunits of the E1 component in *R. prowazekii* and mitochondria (data not shown).

Two paralogous genes code for the dihydrolipoamide dehydrogenase (E3) in *R. prowazekii*. One of these most resembles mitochondrial homologues, whereas the other is most similar to bacterial homologues (data not shown). The presence of several paralogous gene families for pyruvate dehydrogenases complicates attempts to reconstruct a genome phylogeny based on these genes.

ATP production. Genes encoding all enzymes in the TCA cycle are found in *R. prowazekii*. Proton translocation is mediated by NADH dehydrogenase (complex I), cytochrome reductase (complex III) and cytochrome oxidase (complex IV). Several clusters of genes code for components of the NADH dehydrogenase complex. Seven of these genes (*nuoJKLM* and *nuoGHO*) are located near to each other, but the order of genes is inverted relative to the order of this cluster in *E. coli*. An additional set of five genes is grouped in the order *nuoABCDE*, but the single genes *nuoF* and *nuoN* are distant from both of these clusters. Several proteins in the cytochrome *bc*₁ reductase complex, such as ubiquinol–cytochrome *c* reductase

(encoded by *petA*), cytochrome *b* (encoded by *cytb*) and cytochrome *c*₁ (encoded by *fbhC*), are present, as are several subunits of the cytochrome oxidase complex.

The ATP-synthesizing complex is composed of the ATP synthase F₁ component (comprising five polypeptides, α , β , γ , ϵ and δ) and the F_o component, a hydrophobic segment that spans the inner mitochondrial membrane. The genes encoding these components are normally clustered in one of the most highly conserved operon structures in microbial genomes. In *R. prowazekii*, however, the ATP-synthase genes encoding the α , β , γ , δ and ϵ subunits of the F₁ complex (*atpH*, *atpA*, *atpG*, *atpD* and *atpC*) are clustered in the common order, but *atpB*, *atpE* and *atpF*, encoding the A, B and C chains of the F_o complex, are split from this cluster.

Replication, repair and recombination

R. prowazekii has a smaller set of genes involved in DNA replication than do free-living bacteria such as *E. coli*, *Haemophilus influenzae* and *Helicobacter pylori*. Four genes have been identified that code for the core structure of DNA polymerase III, which includes the α (*dnaE*), ϵ (*dnaQ*), β (*dnaN*), γ and θ (*dnaX*) subunits. Extra subunits present in the *E. coli* DNA polymerase III are missing from *R. prowazekii*, as well as from *M. genitalium* and *B. burgdorferi*.

Genes encoding DNA-repair mechanisms are similar in the small genomes of the parasites *R. prowazekii*, *M. genitalium* and *B. burgdorferi*. Thus, genes involved in the repair of ultraviolet-induced DNA damage (*uvrABCD*) have been identified in all three genomes. In *R. prowazekii*, DNA-excision repair probably occurs by a pathway involving endonuclease III, polII and DNA ligase, as in *B. burgdorferi*.

The *R. prowazekii* genome has a limited capacity for mismatch repair. The DNA-mismatch-repair enzymes encoded by *mutL* and *mutS* are present, but *mutH* and *mutY* are not. There is a complete lack of *mut* genes in *M. genitalium*, but *mutL* and *mutHLY* have been identified in *B. burgdorferi* and *Chlamydia trachomatis*. The transcription-repair coupling factor (encoded by *mfd*) is found in *R. prowazekii*, *B. burgdorferi* and *C. trachomatis* but not in *M. genitalium*.

The *R. prowazekii* genome contains several genes involved in homologous recombination, such as *recA*, *recF*, *recJ*, *recN* and *recR*. A similar set of genes has been found in *A. aeolicus*²¹. The *rec* genes

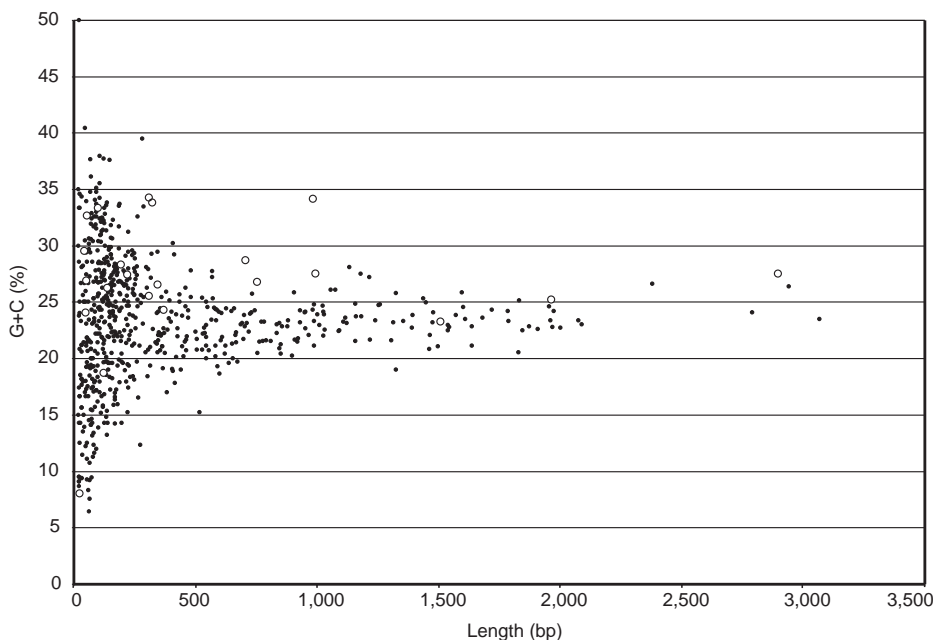


Figure 3 G+C content in intergenic regions longer than 20 bp in the *R. prowazekii* genome. The empty circles correspond to spacer sequences located at 886 to 916 kb, a region with an unusually large fraction of non-coding DNA and pseudogenes.

