

# DOMAIN ARCHAEA

## Table of Contents

1. Glossary of terms
2. Discovery of domain *Archaea*
3. Structural attributes
  - Cell shape & size
  - Cell multiplication
  - Presence & type of cell wall structure
  - Membrane composition
4. Functional attributes
  - Habitat & ecology
  - Nutrition, physiology & metabolism
5. Molecular attributes
  - Genomes
  - Gene organisation in genomes
  - DNA replication
  - Transcription
  - Translation
  - Cloning and expression of archaeal genes
  - Phage and plasmids
6. Kingdom *Crenarchaeota*
  - Section I. Thermophilic and hyperthermophilic crenarchaeotes
    - Order "*Igneococcales*"
    - Order *Sulfolobales*
    - Order *Thermoproteales*
  - Section II. Cold dwelling crenarchaeotes
7. Kingdom *Euryarchaeota*
  - Order *Halobacteriales*, the extreme halophiles
  - Methanogens and its five orders
  - *Thermoplasmatales*, the cell wall-less order
  - *Archaeoglobales*, the sulfate reducers
  - *Thermococcales*, the sulfur respirers

8. Kingdom *Korarchaeota*
9. Evolution and Life at High Temperatures
10. The Limits to microbial existence:
11. Hyperthermophiles Archaea and Microbial Evolution:
12. Comparative genomics of *Archaea*:
13. Some Useful References

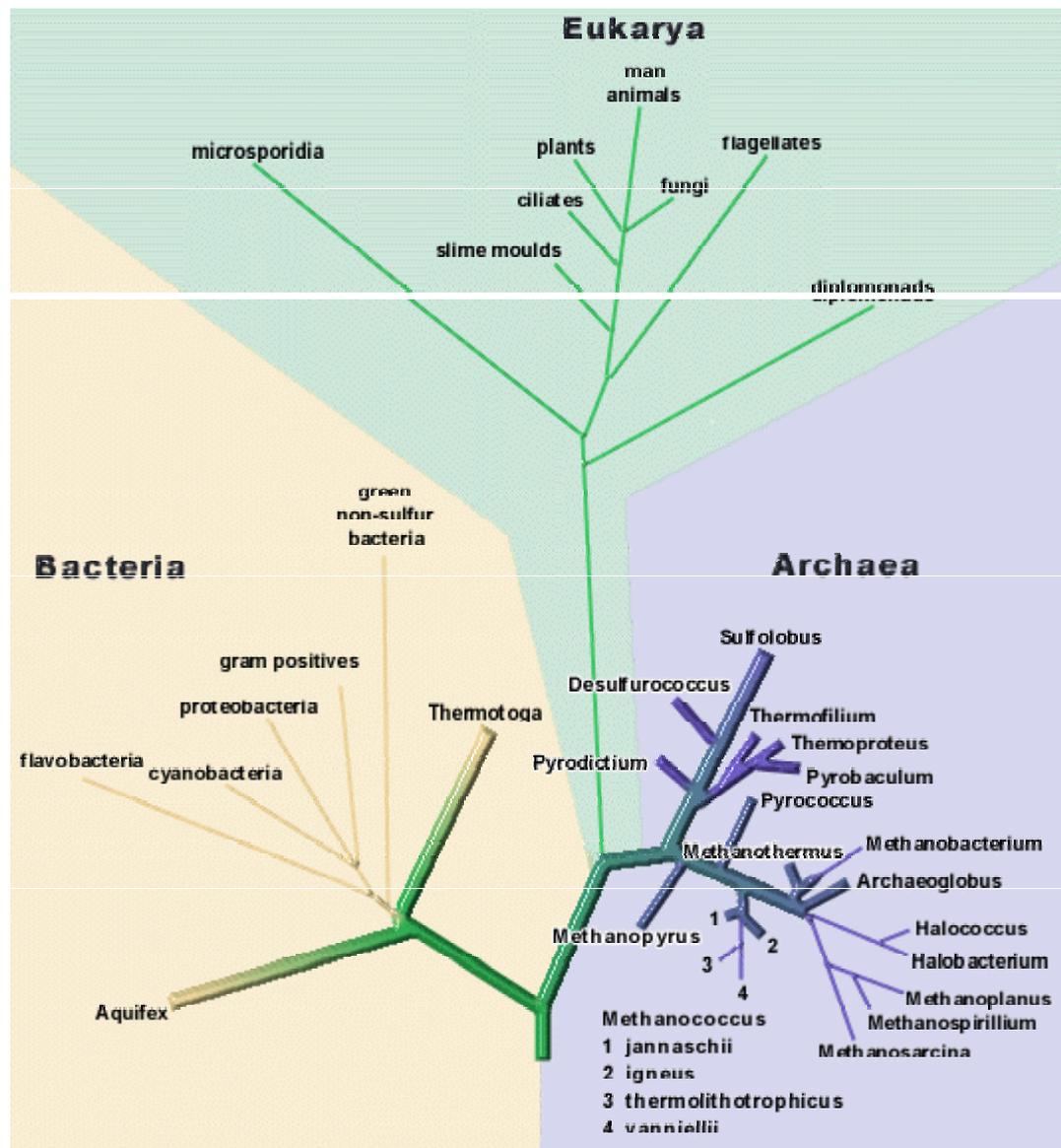
## 1. Glossary:

- Acetotroph: A methanogen which consumes acetate and splits acetate to methane and carbon dioxide during growth.
- Acetyl-CoA (Ljungdahl-Wood) pathway: A autotrophic CO<sub>2</sub> fixing pathway widespread in strict anaerobes (e.g. methanogens, homoacetogens & sulfate-reducing bacteria).
- Compatible Solute: Organic or inorganic substances that accumulate in halophilic cytoplasm for maintaining ionic pressure.
- Crenarchaeota: A kingdom of *Archaea* that contains hyperthermophiles and cold dwelling organisms.
- Euryarchaeota: A kingdom of *Archaea* that contains mainly methanogens, the extreme halophiles and Thermoplasma
- Extreme halophile: An organism whose growth is obligately dependent on high concentrations (> 10%) NaCl.
- Hyperthermophile: A microbe that grows optimally with temperatures > 80°C.
- Korarchaeota: A kingdom of *Archaea* that branches close to the archeal root.
- Reverse DNA gyrase: An enzyme present in hyperthermophiles that introduces positive super coiling into circular DNA.
- Solfatara: A hot, sulfur-rich but generally acidic environment.
- Thermosome: A type of heat shock chaperonin that refolds partially denatured proteins in hyperthermophiles.

## 2. Discovery of domain *Archaea*:

Until 1977, methanogens were regarded as bacteria. Based on 16S and 18S rRNA sequence data, Woese proposed a third kingdom to encompass them [Woese, (1977) PNAS 74:5088-5090]. In 1990, Woese concluded from further 16S rRNA and 18S rRNA sequences that *Halobacterium* regarded previously as a halophilic pseudomonad and *Sulfolobus* regarded as a gram-positive bacterium, were members of domain *Archaea*. He proposed that life on earth is made of 3 primary lineages which he referred to as domains [Woese,(1990) PNAS 87:4576-4579].

- *Eubacteria* (Eu = good or true)
- *Archeae* (Archeae = ancient) and
- *Eukarya* (Carya = nut or kernel)



The evolutionary history of life can be traced to the earliest common ancestor (progenote) for the three domains to be some 3.5 to 4 billion years (Brown and Doolittle (1995) PNAS 1995, 92: 2441-2445; Keeling and Doolittle (1995) PNAS 92:5761-5764; Doolittle (1999) Science 284:2124-2128] and it is expected that as we unravel the mysteries surrounding the evolution of life, the descriptions of life will change.

Members of the domain *Archaea* (aka as *Empire Archaea*) at the time of writing (June 2000) are phylogenetically divided into three kingdoms, namely, *Euryarchaeota*, *Crenarchaeota* and *Korarchaeota* (Fig 1)

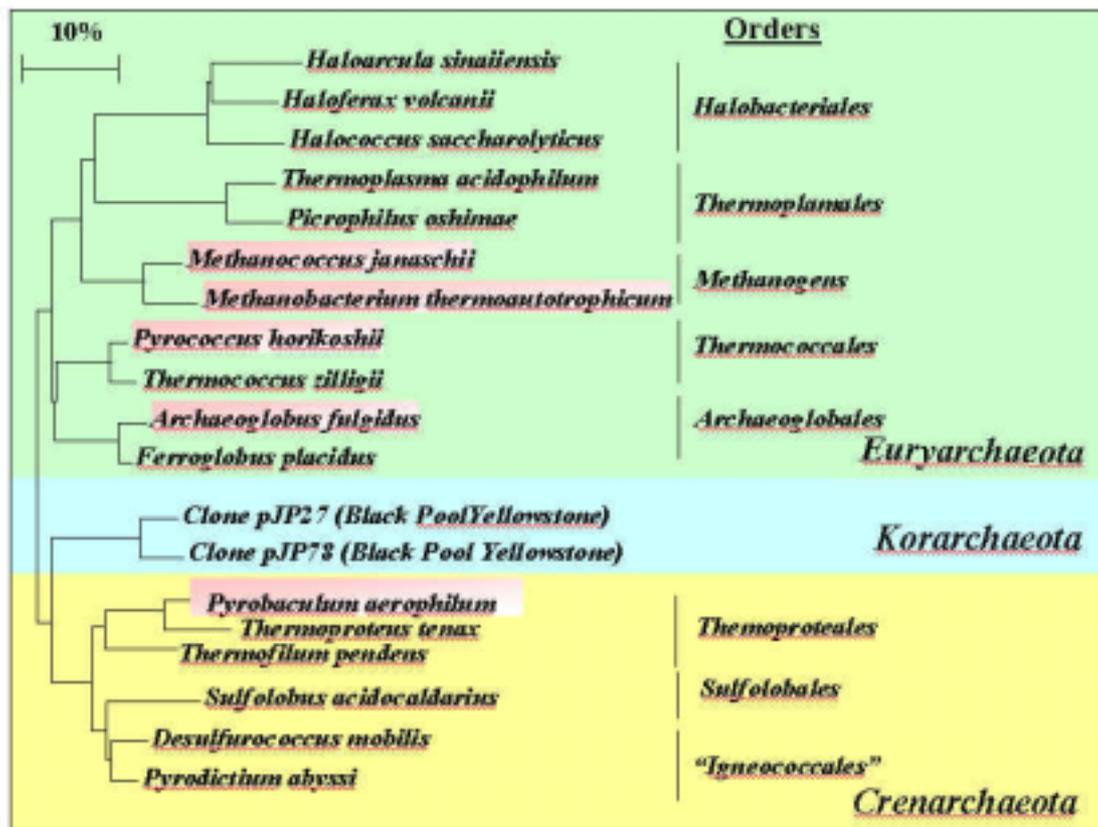


Figure 1 Phylogeny of domain Archaea based on comparison of the 16S rRNA sequences. Representatives of *Euryarchaeota* and *Crenarchaeota* have been cultured but members of the *Korarchaeotas* have yet to be cultured. Genomes of 4 euryarchaeotes and 1 crenarchaeote have been sequenced. The bar indicates nucleotide divergence. Greek *Archaios* = ancient, primitive; Greek *Eurus* = wide (wide distribution); Greek *Crene* = spring, fount (primary habitat).

### 3. Structural Attributes:

#### Shape and size

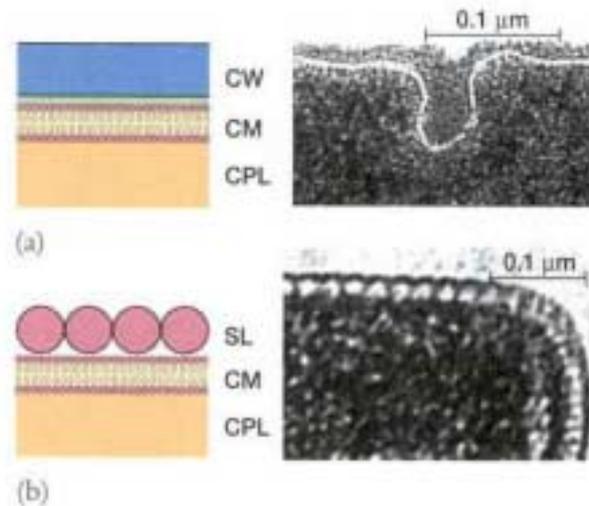
Members of domain *Archaea* are morphologically diverse and include spheres, spiral, rods, lobed, plate-shaped, irregular-shaped or pleomorphic. They may exist as single cells, as aggregates or form filaments. The diameter range is 0.1 to 15  $\mu\text{m}$  and the length can be up to 200  $\mu\text{m}$ .

#### Cell multiplication

Usually by binary fission, but some multiple by budding, fragmentation and as yet unknown mechanisms.

#### Cell Walls

Some *Archaea* such as *Thermoplasma* species do not contain a cell wall whereas most others do contain cell walls. The cell wall-containing *Archaea* can stain Gram-positive or Gram-negative and ultrastructurally are similar to that of members of domain Bacteria.



Schematic representations and electron micrograph of (a) a gram-positive archaeum (e.g. *Methanobacterium*) and a gram-negative archaeum (e.g. *Thermoproteus*). CW = cell wall, CM = cytoplasmic membrane, CPL = cytoplasm and SL = surface layer.

However, the chemistries are very different. In general, *Archaea* also possess more chemical variation in their cell walls than members of domain *Bacteria* do.

(a) Gram-positive archaeal cell walls: The ultrastructure structure of Archaeal Gram-positive shows a thick layer which is similar to the ultrastructure of Gram-positive *Bacteria*.

- *Methanobacterium*, *Methanothermus* and *Methanopyrus* contain pseudomurein (**glycans** [sugars] and **peptides** in their cell walls). **Glycans** are modified sugars viz, N-acetyl talosaminouronic acid (NAT or T) & N-acetyl glucose amine (NAG or G) T and G are linked to each other by a beta 1, 3 glycosidic bond & alternate to form the cell wall backbone. Lysozyme (an enzyme produced by organisms that consume bacteria, and normal body secretions such as tears, saliva, & egg white = protect against would-be pathogenic bacteria) cannot digest beta 1,3 glycosidic bonds. **Peptides** are short amino acid chains attached to T. The amino acids are only of the L-type. Penicillin is ineffective in inhibiting the cell wall peptide bridge formation.
- Some cell walls contain polysaccharides:
  - *Methanosarcina* are non-sulfated polysaccharide. These complex polysaccharides are similar to chondroitin sulfate (aka methanochondroitin) of animal connective tissue.
  - *Halococcus* are sulfated polysaccharides

(b) Gram-negative archaeal cell walls lack the outer membrane and complex lipopolysaccharide found in Gram-negative members of the domain *Bacteria* but instead consists of either a surface glycoprotein or protein subunits.

- *Halobacterium* are made of glycoproteins but also contain negatively charged acidic amino acids which counteract the positive charges of the high  $\text{Na}^+$  environment. Therefore, cells lyse in NaCl concentrations below 15%.

- The cell walls of *Methanobolus*, *Sulfolobus*, *Thermoproteus*, *Desulfurococcus* and *Pyrodictium* are made up of glycoproteins (Glycoprotein S-layer)
- *Methanomicrobium*, *Methanococcus*, *Methanogenium* and *Thermococcus* cell walls are exclusively made up of protein subunits (protein S-layer).
- *Methanospirillum* cell wall consists of a protein sheath.

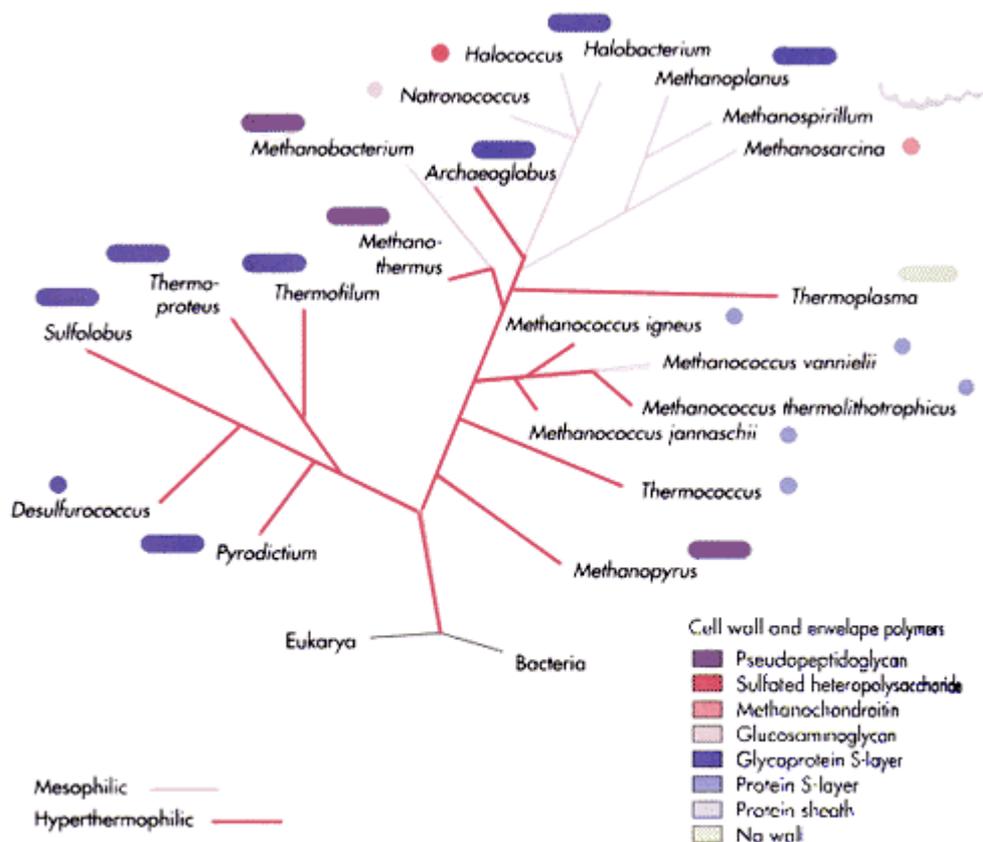


Fig. 18-4 Diversity and Evolutionary Patterns of Archaeal Cell Wall Polymers. Diverse cell wall polymers appear to have evolved independently along various lines of archaeal evolutionary descent.

### Lipids and Cell Membranes

The chemistry of lipids is very different to that of members of domains *Bacteria* and *Archaea* and is perhaps the most distinctive feature of archaeal cells.

Archael glycerol molecules may be linked:

- to a phosphate group (similar to bacteria & eucaryotes) and / or
- to a sulfate and carbohydrates (unlike bacteria & eucaryotes) & therefore phospholipids are not regarded as universal structural lipids.

Archael lipids are hydrocarbons (isoprenoid hydrocarbons) not fatty acids, are branched (straight chain in bacteria & eucaryotes) and linked to glycerol by ether bonds (ester linked in bacterial & eucaryotes).