are scattered in the other small genomes of parasites. The *RecBCD* complex is missing in *R. prowazekii*, *M. genitalium* and *Helicobacter pylori* but it has been identified in *B. burgdorferi*.

Transcription and translation

R. prowazekii has three subunits (α , β and β') of the core RNA polymerase, as well as σ^{70} and one alternative σ factor, σ^{32} , which controls transcription of the genes encoding heat-shock proteins in *E. coli*. Genes involved in transcription elongation and termination, *nusA*, *nusB*, *nusG*, *greA* and *rho*, are also present. The gene encoding σ^{32} is absent in most other small genomes, such as those of *B. burgdorferi*, *Helicobacter pylori*, *M. genitalium* and *C. trachomatis*, although genes for heat-shock proteins are present.

An unusually large number of genes involved in RNA degradation are found in *R. prowazekii*. Of these, only four appear to be common to the bacterial genomes analysed so far (those encoding polyribonucleotide nucleotidyltransferase and ribonucleases HII, III and P). Four more ribonucleases (D, E, HI and PH) are present in *R. prowazekii*, but in none of the other small parasites.

Of the 33 identified tRNA genes, which code for 32 different tRNA isoacceptor species, two code for tRNA^{Phe}. There are two tRNA species for most of the amino acids that are encoded by four-codon boxes; the exceptions are the four-codon boxes for proline and valine, for which we have identified only one isoacceptor-tRNA species, with U in the first anticodon position. *selC*, which codes for tRNA^{Sec}, and *selABD* are missing. *R. prowazekii* has a set of genes coding for tRNA modifications (*tgt*, *queA*, *trmD*, *truA*, *truB* and *miaA*) which resembles that of *Helicobacter pylori*, *C. trachomatis* and *B. burgdorferi*; *M. genitalium* has only *trmD* and *truA*.

In *R. prowazekii*, 21 genes encode 18 of the 20 aminoacyl-tRNA synthetases normally required for protein synthesis. There are two genes (*gltX*) encoding glutamyl-tRNA synthetase. As seen in several bacterial genomes²⁵, the gene coding for glutaminyl-tRNA synthetase, *glnS*, is missing. Three genes encoding subunits of the glutamyl-tRNA amidotransferase are present, indicating that a glutamyl-tRNA charged with glutamic acid may be transamidated to generate Gln-tRNA. The gene coding for asparaginyl-tRNA synthetase, *asnS*, is also missing from the *R. prowazekii* genome as well as from *Helicobacter pylori*, *C. trachomatis* and *A. aeolicus*²⁶. A transamidation process to form Asn-tRNA^{Asn} from Asp-tRNA^{Asn} has been proposed for the archaeon *Haloferax volcanii*²⁷ and this reaction may also occur in *R. prowazekii*. The valyl-tRNA synthetase is 38.3% identical to its homologue in *Methanococcus jannaschii*, but only 27.6% identical to its most similar homologue in bacteria, which is found in *Bacillus stearothermophilus*, possibly indicating a horizontal transfer event. The lysyl-tRNA synthetase (encoded by *lysS*) in *R. prowazekii* is a class I enzyme with no resemblance to the conventional class II lysyl-tRNA synthetases. Class I type of lysyl-tRNA synthetases have been observed previously in only *B. burgdorferi*, *Pyrococcus woesei*, *Methanococcus jannaschii* and a few other methanogens²⁶.

Regulatory systems

As in other genomes of small parasites, *R. prowazekii* has a reduced set of regulatory genes. There are a few members of two-component regulatory systems, such as the proteins encoded by *barA*, *envZ*, *ntrY*, *ntrX*, *ompR* and *phoR*. *spoT*, which is involved in the stringent response, has been identified in *B. burgdorferi*, *Helicobacter pylori* and *M. genitalium*. Only remnants of genes coding for amino-terminal fragments of proteins similar to those encoded by *spoT* and *relA* are identifiable in *R. prowazekii*. No fragments of *spoT* encoding the carboxy-terminal segments of the protein have been identified in the genome.

Cell division and protein secretion

Proteins involved in detoxification, such as superoxide dismutase, and those involved in thiophen and furan oxidation are present in *R.*

prowazekii. Two genes encoding haemolysins have also been identified, and an *R. typhi* homologue of *tlyC* exhibits haemolytic activities when expressed in *E. coli* (S. Radulovic, J. M. Troyer, B. Noden, S.G.E.A. and A. Azad, unpublished observations).

The data indicate that the basic mechanisms of cell division and secretion in *R. prowazekii* are similar to those in free-living proteobacteria. There is a common set of bacterial chaperones (encoded by *dnaK*, *dnaJ*, *hslU*, *hslV*, *groEL*, *groEL*, *groES* and *hspG*) and genes involved in the *secA*-dependent secretory system (*secABDEFGY*, *ffh* and *ftsY*). *R. prowazekii* has a significantly larger set of genes involved in peptide secretion than does *M. genitalium*.

Membrane-protein analysis

Many studies of *R. prowazekii* have focused on outer-surface membrane proteins because of their potential importance in bacterial detection and vaccination. The superficial lipopolysaccharide (LPS) molecule is important in the pathogenesis of *R. prowazekii*. LPS consists of a polysaccharide that is covalently linked to lipid A, the biosynthesis of which is catalysed by products of *lpxABCD*, all of which are present in the *R. prowazekii* genome. These genes are clustered in *E. coli*, but *lpxA* and *lpxD* are separate from *lpxB* and *lpxC* in *R. prowazekii*. Three genes involved in the biosynthesis of the 3-deoxy-D-manno-octulosonic acid (KDO) residues reside in the *R. prowazekii* genome (*kdsA*, *kdsB* and *kdtA*). Only one gene (*rfaJ*) with a putative function in outer-core biosynthesis has been identified.

We have identified a set of genes involved in the biosynthesis of murein and diaminopimelate and a set involved in the biosynthesis of fatty acids. These includes: *fabD*, which is involved in the last step of the initiation phase of fatty-acid biosynthesis; four genes involved in the elongation cycle of fatty-acid biosynthesis (*fabFGHI*); and three genes involved in the first three steps of the synthesis of polar head groups (*cdsA*, *pssA* and *pgsA*). Finally, post-translational processing and addition of lipids to an N-terminal cysteine require the gene products prolipoprotein diacylglycerol transferase (*lgt*), prolipoprotein signal peptidase (*lspA*) and apolipoprotein:phospholipid N-acyl transferase (*lnt*). These are found in the genome with several genes involved in the degradation of fatty acids, such as *fadA* which encodes the 3-ketoacyl-CoA thiolase.

Virulence

The *R. prowazekii* genome contains several homologues of the *VirB* gene operon found in *Agrobacterium tumefaciens*. This gene family encodes proteins that direct the export of the T-DNA-protein complex across the bacterial envelope to the plant nuclei²⁸. *R. prowazekii* has two homologues of *VirB4* and one homologue each of *VirB8*, *VirB9*, *VirB10*, *VirB11* and *VirD4*. The latter five genes are clustered with the gene *trbG*, which is involved in conjugation in *Agrobacterium tumefaciens*. Homologues of the single-stranded DNA-binding proteins *VirD2* and *VirE2* are missing. In *Agrobacterium tumefaciens*, these proteins are bound to the transferred T-DNA, indicating different functions for the homologues of the *VirB* genes in *R. prowazekii*. Indeed, *VirB* proteins are homologous to components of the *E. coli* transport system for plasmids, as well as to components of the Pt1 transport machinery in *Bordetella pertussis*, which exports pertussis toxin²⁸. A set of genes coding for *VirB4* and several other *VirB* proteins has been identified in the *cag* pathogenicity island of *Helicobacter pylori*. In this species, the *VirB* proteins facilitate export of a factor that induces interleukin-8 secretion in gastric epithelial cells²⁸. Thus, *R. prowazekii* may encode components of a transport system for both conjugal DNA transfer and protein export.

The virulence of *Staphylococcus aureus* has been correlated with the production of capsular polysaccharides in phagocytic assays and mouse lethality assays^{29,30}. A cluster of ten capsule genes (*capA-M*) is involved in capsule biosynthesis in *S. aureus* strain M³¹. We have identified three *R. prowazekii* genes with sequence similarities to *S. aureus cap* genes. Two of these (*capD* and *capM*) are separated by ten

genes, most of which are unknown genes or genes involved in the biosynthesis of LPS or teichoic acid. Thus, *R. prowazekii* may produce components of a microcapsular layer that is involved in virulence.

Reductive evolution

Genome sequences of organisms enjoying an endosymbiotic lifestyle are at risk. The activities of homologous nuclear genes may render genes of the endosymbiont expendable and as a consequence they become vulnerable to obliteration by mutation. Good candidates for such purged genes in *Rickettsia* and mitochondria are genes required for amino-acid biosynthesis, nucleoside biosynthesis and anaerobic glycolysis. These and other genes would have been deleted when an ancestral genome first lived in a nucleated cell. Once genes essential to a free-living mode are lost, the endosymbiont becomes an obligate resident of its host.

Likewise, small, bottle-necked populations of bacteria infecting a eukaryotic cell will tend to accumulate deleterious mutations because selection cannot remove them from such clonal populations¹³. The accumulation of such harmful but non-lethal mutations is referred to as ‘Muller’s ratchet’³² or ‘near-neutral evolution’^{33,34}. The consequence of accumulation of these mutations will be the inactivation and eventual deletion of non-essential genes.