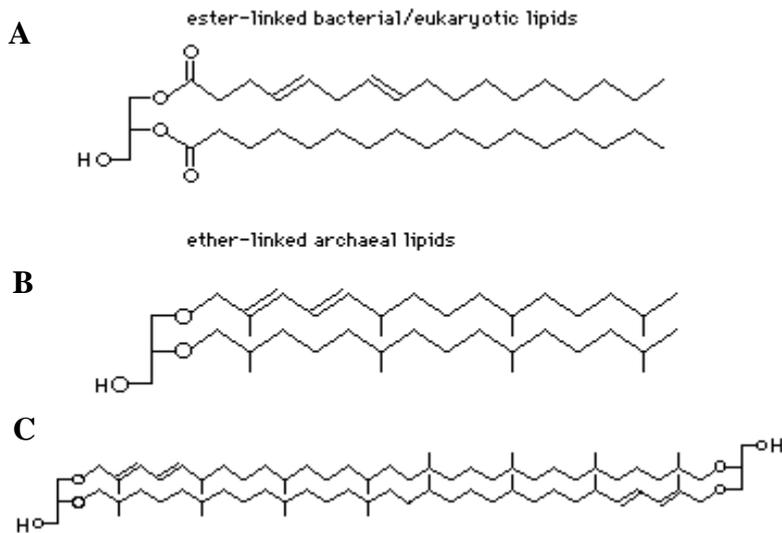
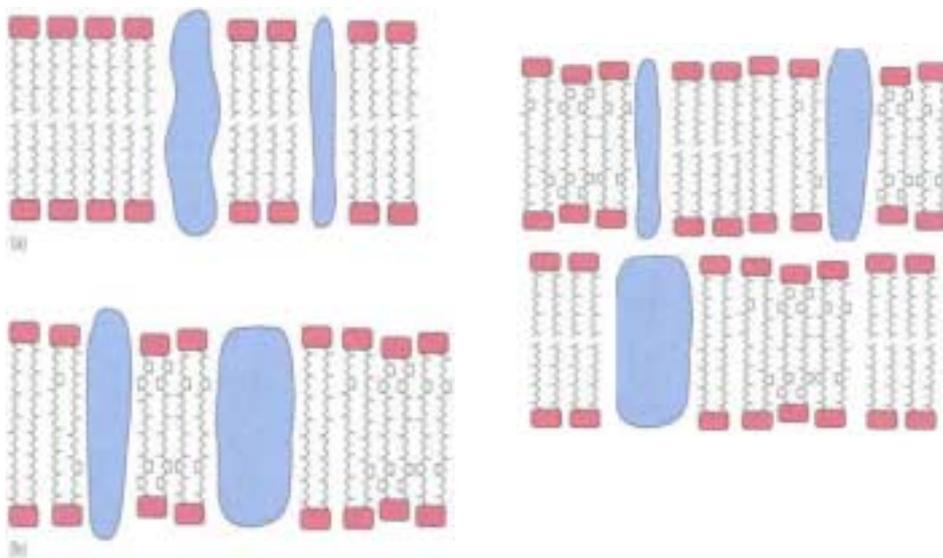


Figure: Bacterial lipids are made of phospholipids - a phosphate group joined to 2 fatty acids by glycerol (glycerol diester) (A) but archaeal lipids are composed of phosphate, sulfate or carbohydrate joined to branched  $C_{20}$  and / or  $C_{40}$  hydrocarbon chains by glycerol diethers (B and C respectively).

Archaeal lipids are diverse in structure:

- Glycerol diether (Glycerol +  $C_{20}$  hydrocarbons)- Bilayered membrane
  - Glycerol tetraether (Glycerol +  $C_{40}$  hydrocarbons)- Monolayered membrane
  - Mixture of di- & tetra- Mono /Bi layered membrane
- Cyclic tetraethers (Glycerol +  $> C_{40}$ )- maintain the 4-5nm membrane thickness



FigureXX. Membranes of archaeae possess integral proteins and a bilayer composed of  $C_{20}$  diethers (a) a rigid monolayer composed of  $C_{40}$  diethers (b) a mono / bilayer consisting of a mixture of  $C_{20}$  and  $C_{40}$  diethers (c).

Diversity of membrane rigidity requirements is related to the diverse habitats that *Archaea* live in and the changes are necessary to maintain membrane fluidity and stability:

- *Sulfolobus* (90°C, pH 2)- branched chain  $C_{40}$  hydrocarbons. Branched chains increase membrane fluidity (unbranched & saturated fatty acids limit sliding of

fatty acid molecules past one another) and is required for growth at high temperatures (upto 110oC, hyperthermophiles)

- *Halobacterium* cope with saturated salts.
- *Thermoplasma* cope with high temperature without cell walls

## 4. Functional Attributes

### Habitat and ecology

They are found in extreme environments and include anaerobic, saline to hypersaline, high temperature and cold temperature habitats (constitutes 34% of microbial biomass in the Antarctic surface waters). Some are symbionts in animal digestive tracts.

### Nutrition, Physiology and Metabolism

Some are aerobes, some strict anaerobes and some facultative anaerobes. Some are chemolithoautotrophs, others chemoorganotrophs and some can grow by changing from one to the other nutritional mode. The growth temperature varies between mesophilic to the hyperthermophilic. The pH growth range is between very acidic (< 0.5) to slightly over neutral. Some are obligate halophiles yet others are halotolerant.

Chemoorganotrophs: Use organic substrates as energy source for growth. Catabolism of glucose occurs via modified Entner-Doudoroff (E-D) pathway. Oxidation of acetate to CO<sub>2</sub> proceeds via the citric acid cycle or by the Acetyl-CoA pathway. Amino acid biosynthesis pathways unknown. Electron transport chains including cytochromes of type a, b and exist in some *Archeae*. Consequently, electrons from organic electron donors enter the electron transport chain leading to the reduction of O<sub>2</sub>, S<sup>0</sup> with concurrent establishment of PMF to drive ATP synthesis through membrane-bound ATPases.

Autotrophy: Widespread and occurs by several different means. Methanogens use Acetyl-CoA or some modification to fix CO<sub>2</sub>. In some *Archeae* (e.g. *Thermoproteus* and *Sulfolobus*), CO<sub>2</sub> fixation occurs via the reverse citric acid cycle and is in common with the Green sulfur bacteria and of the genus *Chlorobium* and *Aquifex* of domain *Bacteria* or via the Calvin Cycle the most common autotrophic pathway in *Bacteria* and *Eucarya*. Very thermostable RubisCO enzyme which catalyses the first step in the Calvin Cycle is found the methanogen *Methanococcus jannaschii* and a *Pyrococcus* species.

In summary, catabolic and anabolic processes are somewhat similar for Bacteria and Archaea but methanogens are very different. More will be said when archael kingdoms as well as their genomes are discussed.

## 5. Molecular Attributes:

### Genomes

- Archaeal chromosome is a single, circular DNA molecule & extrachromosomal elements (eg plasmids) are found in *Archeae*. However, they are smaller than bacterial chromosomes.

- The genomes can vary in the G+C content of their DNA between 21 to 68 mol% indicating a marked genotypic diversity.
- The typical genome of bacterial genome is a single circular DNA molecule. Plasmids are also common. *Streptomyces* spp. and *Borrelia* spp have linear chromosomes, and *Rhodobacter sphaeroides* has two chromosomes.
- A typical eucaryote genome consists of multiple linear DNA molecules (covered in histones and organized in nucleosomes, chloroplast DNA, and mitochondrial DNA. (note the eukaryote dinoflagellate algae has no histones associated with its DNA)
- Archaeal DNA binding proteins (aka Histone-like proteins similar to Eucarya): *Methanosarcinaceae* MC1 and *Methanobacteriales* Hmf share amino acid homology to eucaryal histone proteins
- Organization of DNA in chromatin-like structure  
histone + Eucarya DNA = negative supercoiling & nucleosome  
Hmf + Archaea DNA = positive supercoiling  
HTa: *Thermoplasma*  
HTa-like: *Sulfolobus*
- Genomic resistance to thermal denaturation & genomic structural integrity in extreme halophiles is related to high intracellular salt concentrations (solute)

### **Gene Organisation in genomes**

- Functionally related genes are often organised in operon like structures though the primary sequences of archaeal proteins more often resemble eukaryotic homologues rather than bacterial ones
- Introns have been found in archaeal 23S and 16S rRNA and tRNA genes
- Varying arrangements of genes can be seen in *Archaea*
- Methyl coenzyme M reductase, RNAP and bacteriopsin genes are good Examples

### **DNA Modyfing Enzymes**

- Archaeal DNA polymerases involved in DNA replication have been identified. The primary protein sequences of these enzymes resemble the DNA polymerases from eukaryotes, eukaryal viruses and *E.coli*. Some posses 3'-5' exonuclease (or proofreading) activity
- A *Halobacterium halobium* DNA polymerase / primase has been identified with reverse transcriptase activity
- Topoisomerases, gyrase and restriction endonucleases have also been identified in *Archaea*

## Transcription

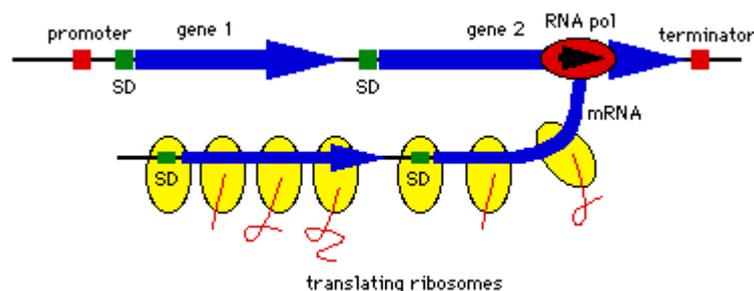
- Bacteria have only one type of RNA polymerase, Eukaryotes have three types namely, RNA polymerases POL I, II and III. and *Archaea* also has one type but it is similar to the eukaryote RNA polymerase POL II. Comparative sequence homology of genes encoding subunits suggests that the RNA polymerase is more closely related to eukaryal polymerases than the bacterial counterpart.
- Archaeal RNA polymerases are complex, consisting of up to 14 subunits (c.f. 5 in *E. coli*)
- Unlike, *E. coli* RNA polymerase, archaeal RNA polymerases are unable to initiate transcription *in vitro*. This is also seen in eukaryotes where general transcription factors are required for initiation
- Archaeal promoters have an A-T rich sequence at -32 to -25 bp upstream of the transcriptional start: the consensus sequence resembles a eukaryotic TATA box.

## Translation

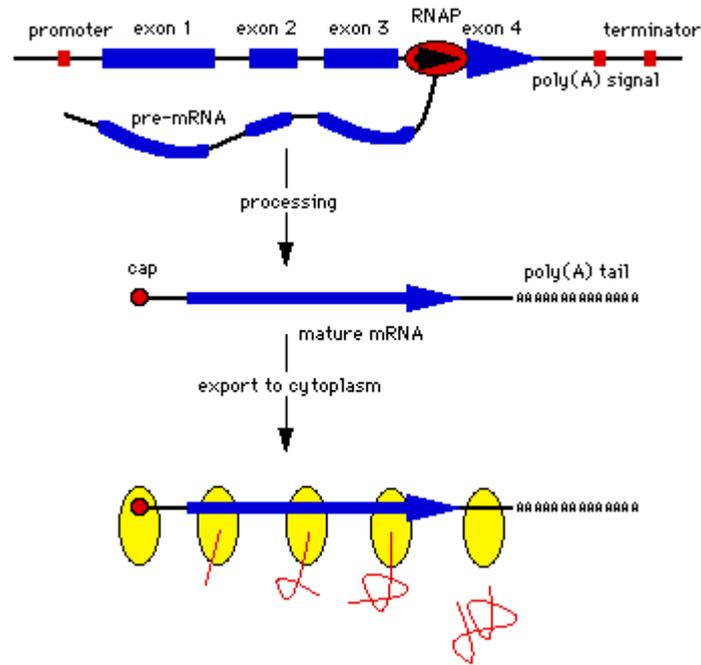
- Translation signals resemble those found in bacteria (i.e. there are short regions of complementarity between the 5' end of the mRNA and the 3' end of the 16S rRNA)
- Complementary 5' mRNA to 3' 16S rRNA similar to bacteria
- Lack of formylmethionine

## Detailed understanding of Transcription and translation (You may also have covered this in some other subjects)

The translational machinery in *Archaea* is generally like those of *Bacteria* with 70S (bacterial-sized) ribosomes. Genes are arranged in co-transcribed clusters called operons. Ribosomes recognize translational start sites and bind to the mRNAs directly at 'Shine-Dalgarno' (SD) sequences just like *Bacteria*. Also like in *Bacteria*, transcription and translation are linked - that is, they occur simultaneously, and failure of an mRNA to be translated causes the RNA polymerase to abort transcription.



Eukaryal genes generally are transcribed separately rather than in clusters. In plants and animals, most genes are segmented into 'exons' separated by 'introns' - introns are spliced out of the mRNAs after transcription. In addition, the 5' end of the mRNA is capped with a modified nucleoside, and the 3' end is cleaved and polyadenylated. Ribosomes generally bind the 5' 'cap' of the mRNA & scan downstream for AUG codons to start translation. Because transcription occurs in the nucleoplasm and translation occurs in the cytoplasm, the mRNA has to be exported through the nuclear membrane after transcription before it can be translated, so transcription and translation are not linked.



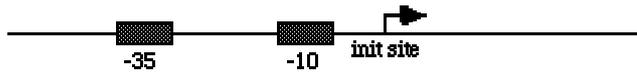
However, in many ways translation in Archaea is like it is in Eukarya. Translation is initiated with methionine (like Eukarya), not formyl-methionine (Bacteria). Translation is inhibited by diphtheria toxin, as are eukaryal ribosomes, but is not inhibited by most bacterial-translation-inhibiting antibiotics. Chimeric archaeal/eukarya ribosomes are functional - bacterial/archaeal & bacteria/eukarya are not functional.

#### Example - RNA polymerase

A typical example of how *Archaea* resemble eukaryotes in some ways and *Bacteria* in others is RNA polymerase:

*Bacteria* contain a single RNA polymerase that transcribes all genes in the cell. The holoenzyme contains 5 polypeptides - 2 copies of alpha, and one each beta, beta prime, and sigma. The sigma subunit of RNA polymerase provides promoter specificity by directing binding to signal sequences 10 and 35 base-pairs upstream of the transcription initiation site. Different sigmas regulate developmental pathways, e.g. sporulation.

Bacterial promoters:



*Eucarya* contain 3 nuclear RNA polymerases. Each is specialized for transcription of specific gene types. These enzymes contain many more subunits than do bacterial RNA polymerase: 3 or 4 large, ~9-14 small. Promoter recognition is not a function of the RNA polymerase, but is provided by transcription factors (TF) that bind directly to the promoter (not RNA polymerase), and the DNA:TF complex is recognized by the RNA polymerase.

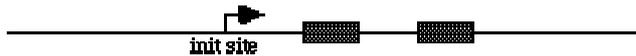
RNA polymerase I transcribes ribosomal RNA genes, and is confined to the nucleolus. Promoter sequences surround the initiation site, defining a binding site for a transcription factor complex.



RNA polymerase II transcribes primarily mRNAs. Promoter sequences vary widely, but generally have a -35 TATA-box that is a binding site for a general RNA polymerase II transcription factor TFIID.



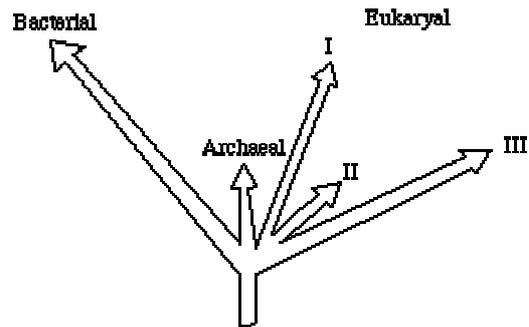
RNA polymerase III transcribes primarily 5S rRNAs and tRNAs. The promoter elements are downstream of the transcription initiation site, and are binding sites for TFIIC or B and TFIID.



*Archaea*, like *Bacteria*, have a single RNA polymerase that transcribes all genes. However, archaeal RNA polymerases are like those of eukaryotes in that they contain 3 or 4 large and many small subunits. Archaeal RNA polymerases are similar in sequence and in antigenicity to eukaryal RNA polymerase II. Promoters in *Archaea* are a -30 TATA-like sequence that is a binding site for a transcription factor (TFB), not sigma-like subunits.



Comparison of RNA polymerase gene sequences:



Archaeal and eukaryal RNA polymerase II are primitive RNA polymerases, whereas bacterial RNA polymerase is the most highly evolved. This is typical of bacterial vs eukaryal vs archaeal evolution; Bacteria keep the mechanisms basically unchanged but hone the parts to perfection, whereas eukaryotes duplicate genes and specialize, each for separate functions. Archaea, on the other hand, appear to have remained unchanged.

#### Archaea as missing links between eukaryotes and *Bacteria*

In many ways, *Archaea* are a 'missing link' between *Bacteria* & *Eukarya*. Archaeal features generally resemble *Eucarya* & *Bacteria* features more than they do each other, a real sign of their primitiveness. Having a third evolutionary group, especially a primitive one, allows the identification of primitive traits, i.e. traits of the Last Common Ancestor (LAC). With only two groups, it's impossible to tell which version of a trait, if either, is primitive. It used to be generally assumed that the bacterial version was primitive, but this is rarely the case.

It has also become clear that the *Archaea* share a common ancestry with the *Eucarya* to the exclusion of the *Bacteria* - in other words, the *Archaea* and *Eucarya* are 'sibling' groups, whereas the *Bacteria* are 'cousins' to the *Archaea* and *Eucarya*. *Archaea*, then, are primitive relatives of *Eucarya*, and as such are ideal organisms to shed light on the complexity of eukaryotic organisms - for example the RNA polymerase we just discussed.

#### Cloning and expression

- Cloning usually follows protein-purification using standard molecular biology techniques
- Expression in heterologous hosts may be complicated by the altered environment in which expression is occurring and differences in translation and post-translational mechanisms

#### Phage and plasmids

A few reports on the occurrence of plasmids and phage have been reported in members of the *Archaea*. However, studies in these areas are progressing slowly due to the very specialised techniques required for culturing especially for the hyperthermophiles.

## 6. Kingdom *Crenarchaeota*:

This kingdom is divided into 2 sections in order to make the discussion easier. The first section deals with the thermophilic and hyperthermophilic crenarchaeotes and the second section deals with the uncultured cold dwelling crenarchaeotes.