The first mutation that inactivates an expendable gene is likely to initiate a sequence of events in which subsequent mutations freely transform it, by degrees, from a pseudogene, to unrecognizable sequence, to small fragments, to extinction. In this sequence, mutations are released from amino-acid-coding constraints. Thus nucleotide substitutions will reflect the mutation bias of the genome. This bias can be estimated roughly by frequencies of third-position bases in the codons. For *R. prowazekii*, the bias of the third-position bases is 18% G+C rather than the 29% G+C average for the genome. So, as sequences age in *R. prowazekii*, their composition should gradually approach the low G+C content of third codon positions. Nearly one-quarter of the *R. prowazekii* genome is composed of non-coding sequences, with a G+C content lower than that of coding sequences (25% G+C compared to 30%; $P < 0.001$). Thus, much of the non-coding sequence may be remnants of coding sequences that are in the process of being eliminated from the genome.

The gene encoding *S*-adenosylmethionine synthetase (*metK*), which catalyses the biosynthesis of *S*-adenosylmethionine (SAM), illustrates the initiation of this process. The *metK* sequence in the strain of *R. prowazekii* studied here has a termination codon within a region of the gene that is otherwise highly conserved among

bacterial species³⁵. However, a closely related strain does not have the termination codon. Many other defects, such as termination codons, insertions, and a preponderance of small deletions, have also been observed in the *metK* genes in several members of the spotted fever group *Rickettsia* (J.O.A. and S.G.E.A., unpublished observations). This random distribution of lethal mutations among some *metK* alleles from different *Rickettsia* species indicates that the gene may have just entered the extinction process. This distribution, and the identification of 11 more pseudogenes for carboxypeptidase (*ypwA*), penicillin-binding protein (*pbpC*), succinyl CoA-transferase (*scoB*), transposase (*tra3*), resolvase (*pin*), conjugative transfer protein (*taxB*), a hypothetical protein (*yfc1*) and four different fragmented open reading frames for (p)ppGpp 3'-pyrophosphohydrolase, indicates that the *R. prowazekii* genome continues to eliminate genes.

Genome sequences can be purged by a more abrupt mechanism. This consists of intrachromosomal recombination at duplicated sequences, which can result in the deletion of intervening sequences, the loss of a sequence duplication and the rearrangement of flanking

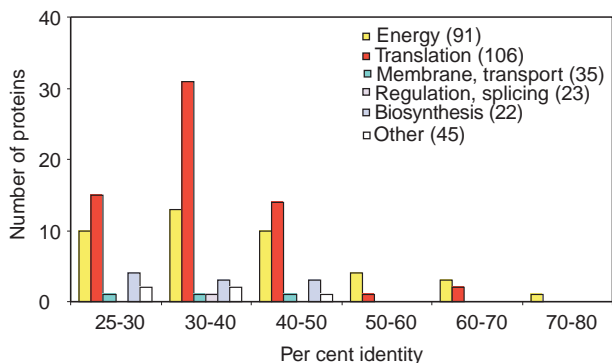


Figure 4 Histogram representation of the similarity of predicted *R. prowazekii* proteins to yeast proteins targeted to the mitochondria. Only protein pairs with per cent identity values greater than 25% are shown. Numbers in parentheses represent the total number of yeast mitochondrial proteins within each category. The yeast mitochondrial protein sequences have been taken from <http://www.proteome.com>.