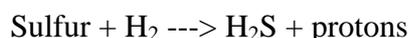
### Section I: Thermophilic & hyperthermophilic *Crenarchaeota*

Not all crenarchaeotes have been cultured and those that have been cultured are all thermophilic or extremely thermophilic (aka hyperthermophiles), with optimal growth temperatures above 80°C. Some of the crenarchaeotes are the most thermophilic organisms known. Most are also acidophilic and autotrophic. This phenotype is also shared by the deepest branches of kingdom *Euryarchaea* and domain *Bacteria*, and therefore it has been suggested that these traits are probably the primitive phenotype of the Last Common Ancestor (LCA).

Most crenarchaeotes, not surprisingly, have been isolated from volcanic geothermal environments rich in elemental sulfur (called solfataras). Volcanic activity is found on the land mass (terrestrial) or on the ocean floor formed along tectonic plates (hydrothermal vents) (Figure 2, Map of volcanic areas of the world, also the tectonic plates). ALSO FIGURE FROM MY PHOTO COLLECTION).

Sulfur metabolism as a key trait: Crenarchaeotes oxidise and / or reduce sulfur by one of 3 biochemical processes. In most cases sulfur compounds such as thiosulfate are also usable in place of sulfur.

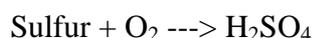
1. Sulfur reduction: These organisms are autotrophic anaerobes that fix carbon from CO<sub>2</sub>. Hydrogen is the electron donor for electron transport and elemental sulfur (or sulfur compounds such as thiosulfate) is the terminal electron acceptor.



2. Sulfur respiration: These organisms are heterotrophic anaerobes. Both carbon and energy are extracted from organic compounds. Organics are the electron donor for electron transport and sulfur (or sulfur compounds) is the terminal electron acceptor. This process is much like 'regular' respiration, except that sulfur compounds take the place of O<sub>2</sub>



3. Sulfur oxidation: These organisms can usually grow heterotrophically, getting fixed carbon from low concentrations of organics in the medium. Most can also be grown autotrophically, fixing carbon from CO<sub>2</sub> via the reverse TCA cycle. All are aerobes, of course, since it is the terminal electron acceptor (sulfur is the electron donor) for electron transport.



Taxonomy: At the time of writing, at least 12 genera belonging to the two orders, namely, *Sulfolobales* and *Thermoproteales* have been validly described, and a third, order "*Igneococcales*", has been proposed (Figure 3).

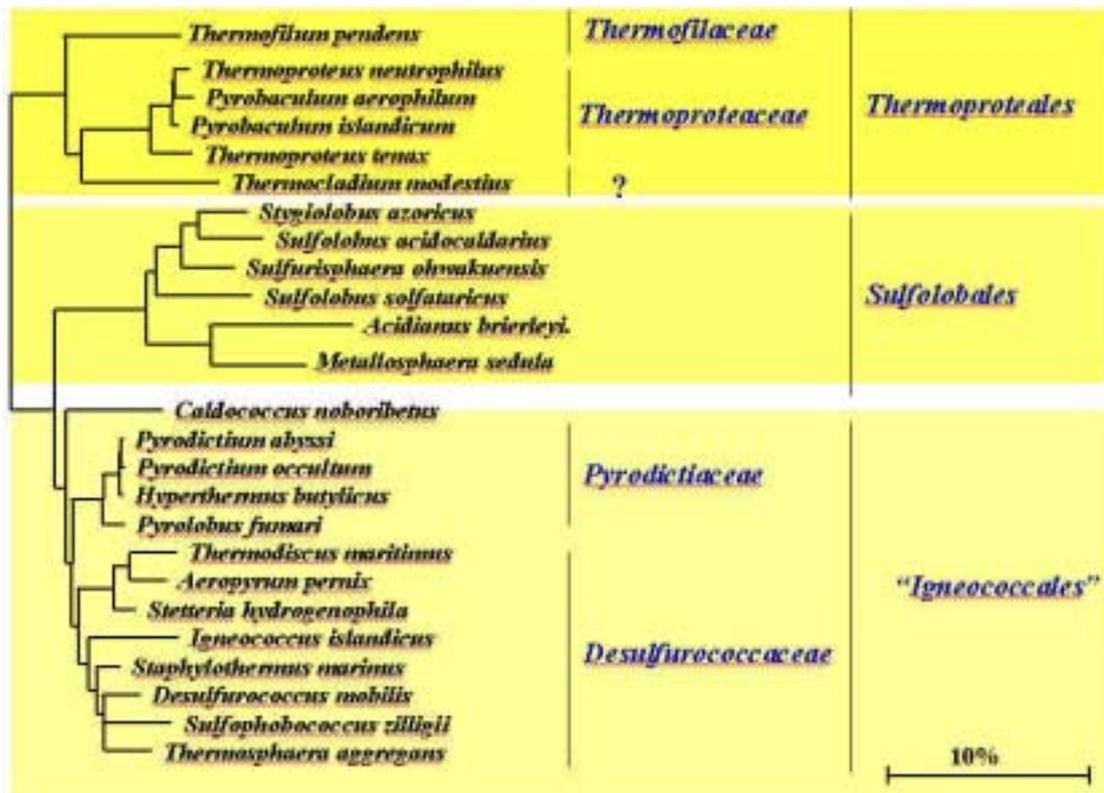


Figure 2. Members of kingdom *Crenarchaeota* are thermophilic to hyperthermophilic microbes and are divided into 3 orders.

The most studied genera include *Thermoproteus* and *Sulfolobus*. The characteristics of some of the members from each of the order are described as examples below.

Order "Igneococcales":

*Pyrodictium*: Most isolates have been cultured from ocean floor volcanoes (hydrothermal vents). All are strict anaerobes. Some species are organotrophs and some are lithotrophs growing on H<sub>2</sub> and S<sup>0</sup>. *P. occultum* is the most thermophilic (hyperthermophile) organisms known to date. The optimum temperature for growth is 105 °C (requires a bar of pressure in culture tubes to avoid medium from boiling) with a minimum growth temperature of 82 °C and a maximum growth temperature of 115 °C.

*Desulfurococcus*: These are cocci shaped and have been isolated from terrestrial and marine volcanic environments. Some members are motile and some are non-motile. Strict anaerobes that reduce sulfur and are also able to respire (See figure XX below)

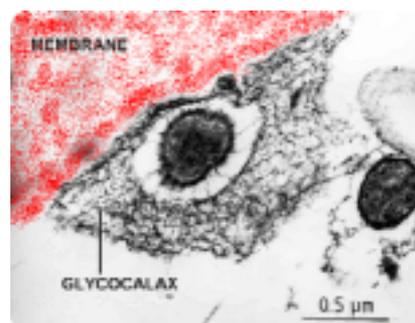


Figure XX. A *Desulfurococcus* species shows extensive glycocalx used in attachment. In this case polycarbonate membrane was immersed for 5 hours in a New Zealand hot spring with a temperature of 90 °C and the sample processed for electron microscopy. The chemical composition of glycocalx is unknown.

Order Sulfolobales:

*Sulfolobus*: Several species have been described and all belong to the order *Sulfolobales*. The cells gram-negative aerobes and have an irregular spherical lobed shape (see Figure XX below). The temperature optimum is between 70 to 80 °C with a pH optimum of between 2 to 3. Some species are able to grow at pH 0 (equal to 0.5M H<sub>2</sub>SO<sub>4</sub>) *Sulfolobus* is sometimes referred to as a thermoacidophile. The cell walls lack peptidoglycan but contain lipoprotein and carbohydrates. Oxygen is the normal electron acceptor but ferric iron can also be used. Some strains are microaerophilic



Figure XX. *Sulfolobus* is not spherical but is lobe-shaped and can be isolated from acidic hot springs.

Some strains grow lithotrophically by oxidising S<sup>0</sup> producing sulfuric acid whereas some other strains grow on sugars and organic acids (eg glutamate) as a carbon and energy source. These organisms, like most acidophiles, are oligotrophic ie high concentrations of organics, especially organic acids, are toxic.

The reason for this is that organic acids are protonated in the external growth medium (pH < 3.5), & hence remain uncharged, enabling them to diffuse freely into the cytoplasm through the lipid membrane. The internal cytoplasmic pH of the cells is pH >5.5 and under this condition, the organic acid ionizes, releasing H<sup>+</sup>. Cytoplasm acidification leads to decomposition of the proton gradient force (i.e. high concentrations of organic acids act as uncouplers). See Fig XX below.

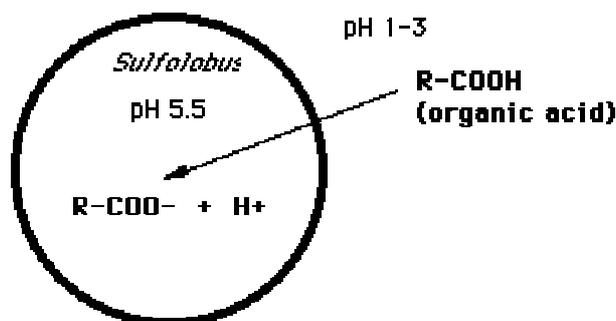


Figure XX: High concentrations of organic acids act as decouplers as they are protonated in the external acidic environment of the growth medium (pH < 3) and are

ionized releasing protons inside the cytoplasm (pH > 5.5). The net result of this is a disruption to the proton gradient leading to death of the cells.

*Sulfolobus acidocaldarius* can be easily isolated from volcanic acid environments (soil and springs). The genome of *Sulfolobus acidocaldarius* has been partially sequenced (REFERENCE TO THE WEB SITE).

Order *Thermoproteales*:

*Thermoproteus*: These are members of the order *Thermoproteales*. Cells are long and thin and can be identified by the presence of occasional branching of the cells and presence of golf ball-like terminal structures. Cell walls are made of glycoprotein subunits. Strict anaerobes which grow between 70 to 90 °C (optimum 85 °C). Lithotrophic growth occurs by S<sup>0</sup> reduction (with H<sub>2</sub>) and CO<sub>2</sub> (or CO) as a sole carbon source. In addition, they can also grow organotrophically by oxidising glucose, amino acids, alcohols and organic acids by anaerobic S<sup>0</sup> respiration. Isolated from neutral volcanic terrestrial and marine hot springs rich in sulfur.

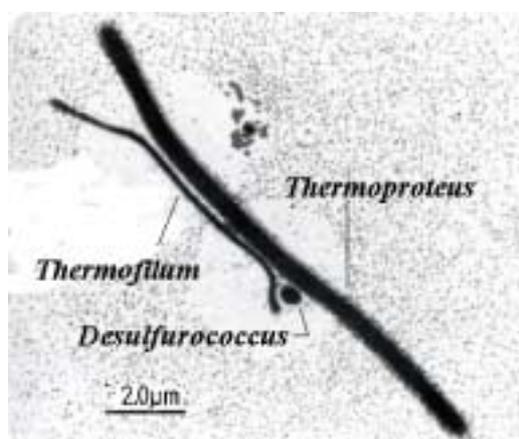


Figure XX. A variety of thermophilic *Archaeae* thrive as mixed populations in terrestrial hot springs and can be visualised using electron microscopy techniques. The crenarchaeotes *Thermofilum* are thin filament, *Thermoproteus* are rod to filaments with a larger diameter whereas *Desulfurococcus* are cocci.

*Thermofilum*: The genus *Thermofilum* is member of the order *Thermoproteales* and two species have so far been described. They require as yet unidentified factors from *Thermoproteus* (by either co-culturing with *Thermoproteus* or with *Thermoproteus* extracts) without which they are unable to grow. The cells are extremely filamentous and very thin (<100 x 0.15-0.3 µm) and can be easily distinguished from *Thermoproteus* based on this trait.

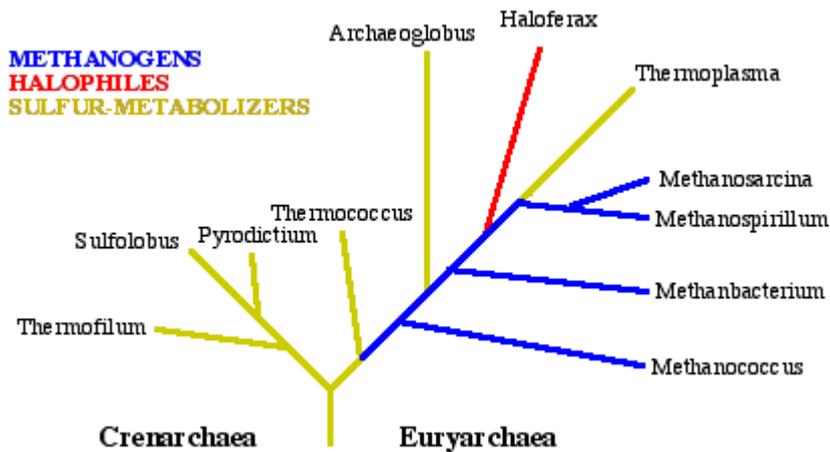
## Section II. The cold dwelling Crenarchaeota

Cold dwelling crenarchaeotes have not as yet been cultured but their presence has been demonstrated by fluorescent phylogenetic staining in marine waters world wide and have also been found in frigid marine waters of the Antarctic. High crenarchaeote concentrations (around 10,000 cells / ml) have been observed in these cold and nutrient deficient waters.

The metabolism and physiology of these organisms remains a mystery, though the cells which were collected after membrane filtration were found to contain ether linked lipids of the diphytanyl tetraether type which have so far been known to occur only in hyperthermophiles. The diphytanyl tetraethers were thought to be the molecular secret behind hyperthermophiles surviving high temperatures but their presence in cold dwellers, leaves this hypothesis to refinement.

### **7. Kingdom Euryarchaeota:**

**Phenotypic traits and evolution:** Most members of kingdom *Euryarchaea* are predominantly methanogens, but there are 2 other phenotypes found in this group - sulfur-metabolizing thermophiles and extreme halophiles. One group of sulfur-metabolizing thermophiles, *Thermococcus* & relatives, seem to have retained that phenotype from the common ancestry of *Euryarchaea* & *Crenarchaea* (and *Bacteria*, for that matter), whereas the other sulfur-metabolizers and halophiles evolved these phenotypes from methanogenic ancestry (see figure below).



**Diversity and Taxonomy:** This is an extremely diverse kingdom and is characterised by 5 distinct orders (SEE FIGURE BELOW):

- (a) *Halobacteriales* (the extreme halophiles)
- (b) Methanogens which are divided into 5 orders
- (c) *Thermoplasmatales*
- (d) *Archaeoglobales* (sulfate reducers) and
- (e) *Thermococcales*

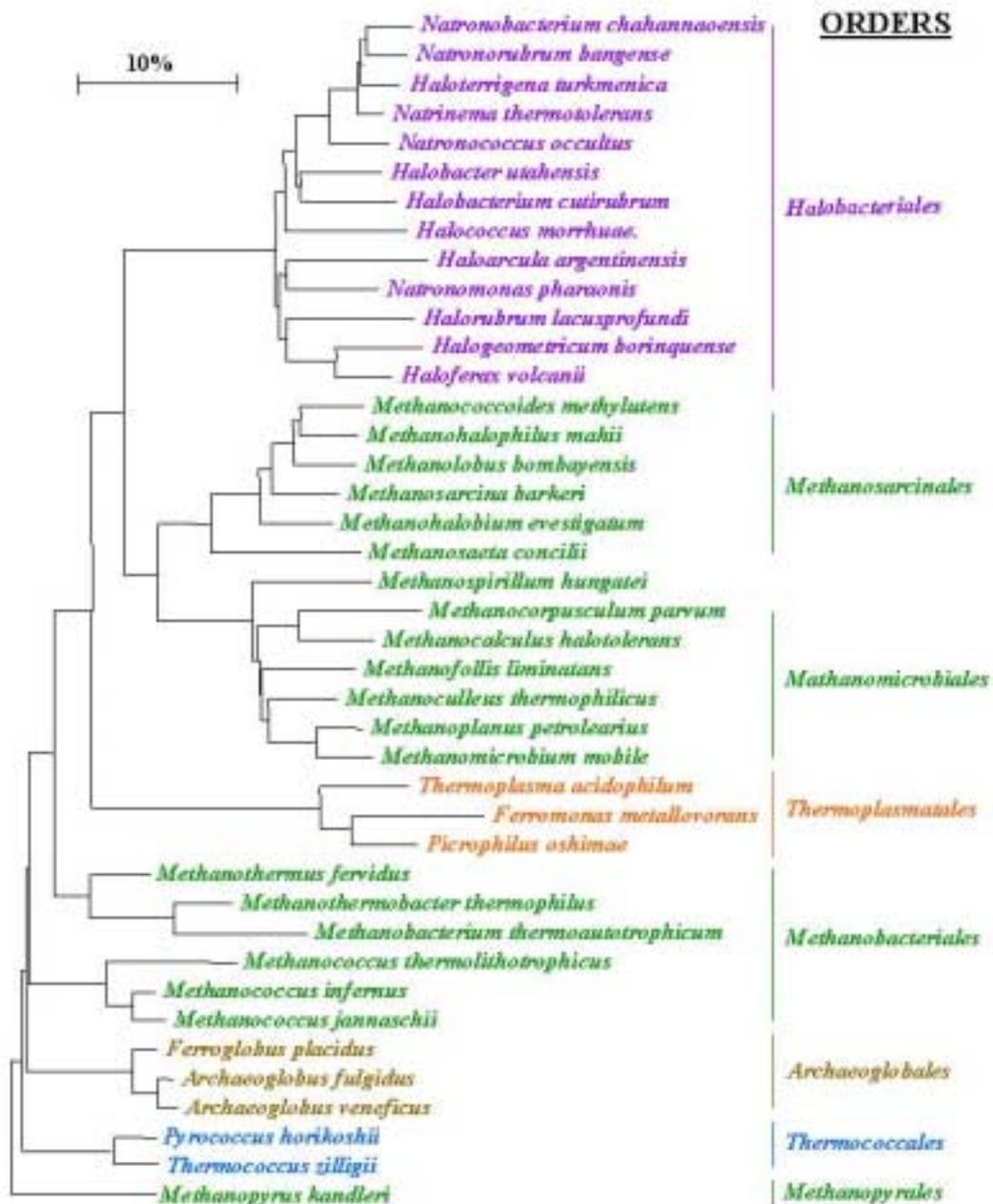


Figure 3. Dendrogram showing the five distinct groups consisting of the orders *Halobacteriales* (the halophiles) *Thermoplasmatales* (lacking cell walls), *Archaeoglobales* (sulfate reducers), *Thermococcales* and methanogens (which are divided into 5 orders) within the Kingdom *Euryarchaeota*.

### **1. Order *Halobacteriales*, the extreme halophiles**

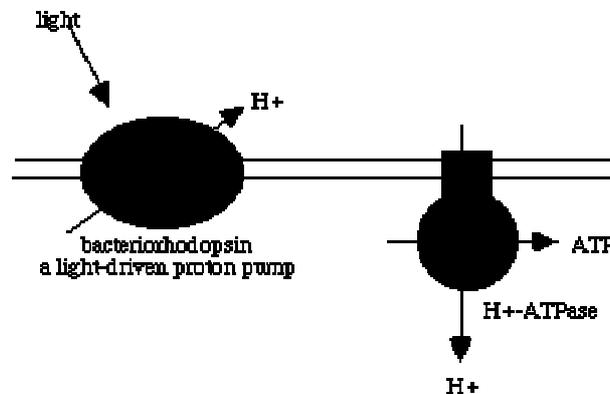
The extremely halophilic *Archaea* require at least 2M NaCl or equivalent ionic strength for growth - most grow in saturated or near-saturated brines. They are the primary, or only, inhabitants of salt lakes. Red pigments make it obvious when large numbers of these organisms are present - blooms often occurs after a rain carries

organic material into a salt lake, and the Red Sea gets its name from such blooms. So does the 'Red Herring', from foul-smelling but hound-decoying salted fish being spoiled by *Halobacterium*. They are common in hypersaline seas, salt evaporation pools, salted meats, dry soil, salt marshes, etc. They are also found in subterranean salt deposits, where micropockets of saturated water 'diffuse' around in the otherwise solid salt.

Other halophilic organisms (e.g. fungi, brine shrimp) have normal cytoplasmic salt concentrations, expending energy to continuously pump salt out of the cell and water into the cell, and contain organic osmolytes like glycerol or sugars. Halophilic Archaea grow at much higher salt concentrations, and don't fight back at all - the internal salt concentrations are as high as they are outside! For this reason, there is little or no net osmotic pressure on the cell wall, and some organisms take advantage of this by adopting high surface-area shapes that are not possible for organisms in 'normal' ionic strength. One example is *Haloarcula*, which comes in squares and triangles with straight edges, sharp corners, and is very flat. Other halophiles are rods or cocci.

Halophiles are mesophilic facultative aerobes. Aerobically, they grow heterotrophically, via respiration, using O<sub>2</sub> as the terminal electron acceptor. Anaerobically, they grow photochemotrophically - they get energy (ATP) from light, but still need organics for carbon.

They do not contain the usual photosystems or electron transport chain for gathering energy from light. Phototrophy is driven by a single protein, bacteriorhodopsin, that is a light-driven proton pump.



This proton pump generates a proton gradient used to make ATP via ATPase, just like in other organisms. It's not nearly as efficient as the bacterial photosystems, but light is rarely limiting for growth in the desert salt lakes where they predominate.

Some halophiles grow at high pH (up to pH10-10.5) i.e. *Natronobacterium* in soda lakes. This is a problem for them (or at least for us, trying to understand how they get away with it). At that pH, any protons pumped to the outside, by electron transport or rhodopsin, are gone forever. Even though the resulting electric potential is still there, it can't be harvested by an ATPase unless it can get protons from the outside. How they get around this is not known, and it is probably this issue that limits the upper pH range or life.

The salient features of the extreme halophiles is shown in the table below.

**Table 18-10 Differentiation of the Genera of Extremely Halophilic Archaea**

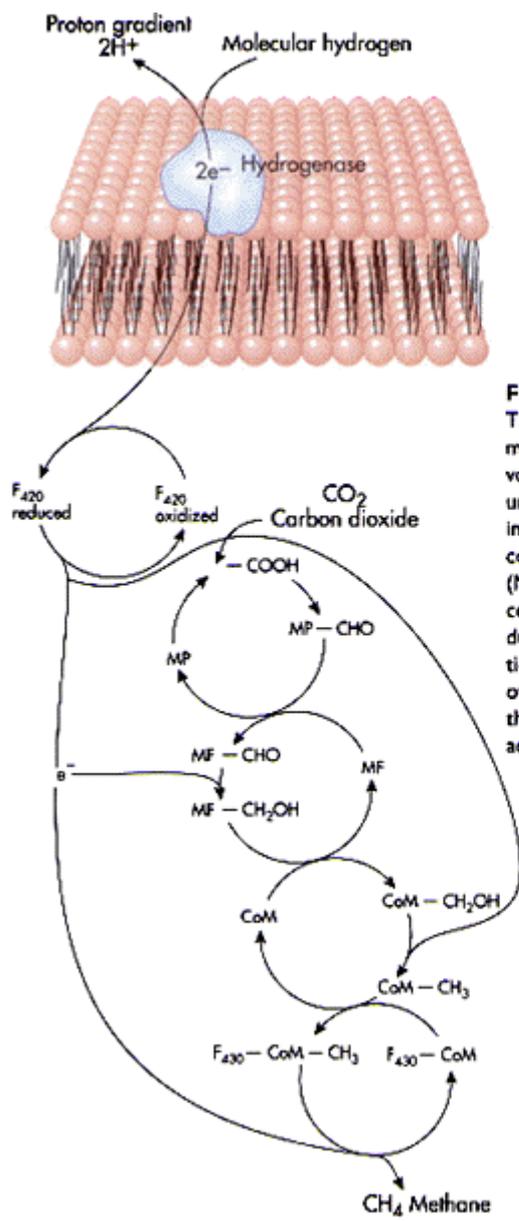
| Characteristic                                  | <i>Haloarcula</i>  | <i>Halobacterium</i>  | <i>Halococcus</i>   | <i>Haloferax</i>  | <i>Natronobacterium</i>   | <i>Natronococcus</i>  |
|---|--|---|---|---|---|---|
| Cell shape                                      | Irregular rods, triangles, rectangles  | Irregular rods  | Cocci   | Irregular rods, disks   | Irregular rods  | Cocci   |
| pH range for growth                             | 5.0-8.0  | 5.0-8.0   | 5.0-8.0   | 5.0-8.0   | 8.5-11.0  | 8.5-11.0  |
| Carbohydrates used as carbon and energy sources | +  | -   | +/-   | +   | +/-   | -   |
| Polar lipids                                    | Contain phosphatidyl glycerol phosphate; glycolipids occur as glucosyl mannosyl glucosyl glycolipid; polar lipids characterized by C <sub>23</sub> , C <sub>26</sub> glycerol ether core | Contain phosphatidyl glycerol phosphate; glycolipids occur as sulfated galactosyl mannosyl glucosyl glycolipid and sulfated digalactosyl mannosyl glucosyl; polar lipids characterized by C <sub>20</sub> , C <sub>20</sub> glycerol ether core | Lack phosphatidyl glycerol phosphate; glycolipids occur as sulfated mannosyl glucosyl glycolipid; polar lipids characterized by C <sub>20</sub> , C <sub>20</sub> and C <sub>20</sub> , C <sub>23</sub> glycerol ether core | Lack phosphatidyl glycerol phosphate; glycolipids occur as sulfated mannosyl glucosyl glycolipid; polar lipids characterized by C <sub>20</sub> , C <sub>20</sub> glycerol ether core | Lack phosphatidyl glycerol phosphate; lack glycolipids; polar lipids characterized by C <sub>20</sub> , C <sub>20</sub> and C <sub>20</sub> , C <sub>23</sub> glycerol ether core | Lack phosphatidyl glycerol phosphate; lack glycolipids; polar lipids characterized by C <sub>19</sub> , C <sub>20</sub> and C <sub>20</sub> , C <sub>27</sub> glycerol ether core |

## 2. Methanogens and it's 5 orders:

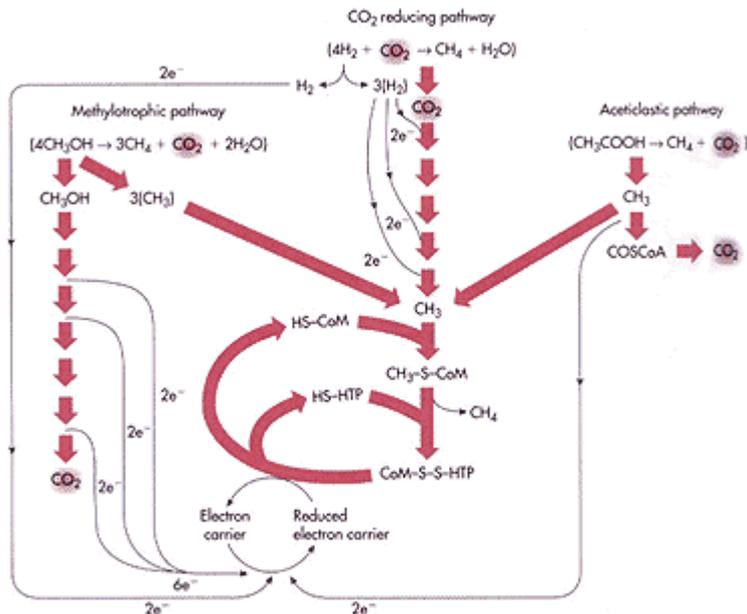
Methanogens make methane (CH<sub>4</sub>) via a unique metabolic pathway with unique enzymes & cofactors, including F<sub>420</sub>, that makes the cells fluorescent.

All methanogens are obligate lithotrophs, that is, they can make energy only by methanogenesis. All can make methane from CO<sub>2</sub> + H<sub>2</sub>, some can also use other one carbon compounds (e.g. CO, formic acid), and a very few can use acetate and methylamines. At least one ATP is generated with the reduction of CO<sub>2</sub> to CH<sub>4</sub> with H<sub>2</sub>.

One carbon compounds (formate, CO<sub>2</sub>, CO) are reduced with hydrogen and attached to methanofuran (MF) in the form of a formyl group. This formyl is transferred to methanopterin (MP) and sequentially reduced through methenyl and methylene to methyl, which is then transferred to coenzyme M, and finally reduced to free methane. Organisms that grow on acetate transfer the methyl group of acetate directly to MP and the carboxy group is released as CO<sub>2</sub>. Organisms growing on methanol transfer the methyl group (indirectly) to CoM.



**Fig. 18-14 Methanogenesis—CO<sub>2</sub> Reduction Pathway.** The conversion of CO<sub>2</sub> to CH<sub>4</sub> (methane) is carried out by methanogenic archaea. This is a strictly anaerobic pathway involving the flow of electrons from a hydrogen donor. Several unique electron carriers are involved in the transfer of electrons in this pathway, including factor 420 (F<sub>420</sub>), factor 430 (F<sub>430</sub>), coenzyme M (CoM), methanopterin (MP), and methanofuran (MF). The oxidation of hydrogen, which occurs outside of the cell, produces hydrogen ions and supplies electrons for the reduction of F<sub>420</sub>, which occurs inside the cell. Because the reduction of F<sub>420</sub> inside the cell consumes protons, and the oxidation of hydrogen produces protons outside the cell, the net result is the establishment of a proton gradient (protonmotive force) across the membrane.



**Fig. 18-13 Methanogenesis.** There are three different metabolic strategies employed by diverse methanogenic archaea. Most methanogens employ the CO<sub>2</sub>-reducing pathway starting with carbon dioxide or formate. Others use a methylotrophic pathway starting with methanol or methylamines. A few methanogens use an aceticlastic pathway starting with acetate.

The enzymes in the methanogenic pathway are very oxygen sensitive, so methanogens are extreme anaerobes.

Methanogens are common organisms, found in all types of anaerobic environments, and are certainly the most prevalent *Archaea* in the 'moderate' world.

Sediments and soils - swamp gas is methane which, because of its low ignition temperature and low threshold concentration, is readily ignited and glows very faintly white in 'will-o-the-wisps' visible at night in swamps. Methanogens are crucial components of the microbial populations of the 'rhizosphere', the plant root environment. Animal guts - especially wood-eating insects and ruminants. African termite mounds are scrupulously aerated by the insects not for oxygen, but to keep methane concentrations low. Termite mounds struck by lightning after a light rain (blocking aeration) can explode spectacularly! Cows also produce large amounts of methane.

Wastewater and landfills - The whole wastewater process works because organics in the wastewater are converted first to biomass (in the early stages of treatment), then digested anaerobically to H<sub>2</sub>, CO<sub>2</sub>, and acetate which in turn is converted by methanogenesis into methane, which floats away. Landfills have to be carefully vented to prevent very messy explosions! Several houses near older unvented landfills have exploded because of the buildup of methane that seeped through the ground into their basements.

Oil deposits - natural gas is methane, produced not geochemically but by methanogens living in the subterranean oil deposit.

Methanogens form a variety of symbioses with plants, animals and protists, but despite these close associations there are no known pathogenic methanogens. None of the other Archaea are pathogens either, but considering the conditions under which they grow, this is not surprising. Methanogens also form close syntrophic associations with heterotrophic Bacteria that generate hydrogen (i.e. use protons as the terminal electron acceptor). Hydrogen-generating heterotrophism is only energetically-favorable where the ambient concentration of hydrogen is extremely low. Methanogens associate with these organisms, utilizing the hydrogen they generate for methanogenesis, and keep the hydrogen concentration low enough for the heterotrophs to make a living. Neither of these organisms could persist in the environment alone, but together they are successful.

Other than the fact that they all make a living the same way, methanogens are a diverse phenotypic and ecological group. The methanogens are divided into 5 orders, the characteristics of which are given below.

### **1. Order *Methanobacteriales***

The order *Methanobacteriales* currently encompasses non-motile methanogens with pseudomurein cell walls and C20 and C40 isoprenyl glycerol ethers in their membranes. The order contains two families namely Family I. *Methanobacteriaceae* and family II. *Methanothermaceae*.

**Family I *Methanobacteriaceae*.** Family *Methanobacteriaceae* contains four morphologically distinct genera. (a) The 13 species of the genus *Methanobacterium* [58-,71] are rod to filamentous cells. Some species are thermophilic and a few are alkaliphilic and found in various freshwater habitats. *Methanobacterium subterraneum*, an isolate of a deep granitic groundwater is alkaliphilic, eurythermic and halotolerant methanogen [68]. Only six species (*M. formicicum*, *M. defluvii*, *M. oryzae*, *M. palustre*, *M. subterraneum*, *M. thermoflexum*) can use formate. Three species (*M. bryantii*, *M. formicicum* and *M. palustre*) can grow on 2-propanol/CO<sub>2</sub>. All species are able to grow on H<sub>2</sub>+CO<sub>2</sub>. Provided the G+C values were correctly determined, the broad range 29 to 62 mol% indicates that the genus *Methanobacterium* is still heterogenous and composed of more than one genus. The type species is *M. formicicum* [58-,60]. (b) The genus *Methanothermobacter* [8] was proposed for the inclusion of thermophilic methanogens such as *M. thermoautotrophicum* [72] and *M. wolfei* [73]. *M. thermoalcaliphilum* [74,75] and *M. thermoformicicum* [76-,78] are considered as synonymous of *M. thermoautotrophicum* [72]. The proposal has now been accepted and a new genus *Methanothermobacter* created to include 3 species namely *Methanothermobacter thermoautotrophicus* comb. nov., *Methanothermobacter wolfeii* comb. nov. and *Methanothermobacter marburgensis* sp. nov. [79]. (c) The seven members of genus *Methanobrevibacter* are neutrophilic mesophilic short rods, often forming pairs or chains and the G+C content varies from between 28 to 32 mol% [80-,85]. Each species inhabits a specialised habitat. *M. ruminantium*, the type species, is the predominant methanogen in the bovine rumen [80]. It requires cofactors for growth, but like *M. smithii* and *M. cuticularis* can use formate. *M. smithii* is abundant in sewage sludge and intestinal tracts of animals and man [85]. *M. arboriphilus* does not use formate and was isolated from wetwood of living trees [81]. *Methanobrevibacter curvatus*, *M. cuticularis* [82], and *M. filiformis* [83] have been isolated from gut of a subterranean termite recently whereas *M. oralis* has been isolated from human

subgingival plaque [84]. (d) The two species of the genus *Methanosphaera* are Gram-positive spherical-shaped organisms which have been isolated from feces of man [86] and rabbit [87] and are generally observed in the digestive tracts of animals. The G+C content is 23 to 26 mol%. Both species require both methanol and H<sub>2</sub> as substrates for methanogenesis and are unable to use H<sub>2</sub> plus CO<sub>2</sub> or formate. Their inability to reduce CO<sub>2</sub> to CH<sub>4</sub> is due to the lack of an active or the presence of an inactive CO<sub>2</sub> reductase system and methyltetrahydromethanopterin:coenzyme M methyltransferase [88]. The type species is *M. stadtmaniae* [86].

**Family II *Methanothermaceae*.** Family *Methanothermaceae* consists of the single genus *Methanothermus* and its 2 species [89,90]. Both the species are extreme thermophiles and have been isolated from a specific habitat (volcanic springs). The temperature optimum is 80°C. The cells are rod shaped, contain a double-layered wall and have a mol G+C content of 33-34%. As hydrogenotrophic methanogens, they use only hydrogen and carbon dioxide with prototrophic growth. The type species is *M. fervidus* [89].

## 2. Order *Methanococcales*

Boone *et al.* [8] proposed a thorough reorganization of this order. The order now contains two families and four genera (Figure 1) of hydrogenotrophic methanogens isolated essentially from marine and coastal environments. All species are irregular cocci, contain proteinaceous cell walls and are motile by a polar tuft of flagella. Cells lyse quickly in detergents. C20 isopranyl glycerol ethers are abundant and C40 ethers are absent excepted in "*Methanocaldococcus jannaschii*". All species use both H<sub>2</sub> and formate as electron donors, and are prototrophs, except the three species of "*Methanocaldococcus*" and "*Methanoignis igneus*" which are unable to utilize formate. Growth is often stimulated by selenium.

**Family I *Methanococcaceae*.** Family *Methanococcaceae* contains two genera. (a) The genus *Methanococcus* includes five mesophilic species (including 1 synonymous) whose G+C content varies between 30 to 41 mol% [91-97]. The type species is *M. vannielii* [91]. *Methanococcus deltae* has been recognized as a synonym of *M. maripaludis* [96]. "*Methanococcus aeolicus*" was included in a genetic study but its characteristics have never been formally described [96]. (b) The genus "*Methanothermococcus*" has been proposed to include the thermophilic species *M. thermolithotrophicus* [8,97].

**Family II "*Methanocaldococcaceae*".** Family "*Methanocaldococcaceae*" has been recently proposed to include two thermophilic genera. The G+C ranges from 31 to 33 mol%. (a) "*Methanocaldococcus jannaschii*" [98], an extreme thermophile isolated from a hydrothermal vent on the East Pacific rise, is the fastest growing methanogen known to date (generation time = 30 min). Two new species, *M. fervens* and *M. vulcanius*, have been recently described in genus *Methanococcus* but is to be reclassified in the genus "*Methanocaldococcus*" [99]. (b) "*Methanoignis igneus*" [100] is the only species in the new genus proposed by Boone *et al.* [8].

## 3. Order *Methanomicrobiales*

The order *Methanomicrobiales* comprises three families and 9 genera [8,101] of hydrogenotrophic methanogens.