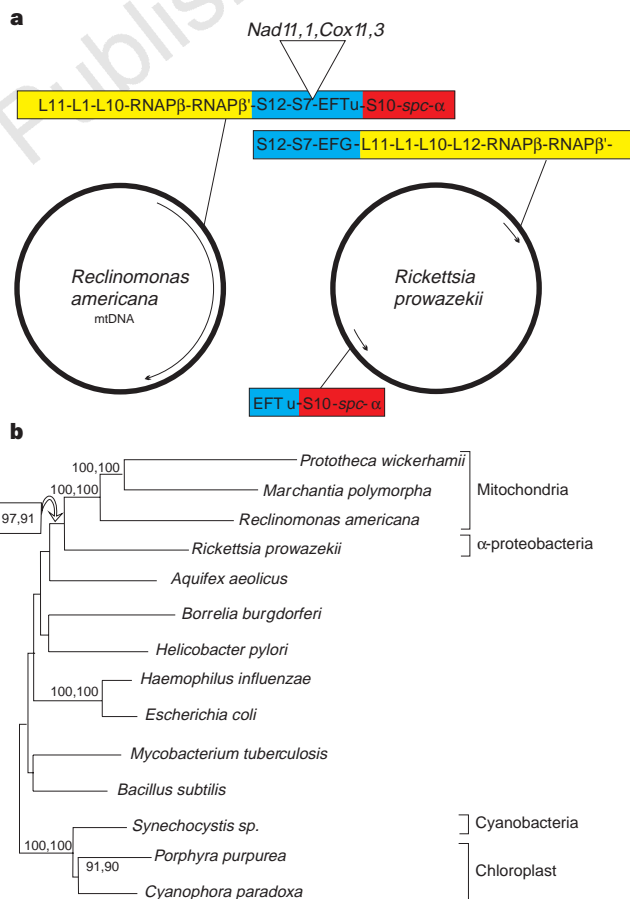


Figure 5 The organization and phylogenetic relationships of gene encoding ribosomal protein from *R. prowazekii* and the mitochondrial genome of *Reclinomonas americana*. **a**, The organization of ribosomal-protein genes. The *S10*, *spc* and α -operons are organized similarly in these two genomes, except that several ribosomal-protein genes³⁸ have been deleted from the mitochondrial genome of *Reclinomonas americana*. **b**, The phylogenetic relationships of mitochondria and bacteria were derived from the combined amino-acid sequences of ribosomal proteins S2, S3, S7, S10, S11, S12, S13, S14, S19, L5, L6 and L16. Neighbour-joining and maximum-parsimony methods gave identical topologies. Branch lengths are proportional to those reconstructed by using the neighbour-joining method. Values at nodes are bootstrap values indicating the degree of support for individual clusters under each method (neighbour-joining, maximum parsimony). Only bootstrap values >90% are shown.

sequences. Such a mechanism will account for the presence in *R. prowazekii* of one, unlinked copy of *rrs* and *rrl*, both of which are surrounded by new flanking sequences³⁶. Likewise, *R. prowazekii* has one *tuf* gene and one *fus* gene in atypical clusters that seem to have been created by intrachromosomal recombination between the two *tuf* genes that are normally found in Gram-negative bacteria³⁷. Indeed, rearranged gene operon structures encoding ribosomal proteins are characteristic of all members of the genus *Rickettsia* (H. Amiri, C.A. and S.G.E.A., unpublished observations).

Conserved operons that are found in free-living bacteria are often dispersed throughout the *Rickettsia* genome (see above). The *R. prowazekii* genome contains an unusually small fraction of repeat sequences (<10% of that observed in free-living bacteria). We suggest that the repeat sequences found in the ancestor to the *Rickettsia* have been 'consumed' by the intrachromosomal-recombination mechanism that generated some of the deletions and rearrangements seen in *R. prowazekii*. Such intrachromosomal recombinants arise at a substantial rate in bacteria growing in culture, but here they are eliminated from the populations by selection. That such remnants of intrachromosomal recombination are retained in *R. prowazekii* indicates that purifying selection has been attenuated in this organism.

Mitochondrial affinities

The reduction in genome size in mitochondria and *Rickettsia* is likely to have occurred independently in the two lineages. Most of

the genes supporting mitochondrial activities are nuclear. Many of the 300 proteins encoded in the nucleus of the yeast *Saccharomyces cerevisiae* but destined for service within the mitochondrion are close homologues of their counterparts in *R. prowazekii*. Nearly one-quarter of these proteins are required for bioenergetic processes and another one-third of them are required for the expression of the genes encoded in the mitochondrial genome. In total, more than 150 nucleus-encoded mitochondrial proteins share significant sequence homology with *R. prowazekii* proteins (Fig. 4).

Another group of 58 nucleus-encoded mitochondrial proteins represents components of the mitochondrial transport machinery and regulatory system (Fig. 4). These include proteins found in the mitochondrial outer membrane and others involved in splicing reactions. Such proteins have probably been secondarily recruited to mitochondria from genomes not necessarily related to that of the α -proteobacterial ancestor.

The mitochondrial genome of the early diverging, freshwater protozoan *Reclinomonas americana* is more like that of a bacterium than any other mitochondrial genome sequenced so far³⁸. This genome contains 67 protein-coding genes, most of which provide components of genetic processes and the bioenergetic system³⁸. Several gene clusters in this mitochondrial genome are reminiscent of those in bacteria (Figs 5a, 6a). Most similarities represent retained, ancestral traits present in the common ancestor of bacteria and mitochondria. For example, the genes *rplKAJL* and *rpoBC* are identically organized in *R. prowazekii* and the mitochondrial genome of *Reclinomonas americana*. Likewise, the genes encoding the S10, *spc* and the α -ribosomal protein operons are organized similarly in the two genomes. The immediate proximity of these two clusters in the *Reclinomonas americana* mitochondrial DNA is reminiscent of the arrangement in free-living bacteria, whereas the physical separation of the two clusters in the *R. prowazekii* genome is atypical. A further rearrangement event is indicated by the fact that the *rpsLrpsGfus* cluster is located upstream of the *rplKAJLrpoBC* cluster in *R. prowazekii*, rather than downstream as it is in the *Reclinomonas americana* mtDNA. Phylogenetic reconstructions based on ribosomal proteins within each of these two clusters indicate that there is a close evolutionary relationship between *R. prowazekii* and mitochondria (Fig. 5b).

Mitochondria and *R. prowazekii* have a similar repertoire of proteins involved in ATP production and transport, including genes encoding components of the TCA cycle, the respiratory-chain complexes, the ATP-synthase complexes and the ATP/ADP translocases. There are some similarities in the gene orders of some functional clusters (Fig. 6a). There are also some rearrangements of clusters that are specific to *Rickettsia*. One example is the inversion of segments corresponding to *nuoJKLM* and *nuoGHI*. Another is the scattered displacement of genes involved in the biogenesis of cytochrome *c*. Nevertheless, phylogenetic reconstructions based on components of the NADH dehydrogenase complexes indicate that there is a close evolutionary relationship between *R. prowazekii* and mitochondria (Fig. 6b).