**Family I Methanomicrobiaceae.** Family *Methanomicrobiaceae* contains 7 genera with a variety of different morphologies which includes small rods, highly irregular cocci, and plane-shaped cells. The cell walls are proteinaceous and the lipids include both C20 and C40 isopranyl glycerol ethers. The G+C range of the family is 39 to 50 mol%. Almost all strains can use formate and some secondary alcohols. (a) The genus *Methanomicrobium* includes the single mesophilic species, *M. mobile* whose G+C content is 49 mol%. [102]. It is a slightly curved rod, sluggishly motile with a polar flagellum. It was isolated from bovine rumen and has a complex nutritional requirement which includes rumen fluid. An unidentified growth factor found in rumen fluid could be replaced by extracts of *Methanobacterium thermoautotrophicum* [103]. (b) The genus *Methanolacinia* has been created to include the reclassified species *Methanomicrobium paynteri* [104] as *Methanolacinia paynteri* [105]. *Methanolacinia paynteri* a short and irregular non motile rod, isolated from marine sediments is unable to use formate. Cells lyse in detergents. The G+C content is 45 mol%. (c) The genus *Methanogenium* contains five species isolated from various environments [106-,109]. Morphologically they are highly irregular cocci, stain Gram-negative and are nonmotile but do possess flagella. Cell walls are composed of regular protein subunits. Cells readily lyse in dilute detergents. They require growth factors and use formate. The G+C content varies from 47 to 52 mol%. Two species can use CO<sub>2</sub> + secondary alcohols to form methane. *Methanogenium frittonii* is a thermophilic species [108] whereas *M. frigidum*, which has an optimum temperature for growth of 15°C was isolated from Ace Lake in Antarctica and is a psychrophile [107]. The type species is *M. cariaci* [106]. (d) The genus *Methanoculleus* [110] consists of five mesophilic species (including 1 synonymous) of highly irregular non motile cocci which stain Gram negative [110-,113] and one thermophilic species [114,115]. Formate is used by five species. The G+C content range is between 49 to 62 mol%. The type species has been proposed as *M. olentangyi* [91,110] and *M. bourgense* [111] as a synonym of *M. olentangyi* [8]. (e) The genus *Methanoplanus* comprises three species of plane-shaped organisms with polar tuft of flagella [116-,118]. The cell walls contain at least one major glycoprotein. Formate is used for methanogenesis. The type species is *M. limicola* [116]. One species is an endosymbiont of marine ciliates and is found in close association with microbodies, and is thought to provide hydrogen to the methanogen [117]. The methanogen functions as an electron sink in the oxidation steps of the carbon flow in the ciliates. These symbiotic relationship is thought to be responsible for a total conversion of metabolites to carbon dioxide and methane in marine sediments. Recently a new species, *M. petrolearius*, has been isolated from an oil well [118]. The G+C range of the family is 39 to 50 mol%. (f) Zellner *et al.* [119] have proposed to reclassify *Methanogenium tationis* [120] and *M. liminatans* [121] in a new genus *Methanofollis*. These species use formate and have a G+C content of 54-60 mol%. (g) *Methanocalculus* is a newly described genus which encompasses the irregular coccoid *M. halotolerans*, an isolate from an offshore oil well [122]. It is a hydrogenotrophic halotolerant methanogen which grows optimally at 5% and tolerates up to 12% NaCl. The 0 to 12% NaCl growth range is the widest reported to date for any hydrogenotrophic methanogen including members of the orders *Methanobacteriales*, *Methanococcales* and *Methanomicrobiales*. Further investigation may lead to the reclassification of this genus to the family *Methanocorpusculaceae*.

**Family II *Methanocorpusculaceae*.** Family *Methanocorpusculaceae* [123] contains one genus, *Methanocorpusculum*, and five species (including 1 synonymous) [123-,127] of mesophilic, small coccoid methanogens with monotrichous flagellation. They use H<sub>2</sub>/CO<sub>2</sub> and formate and some species can use 2-propanol/CO<sub>2</sub>. The type species proposed by Boone *et al.* [8] is *M. parvum*, a tungsten requiring bacterium. It is the first hydrogenotrophic methanogen to possess a cytochrome of *b*- or *c*-type, probably involved in the oxidation of 2-propanol. *Methanocorpusculum aggregans* [125] has been recently recognized as synonymous of *M. parvum* [8]. The mol% G+C of this genus is 48 to 52.

**Family III "*Methanospirillaceae*".** The creation of family "*Methanospirillaceae*" has been proposed recently by Boone *et al.* [8] to include the single genus *Methanospirillum*. Members of the genus are mesophilic and have been reported from various habitats. However, only one species, *Methanospirillum hungatei*, has been described so far [128]. Cells are curved rods and often form filaments several hundred µm in length. Cells present polar, tufted flagella and are sheathed. The cell wall composition contains 70 % amino acids, 11 % lipids, and 6.6 % carbohydrates [129]. The cytoplasmic membrane and cell sheath have also been isolated and their composition determined [130]. The type species uses H<sub>2</sub>+CO<sub>2</sub> and formate, and some strains are able to use 2-propanol and 2-butanol as hydrogen donors for methanogenesis from CO<sub>2</sub> [131,132]. *Methanospirillum hungatei* gave a positive chemotactic response to acetate [133]. The G+C content is 45-49 mol%. This hydrogenotrophic methanogen shows the best affinity to hydrogen and is always utilized for isolation of syntrophic bacteria, when sulfate reducers are not used.

#### **4. Order "*Methanosarcinales*"**

This new order proposed by Boone *et al.* [8] regroups all the acetotrophic and/or methylotrophic methanogens into two families (Figure 1).

**Family I *Methanosarcinaceae*** [134]. Family *Methanosarcinaceae* contains six genera and 21 species (including 1 synonymous). (a) The genus *Methanosarcina* represents the acetotrophic methanogens which predominate in many anaerobic ecosystems where organic matter is completely degraded to CH<sub>4</sub> and CO<sub>2</sub>. They are found in freshwater and marine mud, anoxic soils, animal-waste lagoons, and anaerobic digestors. Some are the most versatile methanogens, able to use H<sub>2</sub>-CO<sub>2</sub>, acetate and methyl compounds (methanol, methylamines), including six mesophilic species (including 1 synonymous), *Methanosarcina barkeri*, the type species [135-,137], *M. acetivorans* [138] *M. mazei* [139-,141], *M. siciliae* [142-,144], and *M. vacuolata* [145] and only one thermophilic species, *M. thermophila* [146]. They share a characteristic pseudosarcina cell arrangement and morphology.

Several isolates have been described which use H<sub>2</sub>-CO<sub>2</sub> and methyl compounds, and have a coccoid morphology as *M. frisia* transferred from *Methanococcus frisius* [147,148] and recognized later as a synonym of *M. mazei* [149]. This intermediate form, *M. mazei*, has a morphology similar to both pseudosarcina and the cocci during different phases of its life cycle [150-,152]. A complex life cycle involving the release of single cells may provide a mechanism for cell dispersal during unfavorable growth conditions, whereas a limited cycle facilitates colony division during growth in favorable conditions. With strain LYC, there is a production of a disaggregatase enzyme that hydrolyses the matrix holding the colony together [151]. Xun *et al.* [152]

have shown that the life cycle of strain S-6 can be controlled by manipulation of growth conditions (magnesium, calcium, and substrate concentration), and by inoculum size.

*Methanosarcina vacuolata* shows the presence of large vacuoles containing gas vesicles but is unable to float in the liquid medium [145]. *Methanosarcina siciliae* was transferred from genus *Methanolobus*. It uses only methyl compounds [142,143] and is similar in this trait to *M. semesiae* [153] but recently an acetoclastic strain of *M. siciliae* was isolated from marine canyon sediments [144]. *Methanosarcina acetivorans* is the only marine species in this genus [138]. *Methanosarcina barkeri* is the most studied acetoclastic methanogen and one of the earliest species of methanogen isolated in axenic culture by Schnell in 1947 [154], but lost and isolated again by Bryant in 1966 (strain MS<sup>T</sup>) and described as the type strain of the species [135-137]. Cells are pseudosarcinae, mostly in small aggregates but sometimes in large masses visible to the unaided eye. They are nonmotile and stain Gram-positive.

The cell wall polymer contains N-acetyl-D-galactosamine and D-glucuronic (or D-galacturonic) acid at a molar ratio of 2:1, as well as a minor amount of D-glucose and traces of D-mannose. Partial hydrolysis of cell wall material yields a disaccharide identical with chondrosine, the N-acetylated and sulfated form of which is known as the repeating unit of animal chondroitin [155]. This unique polymer of *Methanosarcina* is another example of the various eucaryotic resemblances found in *Archaea* and may be termed "methanochondroitin" [155].

Species of *Methanosarcina* contain only C20 isopranyl glycerol ethers. Nutritional requirements vary between species. The hydrogen metabolism during methanogenesis from acetate has been extensively studied [156,157]. H<sub>2</sub> production by the cells appears to be linked to several intracellular redox processes which follow the cleavage of acetate. Belay and Daniels [158] have described the formation of ethane by *M. barkeri* during growth in ethanol supplemented medium; ethanol is converted to ethane using terminal portion of the methanol-to-methane pathway.

In the order "*Methanosarcinales*", the following five remaining genera are obligatory methylotrophic methanogens. These methylotrophs are nonmotile, mostly mesophilic, irregular cocci, using only methanol and methylamines as substrates for methanogenesis. Most biotypes have been isolated from environments with high salt concentrations and some of these are regarded as true hyperhalophilic methanogens.

(b) The genus *Methanolobus* [159] contains five species [160-164]. The type species, *M. tindarius* is an irregular mesophilic coccus isolated from coastal sediments, with a single flagellum, based on electron micrographs [160]. The optimal concentration of NaCl is about 0.5 M. This concentration can reach 1.5 M for *M. oregonensis* [162]. The G+C content range is 39-46 mol%.

(c) The genus *Methanococcoides* includes two species with *M. methylutens* as the type species [165]. Cells lyse readily in SDS. The optimal concentration of NaCl is 0.2-0.6 M, and high concentrations of magnesium (50 mM) are also required. *Methanococcoides burtonii* is a psychrophilic methanogen isolated from Ace Lake in

Antartica; the optimum temperature is 23°C [166]. The mol% G+C of the genus is 40-42.

(d) The genus *Methanohalophilus* encloses four mesophilic, hyperhalophilic species [167-,172]. The type species, *M. mahii* has been isolated from the sediments of the Great Salt Lake, Utah. The optimum salinity for growth is 1-2.5 M NaCl. *Methanohalophilus euhalobius* has been recently transferred from the genus *Methanococcoides* [168,169] and *M. halophilus* from the genus *Methanococcus* [170,171]. The G+C content range of the genus is 38-49 mol%.

(e) The genus "*Methanosalsus*" has been recently proposed to reclassify *Methanohalophilus zhilinae* as "*Methanosalsus zhilinae*" [8,173], an alkaliphilic, halophilic species of methanogen isolated from an Egyptian lake, and able to catabolize dimethylsulfide [173,174]. The mol% G+C is 38.

(f) The genus *Methanohalobium* is represented by only one extremely halophilic species, *M. evestigatum* [175] growing at 25% NaCl and at 50°C.

**Family II "*Methanosaetaceae*".** Family "*Methanosaetaceae*" includes all the obligatory acetotrophic methanogens grouped into the genus *Methanosaeta* currently consists of two species. The type species, *M. concilii*, forms an immunologically cohesive group [176-,179]. The rod-shaped cells, 0.8 x 2 µm, form long sheathed filaments that often form floc-like aggregates. The outer layer of the cell wall consists of proteins. Only C20 isopranyl glycerol ethers are present. Acetate is the sole substrate for methanogenesis, with a doubling time of 4-7 days at 37°C. Formate is split to H<sub>2</sub>-CO<sub>2</sub>. The bacterium was first described as *Methanothrix soehngeni* [180-,183] but the cultures were recognized as non-axenic [177-,179] and as a consequence this name has been rejected [184]. A thermophilic gas-vacuolated species, *M. thermoacetophila*, has been described but culture again was found not to be axenic [185]. Another thermophilic strain has been isolated and characterized [186]. Recently the name *M. thermoacetophila* was rejected and replaced with *M. thermophila* [184,187]. The mol% G+C of the genus is 50-61.

### **5. Order "*Methanopyrales*"**

Boone *et al.* [8] have proposed to include the genus *Methanopyrus* into a new order, "*Methanopyrales*". This order currently represents a novel group of methanogens growing at 110°C and unrelated to all other known methanogens [188-,191]. The single family, "*Methanopyraceae*", includes only one species, *Methanopyrus kandleri* (Figure 1). *Methanopyrus kandleri* is a hydrogenotrophic, hyperthermophilic archaeum which stains Gram positive. It has been isolated from a hydrothermally heated deep sea sediment and from a shallow marine hydrothermal system. In the presence of sulfur, H<sub>2</sub>S is formed and cells tend to lyse. The cell wall consists of a new type of pseudomurein which contains ornithine and lysine but not N-acetylglucosamine. The pseudomurein is covered by a detergent-sensitive protein surface layer. The core lipid consists exclusively of phytanyl diether. The G+C content is 60 mol%.

### **3. Order *Archaeoglobales*, the sulfate reducers.**

Members of the genus *Archaeoglobus* have been isolated only from deep-sea hydrothermal vents & heated marine sediments.

It is a thermophilic (85°C) coccus; some species are motile with tufted flagella (much like *Thermococcus*) and others are nonmotile.

*Archaeoglobus* can be grown either of two ways. It can grow heterotrophically by a unique pathway - it uses the methanogenic pathway in reverse!



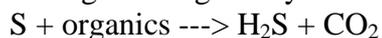
It also grows autotrophically by sulfate reduction:



#### **4. The sulfur-metabolizing order *Thermococcales***

This order includes members of the genera *Thermococcus* & *Pyrococcus*. *T. celer* is the most primitive organism known, i.e. it is closer to the 'root' of the universal tree than any other known living organism. *Thermococcus* is a neutral pH heterotroph, and is thermophilic (75 - 90). *Pyrococcus* is similar but grows at higher temperatures, about 100°C.

These organisms grow by anaerobic sulfur respiration



They are common in hot marine sediments, especially in deep-sea hydrothermal vent areas. They are motile, via a distinctive tuft of polar flagella.

#### **5. The cell wall-less order *Thermoplasmatales***

These organisms are thermoacidophilic heterotrophs. They are facultatively anaerobic - O<sub>2</sub> or sulfur can serve as terminal electron acceptors.

They are acidophilic, most isolates growing best at pH 2, but some grow as low as pH 0.8! (>0.1M HCl). They are also moderately thermophilic, preferring around 60°C for growth.

*Thermoplasma* has been isolated exclusively from smouldering coal refuse piles, and it is presumed that subterranean coal deposits are their natural habitat. All isolates are monotrichously flagellated motile cocci and lack any cell wall. Crosslinking of the carbohydrate chain of membrane glycoproteins provide what little cell rigidity and osmotic tolerance they have.

Like *Archaeoglobus*, they reveal their methanogenic ancestry by containing F420 (the major methanogenic cofactor) & other components of the methanogenic pathway, but it is not known what use, if any, they make of these.

#### **8. Kingdom *Korarchaeota*:**

There is very little information on this kingdom as the cells have only just been cultured. The studies on this kingdom will be very important as they fall very close to the root of the tree of life. Therefore the molecular information contained in them may contribute to our knowledge and understanding of ancient organisms.

#### **9. Evolution and Life at high temperatures**

How do cells cope with high temperatures and what is the maximum temperatures at which life could exist? These questions have yet to be answered conclusively but some pointers are already available for discussion.

(a) Heat stability of Biomolecules: Most proteins from hyperthermophiles have the same structural features. For example the amino acids around the active sites of thermostable enzymes and their mesophilic counterpart are the same. However, thermostable proteins do tend to have highly hydrophobic cores (probably increases internal "sticking" and in general have more "salt bridges" (ionic interactions between amino acids) proteins and vice versa. on the surface. Therefore, it is now thought that subtle changes which is reflected in protein folding, rather than drastic gross changes, are sufficient to render heat labile proteins to heat resistant and vice versa.

Chaperonins (heat shock proteins, aka HSP) refold partially denatured proteins back to "life". In *Pyrococcus*, the HSP, called thermosome, increases in concentration dramatically (to 80% of the cell content) when the culture is grown at 108°C (which is at the temperature limit for growth). This is a type of protection mechanism. The culture with thermosome can withstand autoclaving (121°C) for 1 hour!!!

(b) DNA stability: (i) High temperature depurination of DNA is decreased due to the presence of potassium cyclic 2,3-diphosphoglycerate (K<sup>+</sup> is actually responsible for this) but only in some hyperthermophiles. (ii) Positive supercoiling of DNA (rather than negative coiling) by reverse gyrase stabilises DNA against denaturation at high temperatures. (iii) A minor groove DNA binding protein, termed *Sac7d*, found in *Sulfolobus*, (crenarchaeotes) increases the melting temperature of DNA by 40°C. *Sac7d* also kinks the DNA and is therefore thought to play a role in gene expression. The euryarchaeotes contain highly basic histone-like proteins which have homology with eucaryotic histone proteins. They wind and compact DNA into a nucleosome-like structure and thereby protects the DNA from heat denaturation. The presence of *Sac7d* and histone-like proteins may explain why the transcriptional apparatus is more in line with the eucaryotic system.

(c) Lipid Stability: Contain heat resistant dibiphytanyl diether lipid membranes which forms a monolayer rather than a bilayer (see above). Stops the cell membrane from being pulled apart.

(d) Stability of monomers: The thermostability of monomers is of more significance than macromolecules- dictate the upper temperature of life. (1) ATP and NAD<sup>+</sup> hydrolyse rapidly- half life of 30 mins at 120°C.

## **10. The Limits to microbial existence:**

Extrapolate and hypothesise from the knowledge given above about the limits to life at high temperatures:

- a. Cells require water for life but water above 100°C is steam. Life therefore at temperatures > 100°C will be restricted to hydrothermal vents of the sea floor. Most *Archaea* cultured in the laboratory above 100°C are from such environments:-

- Black smokers- 250 – 350°C form upright metallic sulfide structures called chimneys. The temperature gradient is 250 °C (inside) to 2 °C outside. *Archaea* have been isolated from chimney walls but water >250 °C appears to be sterile.
- b. Laboratory experiments suggest that life can live at the limits of 140 to 150 °C beyond which biomolecules become heat labile. For example, ATP will not be an energy currency but will have to be in some other form.

### **11. Hyperthermophiles Archaea and Microbial Evolution:**

Have *Archaea* adapted to extreme environments or co-evolved and flourished with such extremes during earth's formation, that is, relics of ancient life. Their studies will reveal intersting principles of early life.

- 16S rRNA sequencing and sequence analysis suggests the later. *Aquifex* and *Thermotogales*, domain *Bacteria* which are also extreme (hyperthermophiles) and have slow evolutionary clocks. Their molecular clocks are slower and therefore they evolved at a slower rate than did *Eucarya* and *Bacteria* – a fast clock means adaptation. Phylogenetic trees shows shorter branches.
- Life under extremes are under strong evolutionary pressure to maintain their genes essential for their survival and therefore beyond a certain point additional genetic changes will not be of any further benefits.
- *Korarchaeota* clocks are the slowest and found in volcanic hot springs.
- H<sub>2</sub> as an electron acceptor with S<sup>0</sup>, NO<sub>3</sub><sup>-</sup> and Fe<sup>3+</sup> as electron donors under anaerobic, high temperature and dark subsurface (protected from UV radiation) conditions (primordial early earth conditions) by *Archaea* may be an ancient relic of metabolism.

### **12. Comparative genomics of Archaea:**

Refer to the attached published article

### **13. Some Useful References**

Science 283 p 1476 -5th March 1999 "Mitochondrial Evolution" contains a discussion of the Endosymbiotic theory of eukaryote evolution.

Nature 399 p 323 27 May 1999 Presents evidence for lateral gene transfer between Archaea and Bacteria. Microbiology and Molecular Biology Reviews December 1998 p1436-1491.

Woese, (1977) PNAS 74:5088-5090].

[Woese,(1990) PNAS 87:4576-4579].

Takai-K; Sako-YA molecular view of archaeal diversity in marine and terrestrial hot water environments FEMS-MICROBIOLOGY-ECOLOGY. FEB 1999; 28 (2) : 177-188

Molecular phylogenetic survey of naturally occurring archaeal diversities in hot water environments was carried out by using the PCR-mediated small subunit rRNA gene (SSU rDNA) sequencing. Mixed population DNA was directly extracted from the effluent hot water or sediment of a shallow marine hydrothermal vent at Tachibana Bay, or the acidic hot water of hot spring pools at Mt. Unzen, in Nagasaki Prefecture,

Japan. Based on the partial rDNA sequences amplified with an Archaea-specific primer set, the archaeal populations of hot water environments consisted of phylogenetically and physiologically diverse group of microorganisms. The archaeal populations were varied in each sample and subject to its environmental conditions. In addition, the large number of archaeal rDNA sequences recovered from hot water environments revealed the distant relationship not only to the rDNA sequences of the cultivated thermophilic archaea, but also to the sequences of unidentified archaeal rDNA clones found in other hot water environments. The findings extend our view of archaeal diversity in hot water environments and phylogenetic organization of these organisms. (C) 1999 Federation of European Microbiological Societies. Published by

Kardinahl-S; Schmidt-CL; Hansen-T; Anemuller-S; Petersen-A; Schafer-G RP: Schafer, G . The strict molybdate-dependence of glucose-degradation by the thermoacidophile *Sulfolobus acidocaldarius* reveals the first crenarchaeotic molybdenum containing enzyme - an aldehyde oxidoreductase. EUROPEAN-JOURNAL-OF-BIOCHEMISTRY. MAR 1999; 260 (2) : 540-548

In order to investigate the effects of trace elements on different metabolic pathways, the thermoacidophilic Crenarchaeon *Sulfolobus acidocaldarius* (DSM 639) has been cultivated on various carbon substrates in the presence and absence of molybdate. When grown on glucose (but neither On glutamate nor casein hydrolysate) as sole carbon source, the lack of molybdate results in serious growth inhibition. By analysing cytosolic fractions of glucose adapted cells for molybdenum containing compounds, an aldehyde oxidoreductase was detected that is present in the cytosol to at least 0.4% of the soluble protein. with Cl(2)Ind (2,6-dihlorophenolindophenol) as artificial electron acceptor, the enzyme exhibits oxidizing activity towards glyceraldehyde, glyceraldehyde-3-phosphate, isobutyraldehyde, formaldehyde, acetaldehyde and propionaldehyde. At its pH-optimum (6.7, close to the intracellular pH of *Sulfolobus*, the glyceraldehyde-oxidizing activity is predominant. The protein has an apparent molecular mass of 177 kDa and consists of three subunits of 80.5 kDa (alpha), 32 kDa (beta) and 19.5 kDa (gamma). It contains close to one Mo, four Fe, four acid-labile sulphides and four phosphates per protein molecule. Methanol extraction revealed the existence of 1 FAD per molecule and 1 molybdopterin per molecule, which was identified as molybdopterin guanine dinucleotide on the basis of perchloric acid cleavage and thin layer chromatography. EPR-spectra of the aerobically prepared enzyme exhibit the so-called 'desulpho-inhibited'-signal. known from chemically modified forms of molybdenum containing proteins. Anaerobically prepared samples show bath, the signals arising from the active molybdenum-cofactor as well as from the two [2Fe-2S]-clusters. According to metal-, cofactor-, and subunit-composition, the enzyme resembles the members of the xanthine oxidase family. Nevertheless, the melting point and long-term thermostability of the protein are outstanding and perfectly in tune with the growth temperature of *S. acidocaldarius* (80 degrees C).The findings suggest the enzyme to function as a glyceraldehyde oxidoreductase in the course of the nonphosphorylated Entner-Doudoroff pathway and thereby may attribute a new physiological role to this class of enzyme.

Hopfner-KP; Eichinger-A; Engh-RA; Laue-F; Ankenbauer-W; Huber-R; Angerer-B RP: Hopfner, KP Crystal structure of a thermostable type B DNA polymerase from *Thermococcus gorgonarius*. PROCEEDINGS-OF-THE-NATIONAL-ACADEMY-

Most known archaeal DNA polymerases belong to the type B family, which also includes the DNA replication polymerases of eukaryotes, but maintain high fidelity at extreme conditions. We describe here the 2.5 Angstrom resolution crystal structure of a DNA polymerase from the Archaea *Thermococcus gorgonarius* and identify structural features of the fold and the active site that are likely responsible for its thermostable function. Comparison with the mesophilic B type DNA polymerase gp43 of the bacteriophage RB69 highlights thermophilic adaptations, which include the presence of two disulfide bonds and an enhanced electrostatic complementarity at the DNA-protein interface. In contrast to gp43, several loops in the exonuclease and thumb domains are more closely packed; this apparently blocks primer binding to the exonuclease active site. A physiological role of this "closed" conformation is unknown but may represent a polymerase mode, in contrast to an editing mode with an open exonuclease site. This archaeal 13 DNA polymerase structure provides a starting point for structure-based design of polymerases or ligands with applications in biotechnology and the development of antiviral or anticancer agents.

The most well studied similarity between archaea and eukaryotes is transcription.  
Zillig EMBO J. 2 1291-1294 1983.

Madigan, M.T., Martinko, J.M. and Parker, J. Brock, *Biology of Microorganisms*,  
Prentice Hall, 8th Edition, 1997

M.Ciaramella et al, *Molecular biology of extremophiles*, *World Journal of  
Microbiology and Biotechnology*, Vol 11, pp 71-84, 1995

Danson et al., *Archaeobacteria*, *Biochemistry and Biotechnology*, London:  
Biochemical Society, 1992

Brown, J.W. et al, *Gene Structure, Organisation and Expression in Archaeobacteria*.  
*CRC Critical Reviews in Microbiology*, Vol 16, No. 4, 1989

Questions for thought:

1. How do you suppose that organisms like *Halobacterium* were studied for so long without it being realized that it really wasn't a lot like its supposed relatives, *Pseudomonas*?
2. If *Archaea* are specifically related to eukaryotes to the exclusion of *Bacteria*, why don't we consider *Archaea* to be eukaryotes (even if primitive ones)?

# Comparative genomics of archaea: how much have we learned in six years, and what's next?

Kira S Makarova and Eugene V Koonin

Address: National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA.

Correspondence: Eugene V Koonin. E-mail: koonin@ncbi.nlm.nih.gov

Published: 16 July 2003

*Genome Biology* 2003, **4**:115

The electronic version of this article is the complete one and can be found online at <http://genomebiology.com/2003/4/8/115>

© 2003 BioMed Central Ltd

## Abstract

Archaea comprise one of the three distinct domains of life (with bacteria and eukaryotes). With 16 complete archaeal genomes sequenced to date, comparative genomics has revealed a conserved core of 313 genes that are represented in all sequenced archaeal genomes, plus a variable 'shell' that is prone to lineage-specific gene loss and horizontal gene exchange. The majority of archaeal genes have not been experimentally characterized, but novel functional pathways have been predicted.

"A phylogenetic analysis based upon ribosomal RNA sequence characterization reveals that living systems represent one of three aboriginal lines of descent: (i) the eubacteria, comprising all typical bacteria; (ii) the archaeobacteria, containing methanogenic bacteria; and (iii) the urkaryotes, now represented in the cytoplasmic component of eukaryotic cells."

*CR Woese and GE Fox, 1977 [1]*

## Archaea before and after genomes

The quotation above neatly summarizes what is arguably one of the most important scientific discoveries of the twentieth century (rather remarkably, this quote is the entire abstract of Woese and Fox's groundbreaking article [1]). So profound are its implications that the debate rages to this day: did Carl Woese and George Fox really discover a new domain of life, which is equal in status to bacteria and eukaryotes [2,3], or is it 'merely' an unusual branch of bacteria [4-7]? This debate is reflected even in the different names that, 25 years after their description as a distinct, third line of the evolution of life, are still applied to this group of organisms: on the one hand, archaea, in adherence with the three-domain interpretation, and on the other archaeobacteria, emphasizing the purported affinity with bacteria. Of course, Woese and Fox did not actually discover

these unusual organisms; some of the would-be archaea have been known for decades and their unusual properties, such as extreme halophilic and extreme thermophilic phenotypes, have been described in considerable detail (see, for example, [8-10]). The revolutionary aspect of Woese and Fox's work was subtler and more profound: by comparing certain parts of the genomic sequences of various organisms, they came up with a three-domain classification of life, in which a group of prokaryotes they designated archaeobacteria has been accorded the status of a distinct domain (subsequently renamed archaea, to emphasize the fundamental separation from other domains), on an equal footing with bacteria and eukaryotes. Numerous microbiologists had seen archaea before, but without Woese and Fox's foray into genome analysis no-one recognized these organisms for what they really were. Their way of comparing genome sequences was, by today's standards, extremely crude, as they analyzed not even sequences but oligonucleotide catalogues of rRNA genes. It is all the more astounding that the principal conclusion achieved with this 'primitive' approach stands to this day, 25 years and 16 complete (and several more nearly complete) archaeal genome sequences later (Table 1).

In the years following Woese and Fox's breakthrough [1], many unique features of archaea have become apparent. To

**Table 1****Completely sequenced archaeal genomes**

| Species  | Abbreviation | Optimal growth temperature (°C) | Lifestyle and other features   | Number of proteins* | Number (%) proteins in COGs | Date of genome release | Reference |
|--|--------------|---------------------------------|--|---------------------|-----------------------------|------------------------|-----------|
| <b>Euryarchaeota</b>                                   |              |                                 |  |                     |                             |                        |           |
| <i>Archaeoglobus fulgidus</i> DSM                      | <b>Afu</b>   | 83                              | Anaerobic, sulfate-reducing chemolitho- or chemorgano-autotroph, motile  | 2,420               | 1,953 (81%)                 | 1997                   | [124]     |
| <i>Halobacterium</i> sp. NRC-1                         | <b>Hsp</b>   | 37                              | Aerobic chemorganotroph, obligate halophile, with a cell envelope; motile; two extrachromosomal elements   | 2,622               | 1,809 (69%)                 | 2000                   | [125]     |
| <i>Methanocaldococcus jannaschii</i>                   | <b>Mja</b>   | 85                              | Chemolithoautotroph, strict anaerobe, methanogen, motile; two extrachromosomal elements  | 1,758               | 1,448 (82%)                 | 1996                   | [27]      |
| <i>Methanopyrus kandleri</i> AV19                      | <b>Mka</b>   | 110                             | Chemolithoautotroph, strict anaerobe, methanogen, with high cellular salt concentration  | 1,691               | 1,253 (74%)                 | 2002                   | [45]      |
| <i>Methanosarcina acetivorans</i> C2A                  | <b>Mac</b>   | 37                              | Chemolithoautotroph, anaerobe possibly capable of aerobic growth; nitrogen-fixing, versatile methanogen; motile, and able to form multicellular structures | 4,540               | 3,142 (69%)                 | 2002                   | [55]      |
| <i>Methanosarcina mazei</i> Goe 1                      | <b>Mma</b>   | 37                              | As for <b>Mac</b>  | 3,371               | N/A                         | 2002                   | [54]      |
| <i>Methanothermobacter thermoautotrophicus</i> delta H | <b>Mth</b>   | 65                              | Chemolithoautotroph, strict anaerobe, nitrogen-fixing, methanogen  | 1,873               | 1,500 (80%)                 | 1997                   | [126]     |
| <i>Pyrococcus horikoshii</i>                           | <b>Pho</b>   | 96                              | Anaerobic heterotroph, sulfur enhances growth; motile  | 1,801               | 1,425 (79%)                 | 1998                   | [127]     |
| <i>Pyrococcus abyssi</i>                               | <b>Pab</b>   | 96                              | As for <b>Pho</b>  | 1,769               | 1,506 (85%)                 | 2001                   | [128]     |
| <i>Pyrococcus furiosus</i> DSM 3638                    | <b>Pfu</b>   | 96                              | As for <b>Pho</b>  | 2,065               | N/A                         | 2001                   | [129]     |
| <i>Thermoplasma acidophilum</i>                        | <b>Tac</b>   | 59                              | Facultative anaerobe, chemorganotroph, thermoacidophilic, anaerobically able to metabolize sulfur; motile, with a plasma membrane                          | 1,482               | 1,261 (85%)                 | 2000                   | [96]      |
| <i>Thermoplasma volcanium</i>                          | <b>Tvo</b>   | 60                              | As for <b>Tac</b>  | 1,499               | 1,277 (85%)                 | 2000                   | [130]     |
| <b>Crenarchaeota</b>                                   |              |                                 |  |                     |                             |                        |           |
| <i>Pyrobaculum aerophilum</i>                          | <b>Pae</b>   | 100                             | Facultative nitrate-reducing anaerobe  | 1,840               | 1,236 (67%)                 | 2002                   | [131]     |
| <i>Aeropyrum pernix</i>                                | <b>Ape</b>   | 90                              | Aerobic chemorganotroph; sulfur enhances growth  | 2,605               | 1,529 (59%)                 | 1999                   | [132]     |
| <i>Sulfolobus solfataricus</i>                         | <b>Sso</b>   | 80                              | Aerobe metabolizing sulfur; thermoacidophilic chemorganotroph; motile  | 2,977               | 2,207 (74%)                 | 2001                   | [97]      |
| <i>Sulfolobus tokodaii</i>                             | <b>Sto</b>   | 80                              | As for <b>Sso</b>  | 2,826               | N/A                         | 2001                   | [133]     |

\*According to the original genome annotation.

begin with, many of these organisms thrive under conditions that, by the usual standards of biology, seem unimaginable, such as in the water in the vicinity of the hydrothermal vents called 'black smokers' heated to over-boiling temperatures and saturated with hydrogen sulfide, or in extreme salinity [11-13]. In the most extreme hyperthermophilic habitats, archaea are, in fact, the only detectable life forms. In more moderate environments, archaea coexist with bacteria and eukaryotes, and their ecological importance is being increasingly recognized [14]. The first molecular biological studies showed that archaea are highly unusual and clearly distinct from bacteria at the molecular level. In particular, the structure of the membrane glycerolipids in archaea is different from that of bacterial and eukaryal cells, and archaea do not contain murein, the predominant component of bacterial cell walls [15,16].

But the most striking differences between archaea and bacteria are seen in the organization of their information-processing systems. The structures of ribosomes and chromatin, the presence of histones, and sequence similarity between proteins involved in translation, transcription, replication and DNA repair all point to a closer relationship between archaea and eukaryotes than between either of these and bacteria [17-21]. Moreover, the key components of the DNA replication machinery - such as the polymerases involved in elongation and initiation and the replicative helicases - are not homologous, or at least not orthologous, in archaea and eukaryotes on the one hand, and bacteria on the other [17,22]. This observation led to the hypothesis that replication of double-stranded DNA as the principal form of replication of the genetic material was 'invented' twice, independently: once in bacteria and once in the ancestor of archaea and eukaryotes [22,23]. In contrast many - although not all - of the metabolic pathways of archaea more closely resemble their bacterial rather than eukaryotic counterparts [24-26]. These studies support the status of archaea as a distinct domain of life with specific connections to eukaryotes, and emphasize the unusual and unique nature of archaeal genomes.

The new age of archaea began in 1996 with the whole-genome shotgun sequencing of the first archaeal genome, that of *Methanococcus* (now *Methanocaldococcus*) *jannaschii* [27]. The *Methanococcus* 'genomescape' at first looked largely mysterious, with clear functional assignments produced for only 38% of the genes [27]. A more detailed computational analysis that pushed the methodology available at the time to its limits yielded general functional predictions for up to 70% of the genes, showing that a solid connection between the genomes of archaea and those of other, better known forms of life did exist [24]. Nevertheless, the fact remained that, more than anything, the first sequenced archaeal genome revealed the depth of our ignorance of the biology of this remarkable group of organisms. Subsequent genome sequencing, while certainly less extensive than the devoted 'archaeologists' would wish, produced

a rich sampling of genomes of taxonomically diverse archaea (Table 1). This set of completely sequenced genomes includes multiple representatives of the two major divisions of the archaea established by phylogenetic analysis of rRNA, namely the Euryarchaeota and the Crenarchaeota [3], as well as the principal ecological types of archaea, such as hyperthermophiles, moderate thermophiles, and mesophiles, as well as halophiles and methanogens; autotrophic and heterotrophic forms, and anaerobes and aerobes are also represented by multiple species (Table 1).

Some potentially important branches of archaea are still missing from sequence databases, however, such as the mysterious Korarchaeota, which might have branched off the trunk of the phylogenetic tree prior to the divergence of the remainder of the archaea [28], and the equally intriguing Nanoarchaea that so far seem to have the smallest genomes of all known cellular life forms [29,30]. These lacunae notwithstanding, the available sampling of archaeal genomes is substantial and is complemented by an even greater diversity of bacterial and eukaryotic genomes that are available for comparative analysis. This article critically assesses the contribution of comparative genomics to our understanding of the functional systems of archaeal cells and their evolution. We pose the following question: what have we learned from comparisons of archaeal genomes that could not easily have been learned by other, more traditional approaches? We suggest some tentative answers, as we see them. What follows is a viewpoint from behind a computer terminal; we realize that, from the experimenter's bench, the perspective might be somewhat different.