We have identified as many as five genes coding for ATP/ADP transporters, all of which are expressed (R.M.P. *et al.*, unpublished observations). The *Rickettsia* ATP/ADP translocases are monomers with 12 transmembrane regions each, whereas the mitochondrial translocases are dimers with six transmembrane regions per dimer. We found no relationship between the primary structures of the mitochondrial and *Rickettsia* ATP/ADP translocases, indicating that these transport systems may have originated independently.

The study of the *R. prowazekii* genome sequence supports the idea that aerobic respiration in eukaryotes originated from an ancestor of the *Rickettsia*, as indicated previously by phylogenetic reconstructions based on the rRNA gene sequences^{7,9}. Phylogenetic analyses of the *petB* and *coxA* genes indicate that the respiration systems of *Rickettsia* and mitochondria diverged ~1,500–2,000 million years ago¹⁰, shortly after the amount of oxygen in the atmosphere began

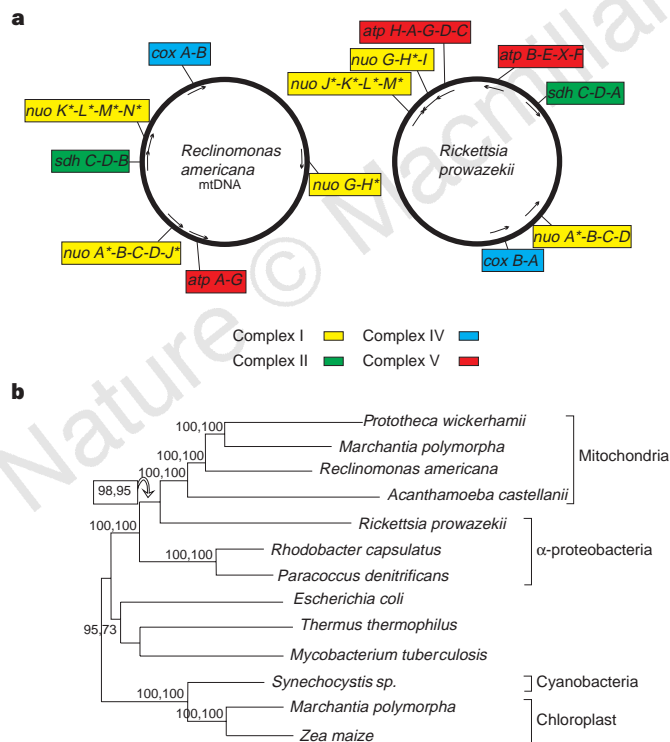


Figure 6 The organization and phylogenetic relationships of genes involved in ATP synthesis from *R. prowazekii* and the mitochondrial genome of *Reclinomonas americana*. **a**, The organization of bioenergetic genes. **b**, The phylogenetic relationships of mitochondria and bacteria were derived from the combined amino-acid sequences of NADH dehydrogenase I chains A, J, K, L, M and N, which are encoded by the genes *nuoA*, *J*, *K*, *L*, *M*, *N*. These genes are highlighted by asterisks in **a**. Neighbour-joining and maximum-parsimony methods gave identical topologies. Branch lengths are proportional to those reconstructed using the neighbour-joining method. Values at nodes are bootstrap values indicating the degree of support for individual clusters under each method (neighbour-joining, maximum parsimony). Only bootstrap values >90% are shown.

to increase. The finding that the ATP/ADP translocases in *R. prowazekii* and mitochondria are of different evolutionary origin is problematic (R.M.P. *et al.*, unpublished observations). Free-living bacteria do not seem to have homologues of ATP/ADP translocases, which are found only in organelles and in two obligate intracellular parasites, *Rickettsia* and *Chlamydia*. Thus it is not known whether the original endosymbiont was capable of efficient exchange of adenosine nucleotides with its host cell. More detailed comparative analysis of the genomes of α -proteobacteria may refine our understanding of the origins of mitochondria. □

Methods

Genome sequencing. We prepared genomic DNA from the Madrid E strain of *R. prowazekii*, which was originally isolated in Madrid from a patient who died in 1941 with epidemic typhus. We propagated *R. prowazekii* in the yolk sac of embryonated hen eggs and purified DNA according to standard procedures³⁹. We sequenced the *R. prowazekii* genome by a whole-genome shotgun approach in combination with shotgun sequencing of a selected set of clones from a cosmid library (A.Z. *et al.*, unpublished observations). Genomic and cosmid DNA was sheared by nebulization to an average size of ~2 kb. The random fragments were cloned into a modified M13 vector using the double adaptor method⁴⁰. We collected 19,078 sequence reads during the random sequencing phase using Applied Biosystems 377 DNA sequencers (Perkin-Elmer).

The sequences were assembled and the consensus sequence was edited using the STADEN program⁴¹. We verified the structure of the assembled sequence by end-sequencing of 3-kb-insert λ Zap II clones³⁶, 10-kb λ clones and 30-kb cosmid clones. More than 97% of the genome was covered by clones from the three different libraries (A.Z. *et al.*, unpublished observations). Gaps between contigs were closed by direct sequencing of clones from the three libraries or of polymerase chain reaction (PCR) products. The final four gaps were closed by direct sequencing of PCR products generated with the Long Range PCR system (Gene Amp). Regions of ambiguity were identified by visual inspection of the assembly and resequenced. The final assembly contains ~20,000 sequences. The genome sequence has eightfold coverage on average and no single region has less than twofold coverage. We estimate the overall error frequency to be $< 1 \times 10^{-5}$.

Informatics. Sequence analysis and annotation was managed by CapDB (T.S.-P. *et al.*, unpublished observations). We identified open reading frames of more than 50 codons as genes on the basis of their characteristic patterns in nucleotide-frequency statistics¹⁴ using BioWish⁴². The identified genes were analysed using the program BLASTX⁴³ to search for sequence similarities in EMBL, TrEMBL, SwissProt and in-house databases. We identified tRNA genes with the program tRNA scan-SE⁴⁴. Remaining frameshifts were considered to be authentic and annotated as pseudogenes. Families of paralogues were constructed using BLAST to search for sequence similarities within the *R. prowazekii* genome. Multiple alignments and phylogenetic trees for genes with significant sequence similarities to genes in the public databases were constructed automatically using CLUSTAL-W⁴⁵, Phylo_win⁴⁶ and GRS⁴⁷. The final annotation was based on manual inspection of the phylogenetic placement of *R. prowazekii* in the resulting gene trees.

Received 21 July; accepted 24 September 1998.

1. Gross, L. How Charles Nicolle of the Pasteur Institute discovered that epidemic typhus is transmitted by lice: reminiscences from my years at the Pasteur Institute in Paris. *Proc. Natl Acad. Sci. USA* **93**, 10539–10540 (1996).
2. Weisburg, W. G., Woese, C. R., Dobson, M. E. & Weiss, E. A common origin of Rickettsiae and certain plant pathogens. *Science* **230**, 556–558 (1985).
3. Woese, C. R. Bacterial evolution. *Microbiol. Rev.* **51**, 221–227 (1987).
4. Weisburg, W. G. *et al.* Phylogenetic diversity of the rickettsias. *J. Bacteriol.* **171**, 4202–4206 (1989).
5. Yang, D., Oyaizu, Y., Oyaizu, H., Olsen, G. J. & Woese, C. R. Mitochondrial origins. *Proc. Natl Acad. Sci. USA* **82**, 4443–4447 (1985).
6. Gray, M. W., Cedergren, R., Abel, Y. & Sankoff, D. On the evolutionary origin of the plant mitochondrion and its genome. *Proc. Natl Acad. Sci. USA* **86**, 2267–2271 (1989).
7. Olsen, G. J., Woese, C. R. & Overbeek, R. The winds of (evolutionary) change: breathing new life into microbiology. *J. Bacteriol.* **176**, 1–6 (1994).
8. Viale, A. & Arakaki, A. K. The chaperone connection to the origins of the eukaryotic organelles. *FEBS Lett.* **341**, 146–151 (1994).