### Evolutionary archaeogenomics

From the beginning of comparative genomics, it has been obvious that genome comparisons will yield valuable functional and evolutionary information only within a framework of the rational classification of genes and proteins. In our view, perhaps the most natural form of such a classification is a system of orthologous gene sets, which allows a researcher to analyze the evolutionary fate of each individual gene [31]. Orthologs are homologous genes that evolved from a single ancestral gene in the last common ancestor of the compared genomes, whereas paralogs are genes related via duplication within a genome [32-34]. When duplication(s) succeeds speciation, a family of paralogs in one species should be considered orthologous to the corresponding family in the other species [34]. Inasmuch as orthologous relationships are correctly defined, phyletic (or phylogenetic) patterns of orthologous gene sets help in the prediction of gene functions and provide clues to the prevailing trends in genome evolution (a phyletic pattern is defined, simply, as the pattern of representation of genomes in each orthologous set) [26,31,35,36]. These phyletic patterns are captured in the database of Clusters of Orthologous Groups of proteins (COGs) [37], and here we use COGs for a

**Table 2****The top 15 phyletic patterns in proteins from archaea**

| Pattern*  | Number of COGs<br>(and of the<br>complementary<br>pattern, CP) | Comments and examples  |
|---|--|--|
| AHMMMMTTPPPSA<br>fbatjkavhaasp<br>uschaacoobeoe |  |  |
| +++++   | <b>313</b> (0)   | Archaeal core, including 200 COGs present in both <b>B</b> <sup>†</sup> and <b>E</b> , 34 present in at least one <b>B</b> , 63 present in at least one <b>E</b> , 16 unique for <b>A</b><br><b>CP</b> : Only COG0564, pseudouridylate synthase, 23S RNA-specific pseudouridylate synthase present in all <b>E</b> (in which it has an apparently mitochondrial origin) and <b>B</b> , but not in <b>A</b> . In all <b>A</b> another specific pseudouridylate synthase is present (COG1258)  |
| ---+-----                                       | <b>163</b> (3)   | This pattern reflects a large number of genes acquired via HGT <sup>†</sup> in <b>Mac</b> (see [55]), including F <sub>0</sub> F <sub>1</sub> -type ATP synthase and NADH:ubiquinone oxidoreductase, and a specific signal transduction system based on several apoptosis-related domains<br><b>CP</b> : The small number of such COGs indicates that the archaeal core is almost fully conserved in <b>Mac</b>  |
| -+-----   | <b>79</b> (14)   | This pattern reflects a substantial amount of HGT in <b>Hsp</b> ; see [125]  |
| +-----  | <b>47</b> (7)  | This pattern consists of COGs including four methanogens and <b>Afu</b> ; these organisms specifically share several metabolic pathways (see [45]). The set includes subunits of coenzyme F <sub>420</sub> -reducing hydrogenase, formylmethanofuran dehydrogenase, CO dehydrogenase/acetyl-CoA synthase and other enzymes of energy metabolism. These might have originally evolved in methanogens and subsequently transferred to <b>Afu</b><br><b>CP</b> : Sugar ABC transporter and some fatty acid biosynthesis enzymes are missing from methanogens and <b>Afu</b>   |
| -----   | <b>40</b> (2)  | This pattern is specific for four methanogens, including unique pathways for coenzyme M biosynthesis and reduction and 14 uncharacterized proteins, many of which are likely to be unique enzymes involved in biosynthesis of other specific coenzymes and their utilization<br><b>CP</b> : COG2096, cob(I)alamin adenosyltransferase and COG1058, predicted nucleotide-utilizing enzyme related to molybdopterin-biosynthesis enzyme MoeA, for which functional substitutes remain to be identified   |
| -----+  | <b>33</b> (16)   | A pattern specific for thermophilic methanogens ( <b>Mth</b> , <b>Mja</b> and <b>Mka</b> ), comprising mostly uncharacterized COGs, it includes a specific membrane complex EhaA-EhaP (approximately 18 components) involved in hydrogen production and possibly electron transfer [45,134]<br><b>CP</b> : Specific gene loss: peptide ABC-type transporter, NADH:ubiquinone oxidoreductase, malic enzyme (COG0281), and cysteinyl-tRNA synthetase (COG0215; see text)   |
| -----+  | <b>28</b> (6)  | This pattern reflects a substantial amount of HGT in <b>Sso</b> , including several enzymes of carbohydrate metabolism (beta-glucosidase, alpha-L-fucosidase, and malto-oligosyl trehalose synthase) [97]  |
| +-----  | <b>27</b> (1)  | This reflects a substantial amount of HGT in <b>Afu</b><br><b>CP</b> : COG0449, glucosamine 6-phosphate synthetase, which catalyzes the first step in hexosamine metabolism. A functional substitute remains to be identified  |
| -+-----   | <b>25</b> (4)  | A pattern specific for two mesophilic archaea, probably resulting from independent HGT   |
| ----+-----                                      | <b>23</b> (7)  | This pattern includes genes that might have been acquired via HGT in <b>Mja</b> , in particular three enzymes of biotin biosynthesis: pimeloyl-CoA synthetase (COG1424), dethiobiotin synthetase (COG0132), and adenosylmethionine-8-amino-7-oxononanoate aminotransferase (COG0161)   |
| -----++   | <b>21</b> (13)   | A crenarchaea-specific pattern, including 11 COGs that do not have orthologs outside this lineage. Among genes shared with bacteria but not euryarchaeota are three subunits of aerobic-type CO dehydrogenase and CO dehydrogenase maturation factor. Genes specifically shared with eukaryotes are three ribosomal proteins (S30, S25 and L13E)<br><b>CP</b> : Euryarchaea-specific pattern, including two subunits of archaeal DNA polymerase II and ERCC4-like helicase, division GTPase FtsZ (COG0206) and ATP-dependent protease LonB (COG1067) plus six COGs that do not have orthologs outside this lineage |
| +-----  | <b>20</b> (0)  | Apparent independent HGT to <b>Mac</b> and <b>Afu</b>  |
| ++++-----                                       | <b>19</b> (16)   | Apparent specific gene loss in the <i>Thermoplasma</i> lineage: two subunits of topoisomerase VI (COG1389, 1697), adenylate cyclase of class 2 (COG1437), and predicted exosome subunits (COG1325, COG1931).<br><b>CP</b> : genes apparently acquired via HGT in <i>Thermoplasma</i> , including bacterial nucleoid DNA-binding protein HU (COG0776). See also [96]  |
| ++++-----                                       | <b>18</b> (6)  | Apparent gene loss in <b>Ape</b> , including 9 enzymes of purine biosynthesis [135]  |
| -----++   | <b>17</b> (11)   | Apparent HGT in <i>Pyrococci</i> . Includes two subunits of allophanate hydrolase (COG1984, 2049), two enzymes of carbohydrate metabolism, β-galactosidase (COG1874) and endoglucanase (COG2730)<br><b>CP</b> : Specific gene loss in the <i>Pyrococcus</i> lineage includes five enzymes of heme biosynthesis   |

\*The pattern of appearance within the 13 sequenced archaeal species currently available in the COG database. Species abbreviations are as given in Table 1 and are written vertically. †Abbreviations: **A**, archaea; **B**, bacteria; **E**, eukaryotes; **CP**, complementary pattern; HGT, horizontal gene transfer.

systematic survey of archaeal genomes (most of the phyletic pattern analyses can be done directly on the COG website by using the phyletic pattern search tool [38]).

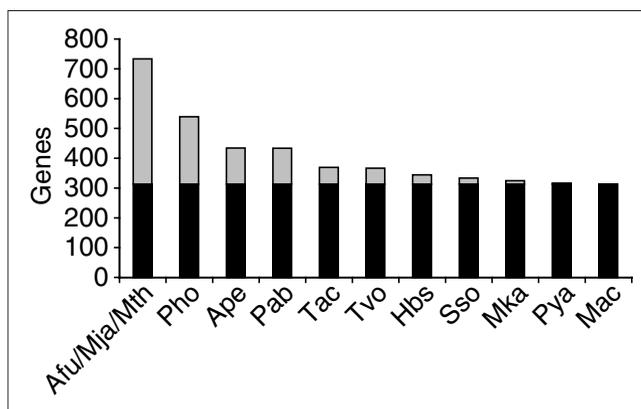
The most common phyletic patterns found in archaea are shown in Table 2. Not unpredictably, the top pattern consists of the 313 COGs that are represented in all archaeal genomes sequenced so far. What is more remarkable is that this apparent conserved core of archaeal genomes has undergone only limited shrinkage since the time it was first defined by comparative analysis of four archaeal genomes [39] (Figure 1). Extrapolating from the effect (or rather the near lack thereof) of the latest additions to the collection of archaeal genomes on the size of the conserved core of archaeal genes, we are compelled to conclude that around 300 genes are shared by all archaea, encode essential functions and have not been subject to non-orthologous gene displacement during archaeal evolution (non-orthologous gene displacement is a widespread phenomenon whereby a gene responsible for an essential function is displaced by an unrelated or distantly related gene responsible for the same function [40]).

Of the COGs represented in all archaea, 16 so far have no members from other domains of life and comprise a unique archaeal genomic signature, whereas 61 are exclusively archaeo-eukaryotic. The majority of the pan-archaeal genes are known to be involved in, or are implicated in, information processing, particularly translation and RNA modification (Figure 2). Strikingly, among the 61 COGs that are uniquely shared by archaea and eukaryotes, only two do not, technically, belong to the information-processing machinery (COG1936, a nucleotide kinase, and COG3642, a protein

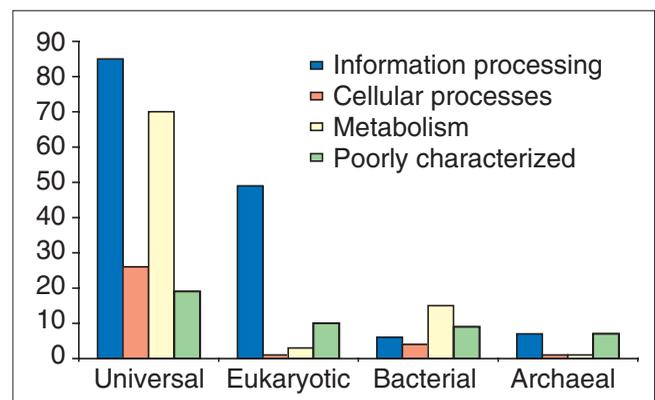
kinase typically fused to a metalloprotease domain); the 10 uncharacterized COGs in this category consist of proteins whose predicted biochemical activity (GTPase, methyltransferase or RNA-binding protein) suggests a role in translation or RNA modification.

Thus, phyletic pattern analysis strongly supports the identity of archaea as a distinct group of organisms with a stable, conserved core of genes that primarily encodes proteins involved in the replication and expression of the genome. Furthermore, there is clearly a subset of genes, again primarily associated with information processing, that is shared by archaea and eukaryotes, to the exclusion of bacteria; this is compatible with the archaeo-eukaryotic affinity suggested by phylogenetic analyses of rRNA and proteins involved in translation, transcription and replication. The fact that this archaeo-eukaryotic component is quantitatively small, however, shows that the process of evolution has been more complex than simple vertical inheritance and has involved extensive horizontal gene transfer (HGT) between archaea and bacteria, at least outside the core gene set [24,25,41]. An intensely mixing pool of genes coding for metabolic enzymes, structural components of the cell and other proteins outside the central information-processing machinery might have existed after the divergence of bacteria and archaea but prior to the separation of the major archaeal and bacterial lineages.

More recent HGT, which has emerged as a major aspect of prokaryotic evolution in general [26,42-44], was apparently prominent in all archaea, although gene exchange with bacteria seems to have been much less extensive in hyperthermophiles than in mesophiles such as *Methanosarcina* or even *Halobacterium* [44,45]. Apparent preferential HGT has been noticed between archaea and hyperthermophilic



**Figure 1**  
The archaeal gene core: changes resulting from the appearance of new genome sequences. Black bars indicate the current set of pan-archaeal genes (313 COGs); gray indicates COGs that are not part of the current pan-archaeal core but are seen to be conserved after the addition of the given genome sequence. The genomes are listed from left to right in chronological order of release of the complete sequence; species name abbreviations are as in Table 1.



**Figure 2**  
Functional breakdown of genes within the conserved archaeal core. 'Universal' indicates genes with orthologs in both bacteria and eukaryotes; 'eukaryotic', genes with orthologs only in eukaryotes; 'bacterial', genes with orthologs only in bacteria; 'archaeal', genes without non-archaeal orthologs. The data on orthology and functional classification are derived from the COGs.

bacteria, such as *Aquifex* and *Thermotoga*; when compared to bacterial mesophiles these bacteria have many more proteins with greater similarity to archaeal than to bacterial homologs [46,47]. With HGT, or more precisely the pivotal role of HGT in evolution, remaining a controversial subject [48], this conclusion has been disputed on the grounds that *Aquifex* and *Thermotoga* might be early-branching bacteria retaining ancestral features in many protein sequences [49]. But this argument seems untenable simply because of the obvious split of the gene complements of these bacteria into 'garden variety' bacterial genes and 'archaeal' genes [50]. The reality of horizontal gene flow from archaea to thermophilic bacteria becomes even more tangible upon examination of the proteins encoded in the genome of *Thermoanaerobacter tengcongensis* [51,52], which contains many more 'archaeal' genes than appear in other bacteria of the *Bacillus-Clostridium* group and to which the early-branching argument would not apply.

Although archaeal hyperthermophiles do not appear to have many genes acquired via HGT from bacteria, at least after the divergence of the archaeal lineages, horizontal gene exchange between archaea themselves might have been extensive. Strikingly, even within the conserved core of archaeal genes, major diversity of phylogenetic tree topologies has been observed ([53] and Y.I. Wolf and E.V.K., unpublished observations). As noted by Nesbo and coworkers [53], "the notion that there is a core of non-transferable genes...has not been proven and may be unprovable". These findings do not invalidate the notion of a core of indispensable genes that are conserved across archaea but suggest a wide spread of xenologous gene displacement, whereby an essential gene is displaced by an ortholog from a distant lineage, typically via an intermediate stage of redundancy [44].

Other phyletic patterns that are common among archaea seem primarily to reflect HGT or gene loss prevalent in individual archaeal lineages (Table 2). Thus, *Methanosarcina*, a mesophile with by far the largest genome among the sequenced archaeal genomes, is represented in numerous COGs that have no other archaeal members but are present in various groups of bacteria. This organism, which coexists with a diverse bacterial biota, appears to be a veritable sink for horizontally acquired bacterial genes [54,55]. Similar, if less dramatic, evidence of apparent horizontal gene transfer was seen in *Halobacterium*, *Sulfolobus*, and *A. fulgidus* (Table 2; [44]). Of further note are the patterns of genes that are ubiquitous in one of the major branches of archaea, namely Euryarchaeota or Crenarchaeota, but are missing from the other branch. While quantitatively small, the set of euryarchaea-specific genes includes those for several crucial cellular functions, such as the two subunits of DNA polymerase II and the FtsZ GTPase that is required for cell division in Euryarchaeota and bacteria but missing from Crenarchaeota and eukaryotes.

Phyletic patterns can be used for interesting and potentially useful forays into functional genomics - more specifically for the identification of the genomic cognates of particular phenotypes. The most dramatic phenotypic characteristic of archaea is hyperthermophily, and attempts have been made to use the phyletic pattern approach to identify a gene set typical of hyperthermophiles. Strikingly, there is only one COG that is represented in all hyperthermophiles (both bacteria and archaea) but not in any other sequenced genomes, the reverse gyrase ([56]; COG1110). Reverse gyrase consists of a topoisomerase and a helicase domain and functions to introduce negative supercoiling into DNA; this activity is apparently required for DNA replication and gene expression at extreme high temperatures [57]. But 'clean' phyletic patterns that have an unequivocal association with a given phenotype are an exception rather than the rule, so flexible pattern selection approaches have been employed. Our recent analysis of phyletic patterns enriched in archaeal and bacterial hyperthermophiles yielded around 60 COGs potentially related to this phenotype [58]. About one quarter of these COGs encode parts of a predicted DNA repair system that is largely characteristic of thermophiles ([59] and see below). The remaining COGs in this set suggest the existence of a transcriptional regulator that might be involved in adaptation to hyperthermal environments, and a distinct class of enzymes, the S-adenosyl methionine (SAM)-radical enzymes, whose chemistry is likely to be particularly efficient under these conditions [58]. Finally, a substantial number of COGs are specific for methanogens or shared by the methanogens and *A. fulgidus* (Table 2 and [45]). Many of these include known or predicted enzymes involved in methanogenesis and associated metabolic pathways [45,60]; others remain to be characterized and are likely to encode additional components of these pathways.

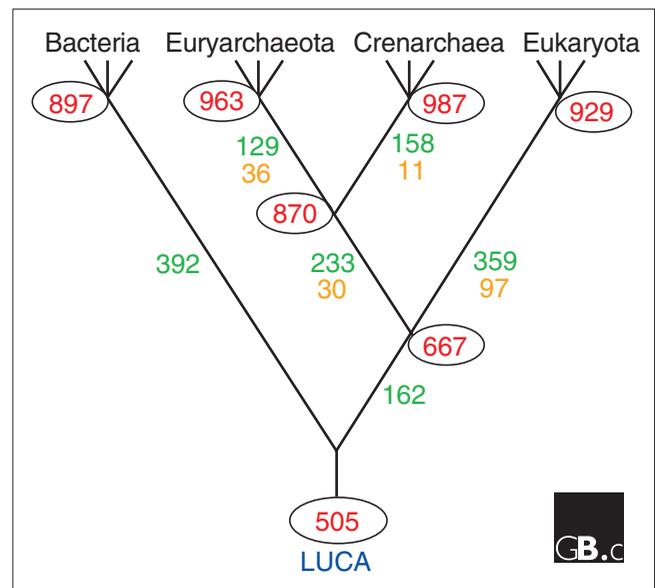
Further functional and evolutionary information can be extracted from complementary phyletic patterns, which are the signature of non-orthologous gene displacement [26,61]. Although the complementarity is, most often, only partially due to redundancy in some species, several cases of near-perfect complementarity among archaea are notable, such as the two classes of unrelated lysyl-tRNA synthetases [62,63], and two forms of thymidylate synthase that are also unrelated to each other [61,64]. Below, when discussing functional genomics of the archaea, we return to the use of conserved and complementary phyletic patterns for functional prediction.