9. Gray, M. W. & Spencer, D. F. in *Evolution of Microbial Life* (eds Roberts, D. M., Sharp, P. M., Alderson, G. & Spencer, D. F.) 109–126 (Cambridge Univ. Press, Cambridge, 1996).
10. Sicheritz-Pontén, T., Kurland, C. G. & Andersson, S. G. E. A phylogenetic analysis of the cytochrome *b* and cytochrome *c* oxidase I genes supports an origin of mitochondria from within the Rickettsiaceae. *Biochim. Biophys. Acta* **1365**, 545–551 (1998).
11. Margulis, L. *Origin of Eukaryotic Cells* (Yale Univ. Press, New Haven, 1970).
12. Kurland, C. G. Evolution of mitochondrial genomes and the genetic code. *Bioessays* **14**, 709–714 (1992).
13. Andersson, S. G. E. & Kurland, C. G. Reductive evolution of resident genomes. *Trends Microbiol.* **6**, 263–268 (1998).
14. Andersson, S. G. E. & Sharp, P. M. Codon usage and base composition in *Rickettsia prowazekii*. *J. Mol. Evol.* **42**, 525–536 (1996).
15. Fleischmann, R. D. *et al.* Whole-genome random sequencing and assembly of *Haemophilus influenzae* Rd. *Science* **269**, 496–511 (1995).
16. Fraser, C. M. *et al.* The *Mycoplasma genitalium* genome reveals a minimal gene complement. *Science* **270**, 397–403 (1995).
17. Himmelreich, R. *et al.* Complete genome sequence analysis of the genome of the bacterium *Mycoplasma pneumoniae*. *Nucleic Acids Res.* **3**, 109–136 (1996).
18. Fraser, C. M. *et al.* Genomic sequence of a Lyme disease spirochaete, *Borrelia burgdorferi*. *Nature* **390**, 580–586 (1997).
19. Tomb, J.-F. *et al.* The complete genome sequence of the gastric pathogen *Helicobacter pylori*. *Nature* **388**, 539–547 (1997).
20. Blattner, F. R. *et al.* The complete genome sequence of *Escherichia coli* K-12. *Science* **277**, 1453–1462 (1997).
21. Deckert, G. The complete genome of the hyperthermophilic bacterium *Aquifex aeolicus*. *Nature* **392**, 353–358 (1998).
22. Cole, S. T. *et al.* Deciphering the biology of *Mycobacterium tuberculosis* from the complete genome sequence. *Nature* **393**, 537–544 (1998).
23. Fraser, C. M. *et al.* Complete genome sequence of *Treponema pallidum*, the syphilis spirochete. *Science* **281**, 375–388 (1998).
24. Lobry, J. R. Asymmetric substitution patterns in the two DNA strands of bacteria. *Mol. Biol. Evol.* **13**, 660–665 (1996).
25. Ebeling, S., Kündig, C. & Hennecke, H. Discovery of a rhizobial RNA that is essential for symbiotic root nodule development. *J. Bacteriol.* **173**, 6373–6382 (1991).
26. Koonin, E. V. & Aravind, L. Genomics: re-evaluation of translation machinery evolution. *Curr. Biol.* **8**, 266–269 (1998).
27. Curnow, A. W., Ibba, M. & Soll, D. tRNA-dependent asparagine formation. *Nature* **382**, 589–590 (1996).
28. Christie, P. J. *Agrobacterium tumefaciens* T-complex transport apparatus: a paradigm for a new family of multifunctional transporters in eubacteria. *J. Bacteriol.* **179**, 3085–3094 (1997).
29. Melly, M. A., Duke, L. J., Liau, D.-F. & Hash, J. H. Biological properties of the encapsulated *Staphylococcus aureus*. *M. Infect. Immun.* **10**, 389–397 (1974).
30. Peterson, P. K., Wilkinson, B. J., Kim, Y., Schmeling, D. & Quie, P. G. Influence of encapsulation on staphylococcal opsonization and phagocytosis by human polymorphonuclear leukocytes. *Infect. Immun.* **19**, 943–949 (1978).
31. Lin, W. S., Cumen, T. & Lee, C. Y. Sequence analysis and molecular characterization of genes required for the biosynthesis of type I capsular polysaccharide in *Staphylococcus aureus*. *J. Bacteriol.* **176**, 7005–7016 (1994).
32. Felsenstein, J. The evolutionary advantage of recombination. *Genetics* **78**, 157–159 (1977).
33. Ohta, T. & Kimura, M. On the constancy of the evolutionary rate of cistrons. *J. Mol. Evol.* **1**, 18–25 (1971).
34. Ohta, T. Evolutionary rate of cistrons and DNA divergence. *J. Mol. Evol.* **1**, 150–157 (1972).
35. Andersson, J. O. & Andersson, S. G. E. Genomic rearrangements during evolution of the obligate intracellular parasite *Rickettsia prowazekii* as inferred from an analysis of 52 015 bp nucleotide sequence. *Microbiology* **143**, 2783–2795 (1997).
36. Andersson, S. G. E., Zomorodipour, A., Winkler, H. H. & Kurland, C. G. Unusual organization of the rRNA genes in *Rickettsia prowazekii*. *J. Bacteriol.* **177**, 4171–4175 (1995).
37. Andersson, S. G. E. & Kurland, C. G. Genomic evolution drives the evolution of the translation system. *Biochem. Cell Biol.* **73**, 775–787 (1995).
38. Lang, B. F. *et al.* An ancestral mitochondrial DNA resembling a eubacterial genome in miniature. *Nature* **387**, 493–497 (1997).
39. Winkler, H. H. Rickettsia permeability: an ATP/ADP transport system. *J. Biol. Chem.* **251**, 389–396 (1976).
40. Andersson, B. *et al.* A 'double adaptor' method for improved shotgun library construction. *Anal. Biochem.* **236**, 107–113 (1996).
41. Staden, R. The Staden sequence analysis package. *Mol. Biotech.* **5**, 233–241 (1996).
42. Sicheritz-Pontén, T. BioWish: a molecular biology command extension to Tcl/Tk. *Comput. Appl. Biosci.* **13**, 621–622 (1997).
43. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
44. Low, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
45. Thompson, J. D., Higgins, D. G. & Gibson, T. J. CLUSTAL W: improving the sensitivity of progressive multiple alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**, 4673–4680 (1994).
46. Galtier, N., Gouy, M. & Gautier, C. SeaView and Phylo_win, two graphic tools for sequence alignment and molecular phylogeny. *Comput. Appl. Biosci.* **12**, 543–548 (1996).
47. Sicheritz-Pontén, T. & Andersson, S. G. E. GRS: a graphic tool for genome retrieval and segment analysis. *Microb. Comp. Genomics* **2**, 123–139 (1997).

Acknowledgements. We thank C. Woese for discussions; M. Andersen for computer system support; and B. Andersson, K. Andersson, I. Tamas, B. Canbäck, A. Jamal, H. Amiri and S. Jossan for technical advice and assistance. This work was supported by the Swedish Foundation for Strategic Research, the Swedish Natural Sciences Research Council, the Knut and Alice Wallenberg Foundation and the European Commission.

Correspondence and requests for materials should be addressed to C.G.K. (e-mail: chuck@xray.bmc.uu.se).