### Genome-wide phylogeny of archaea and reconstruction of archaeal ancestors

Comparative genomics nowadays includes a new variety of phylogenetic analysis, which for short has been dubbed genome-tree construction. Under this approach, phylogenetic trees are built not from the sequences of a single gene (such as an rRNA) but from concatenated sequences of

multiple genes (proteins), from other, integral measures of the evolutionary distance between genomes (for example, the median of the distribution of evolutionary rates between orthologs), or from non-sequence-based measures such as the similarity of gene repertoire and gene orders [65]. Generally, it appears that trees produced from concatenated alignments of gene products that are not particularly prone to HGT yield the best resolution [66-68]. All genome-tree analyses unequivocally supported the monophyly of archaea and the monophyly of Crenarchaeota. Beyond that, however, the genome-tree topology is not necessarily compatible with that of rRNA-based trees. Thus, genome-tree analysis cast doubt on the bifurcation of Euryarchaeota and Crenarchaeota being the first split in archaeal evolution; in some of these analyses, *Halobacterium* and *Thermoplasma* branch off first, suggesting that Crenarchaeota are a highly derived lineage that evolved from within Euryarchaeota [66]. The same versions of genome-trees strongly suggest monophyly of methanogens, which is compatible with their distinct gene repertoire and life style [45]; but alternative trees constructed from concatenated multiple alignments of a different assortment of translation machinery components support the original divergence of Crenarchaeota and Euryarchaeota but reject the monophyly of methanogens [21,69]. It appears that a robust phylogeny of archaea will require many additional genome sequences and perhaps also further refinement of phylogenetic methods dealing with long branches and with large amounts of data. The reconstruction of the best approximation of archaeal phylogeny is of interest not so much in and of itself, but more in terms of clarifying the tempo and mode of evolution of this remarkable group of organisms. A definitive tree topology will help answer fundamental questions, such as whether methanogenesis evolved only once or several times, whether the role of histones in chromatin formation is ancestral or derived in the archaeo-eukaryotic lineage, and even the exact evolutionary relationship between archaea and eukaryotes.

Phylogenetic trees can also be employed for reconstruction of the gene sets of ancestral life forms. Given a species tree topology and phyletic patterns of the maximum possible number of orthologous gene sets (or COGs), the most parsimonious evolutionary scenario, which includes the minimum possible number of elementary events, can be reconstructed using various parsimony algorithms [70,71]. The elementary events included in this type of analysis are gene gain and gene loss. Gene gain in a given lineage may occur either as emergence of new genes (COGs), primarily via duplication with subsequent radical divergence, or as HGT from other lineages. The relative likelihood of gene loss and gene gain (the gain penalty) substantially affects the reconstructed evolutionary scenario and the gene composition of the reconstructed ancestral genomes - but this parameter is a major unknown. Nevertheless, examination of the gene sets for the last universal common ancestor (LUCA) derived with



**Figure 3**

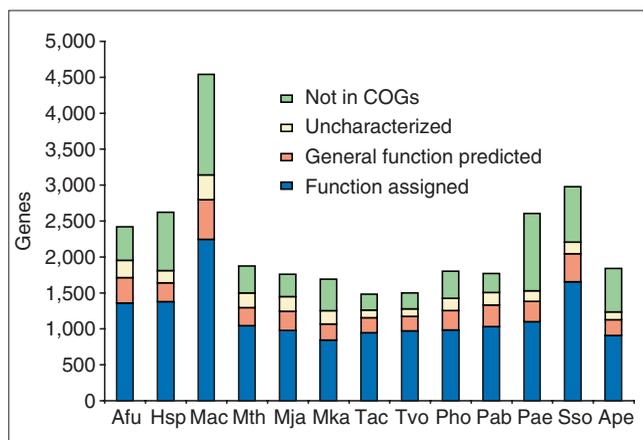
The most parsimonious scenario for the evolution of the main lineages of life. The red numbers in ovals near the internal nodes show the size of the reconstructed gene sets of the respective ancestral forms. Green numbers show gene gains and brown numbers gene losses assigned to each of the branches in the tree. LUCA, last universal common ancestor.

different gain penalties showed, perhaps rather unexpectedly, that the assumption of equal probabilities of gains and losses (a gain penalty of 1) yields a reasonable reconstruction of the main functional systems of the cell [71].

We therefore applied our version of the weighted parsimony algorithm [70], with that assumption, to the updated set of bacterial, archaeal and eukaryotic genomes (also assuming the dichotomy of Euryarchaeota and Crenarchaeota suggested by rRNA trees and some of the genome-trees) and the results are schematically shown in Figure 3 (see also additional data file). This reconstruction suggests that the common ancestor of archaea could have had around 900 genes, with substantial gene gain but only minimal gene loss compared to the more ancient common ancestor of the archaeo-eukaryotic lineage. Obviously, the conserved core of the pan-archaeal genes is a subset of the reconstructed ancestral gene set, but it seems striking that approximately two thirds of the ancestral genes have been lost from at least one of the sequenced archaeal genomes (Figure 3).

### From genome comparisons to functional and structural genomics of the archaea

In the era of comparative genomics, experimental studies on a genomic scale lag woefully behind computational studies. The great majority of the genes in most species will never be studied experimentally, and our understanding of the biochemistry and physiology of the respective organisms therefore



**Figure 4**  
Functional breakdown of genes in each of the sequenced archaeal genomes. The data are from COGs; species name abbreviations are as in Table 1.

depends on the transfer of information from functionally characterized orthologs [26,72]. For both bacteria and eukaryotes, such transfer is facilitated by the availability of a vast body of experimental data on model organisms, such as *Escherichia coli*, *Bacillus subtilis*, the yeast *Saccharomyces cerevisiae* or the fruit fly *Drosophila melanogaster*. The situation is quite different for archaea because, some genetic studies of mesophilic archaeal species notwithstanding [73], there is, so far, no satisfactory model system; this results primarily from the fact that most of these organisms grow slowly and are hard to cultivate. The functions of most of the archaeal genes have therefore been predicted by sequence analysis. Moreover, on many occasions the similarity between an archaeal protein and its functionally characterized homolog is so low that computational methods for sequence analysis have to be extended to the limit of their power.

A substantial fraction of the functional predictions for archaeal proteins appear 'trivial' in the sense that the respective proteins are highly conserved orthologs of well-characterized proteins from model organisms and, for all practical purposes, the validity of the prediction is beyond reasonable doubt (which is not to say that there are no important details of the functions of these proteins that can be uncovered only by experiment). For many other proteins, however, the prediction remains only a pointer to the probable biochemical function while the biology remains a mystery. A rough breakdown of the state of functional characterization of several archaea with sequenced genomes is given in Figure 4. The substantial fraction of genes for which only general, typically biochemical, prediction is available, is testimony to the current limited understanding of archaeal biology. Moreover, even some of the more definitive predictions only serve to emphasize the biological differences between archaea and the bacterial or eukaryotic

models from which the predictions are inferred (Table 3). A good example is the archaeal ortholog of the bacterial DNA primase (DnaG), which is a highly conserved protein present in all archaea [24]. The discovery of a predicted bacterial-type primase in archaea was unexpected, given that the archaeal replication system is orthologous to that of eukaryotes and, in particular, archaea encode the two subunits of the eukaryotic-type primase (COG1467 and COG2219; it should be noted parenthetically that detection of the large primase subunit itself required extremely careful sequence analysis due to the low similarity to the eukaryotic ortholog [22]). Given that the niche of the replicative primase seems to be occupied by the eukaryotic-type enzyme [74,75], the DnaG ortholog is likely to have a critical role in repair, but beyond this general idea its function has yet to be determined by direct experimentation; such experiments have the potential to reveal completely new repair systems and pathways. Other proteins implicated in repair as a result of exhaustive sequence analysis, such as the putative nucleases encoded by COG1833 and COG1628 (Table 3), illustrate the same point: the biochemical activities are predicted but the biology remains to be investigated experimentally.

Some of the other functional predictions inferred from sequence analysis directly help filling glaring gaps in otherwise well-characterized pathways of archaeal metabolism. A good example of such focused prediction is the identification of an archaeal fructose-1,6-bisphosphate aldolase, an indispensable glycolytic enzyme, which was first predicted computationally to be a member of the DhnA family of aldolases by our group [76] and subsequently identified experimentally [77]. In the same vein, during work for this article, we predicted the missing archaeal aconitase, an essential enzyme of the tricarboxylic acid cycle (Table 3; K.S.M. and E.V.K., unpublished observations).

The identities of a considerable number of proteins responsible for essential functions in archaea remain a mystery. Perhaps the most notable case is the missing cysteinyl-tRNA synthetase of thermophilic methanogens. Cysteine is incorporated into the proteins of these organisms as readily as in any others, but they lack an ortholog of cysteinyl-tRNA synthetase. Two different solutions for this paradox have been proposed, one involving an uncharacterized protein that has been proposed to be a 'third class' of aminoacyl-tRNA synthetases [78], and the other based on the apparent ability of the archaeal prolyl-tRNA synthetase to couple tRNA<sup>Cys</sup> with cysteine [79]. The first hypothesis has been refuted by our group upon more detailed sequence analysis [80], however, and the second did not seem to be compatible with subsequent structural studies [81]. The real cysteinyl-tRNA synthetase of methanogens seems still to be hiding among uncharacterized proteins. Gaping holes also remain in archaeal pathways of isoleucine biosynthesis [82], heme biosynthesis [83], biotin biosynthesis [26], and several others.

**Table 3****Examples of computational and experimental discovery of unexpected functions in archaea**

| COG numbers [37,38]  | Function and comments   | References  |
|--|---|---|
| <b>Computational predictions</b>   |   |   |
| 0012, 1325, 1603, 1369, 0638, 1500, 1097, 689, 2123, 1996, 2136, 2892, 0618, 1782, 1096, 3286, 1761 and more           | Archaeal exosome. Orthologs of eukaryotic exosome subunits form the largest conserved superoperon in archaea, after the ribosomal superoperon, suggesting the existence of a physical complex | [88]  |
| 1769, 1336, 3337, 1583, 1367, 1604, 1517, 1857, 1688, 1203, 1468, 1518, 2254, 1343, 1353, 1421, 1337, 1567, 1332, 4343 | DNA repair system represented primarily in thermophiles   | [59]  |
| 0358   | Bacterial-type DNA primase (DnaG orthologs)   | [24]  |
| 1311   | Small subunit of euryarchaeal DNA polymerase II, predicted PHP family phosphohydrolase (probably phosphatase); eukaryotic homologs appear to be inactivated                                   | [123]   |
| 1833   | Uri superfamily endonuclease  | [136]   |
| 1628   | Endonuclease V homologs   | K.S.M. and E.V.K., unpublished observations           |
| 1679,1786  | Aconitase catalytic core and an interacting 'swiveling domain'  | K.S.M. and E.V.K., unpublished observations           |
| 1711   | Possible subunit of the DNA replication machinery   | K.S.M. and E.V.K., unpublished observations           |
| 1310   | Zn <sup>2+</sup> -dependent hydrolase homologous to the eukaryotic ubiquitin isopeptidase contained in the proteasome and COP9 signalosome  | [137,138]   |
| <b>Computational predictions validated by experiments</b>  |   |   |
| 1708   | 'Minimal' nucleotidyltransferases   | [100,139]   |
| 1830   | Fructose-1,6-bisphosphate aldolases (DhnA family)   | [76,77]   |
| 1351   | Thymidylate synthase  | [61,64]   |
| 1685   | Shikimate kinase (predicted on the basis of operon organization)  | [140]   |
| 3635   | Phosphoglycerate mutase   | [24,141]  |
| <b>Experimental discovery of unexpected protein functions in archaea</b>   |   |   |
| 1384   | Class I lysyl-tRNA synthetase   | [62]  |
| 1933   | DNA polymerase II   | [104]   |
| 1980   | Fructose 1,6-bisphosphatase   | [142]   |
| 1630   | NurA, a novel 5'-3' nuclease encoded next to Rad50 and MreI I orthologs; present in all sequenced archaeal genomes and some bacteria  | [143] and K.S.M. and E.V.K., unpublished observations |
| 1812   | S-adenosylmethionine synthetase, was identified by mass tags  | [144]   |
| 1591   | Holliday junction resolvase   | [101]   |
| 1581   | Alba, a major DNA-binding chromatin protein in Crenarchaeota  | [106]   |
| 1945   | Pyruvoyl-dependent arginine decarboxylase (PvIArgDC), involved in polyamine biosynthesis  | [145]   |

Beyond straightforward (even if highly sensitive) sequence analysis, a powerful approach to the prediction of functions involves analysis of various forms of genomic context, or establishing 'guilt by association' [26,84-87]. The associations employed to infer gene functions may be manifest at different levels, including the phyletic patterns discussed

above, juxtaposition of domains in multidomain proteins, clustering of genes in (predicted) operons, co-expression, and protein-protein interaction. The last two of these types of data, obtained through transcriptomic and proteomic efforts, are becoming increasingly important in the functional genomics of eukaryotes and, to a somewhat lesser

extent, bacteria, but are so far unavailable for archaea. The main type of context information in archaea has therefore been obtained by analyzing conserved elements of gene order and multidomain proteins. Only a relatively small fraction (10-15%) of each archaeal genome is covered by evolutionarily conserved gene strings that can be predicted to form operons [87]. Nevertheless, by comparing gene orders in multiple genomes, partially conserved gene neighborhoods can be reconstructed and examination of some of these leads to predictions of functional systems whose existence has not previously been suspected (Table 3).

The most notable illustrations of this approach (both from our own group) are the prediction of the archaeal exosome [88] and a potential new repair system typical of archaeal and bacterial thermophiles [59]. The eukaryotic exosome is a multisubunit complex that consists of RNases, helicases and RNA-binding proteins and is involved in the exonucleolytic degradation of various classes of RNA [89-91]. During comparative analysis of gene order in prokaryotic genomes, it was observed that a distinct set of genes, some of which encode orthologs of eukaryotic exosome components, form a partially conserved predicted superoperon, which includes in total over 15 genes (although none of the archaeal genomes contains every one of these within the predicted superoperon). In addition to RNases and RNA-binding proteins (with an RNA helicase apparently encoded in a separate operon), the exosomal superoperon also encodes a proteasome subunit and a subunit of prefoldin, a co-translational molecular chaperone ([88] and Figure 5a). Thus, these observations point to the existence of a multifunctional macromolecular complex that could couple post-translational protein folding with regulated, ATP-dependent degradation of RNA and proteins. This complex remains to be discovered experimentally, and the potential implications for new functional and physical interactions in eukaryotes are also open to experimental study.

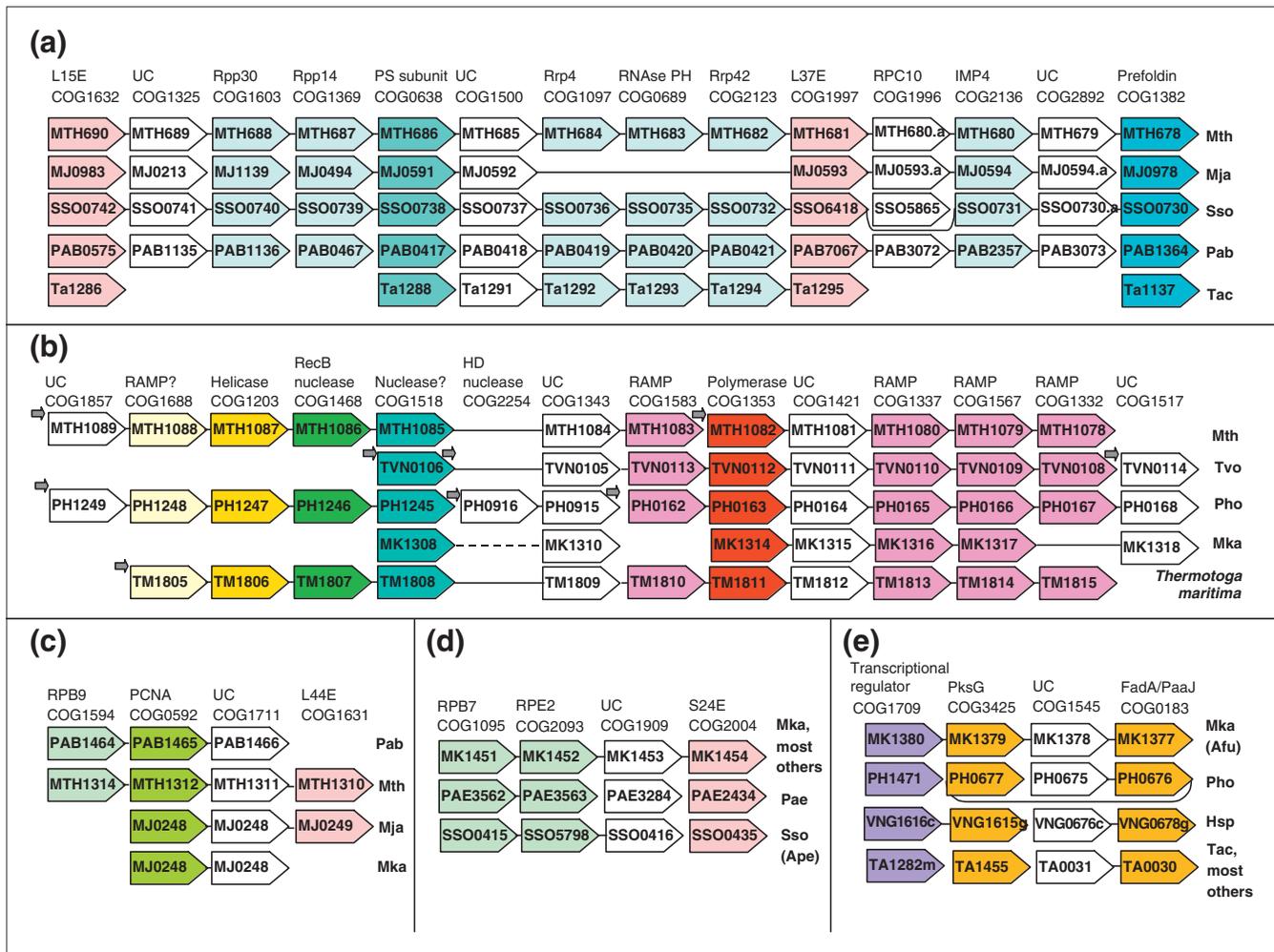
A more sophisticated comparison of gene orders, which required special algorithms for delineation of partially conserved genomic neighborhoods [92], led us to predict a distinct DNA repair system that is most prevalent in thermophiles and includes genes for a predicted novel DNA polymerase, a helicase, two nucleases and several uncharacterized genes, at least one of which could encode a novel nuclease ([59] and Figure 5b). Furthermore, this neighborhood contains multiple, diverged versions of a gene coding for a protein with a probable structural role dubbed RAMP (repair-associated mysterious protein). The proliferation of RAMP genes (Figure 5b) is an example of a potentially adaptive lineage-specific expansion of a gene family; such expansions are discussed below in greater detail.

Additional, simpler cases of functional prediction via 'guilt by association' are illustrated in Figure 5c-e. The gene for the uncharacterized protein represented by COG1711 (Figure 5c)

forms an evolutionarily highly conserved gene pair with the gene for the clamp subunit of DNA polymerase (ortholog of the eukaryotic PCNA). The orthologs of COG1711 proteins are conserved in all eukaryotes, and this protein might be an essential but still uncharacterized component of the archaeo-eukaryotic DNA replication machinery (K.S.M. and E.V.K., unpublished observations). The gene represented by uncharacterized COG1909 is squeezed between genes for RNA polymerase subunits and that for a ribosomal protein (Figure 5d). Examination of the multiple alignments that lead to this COG shows conservation of polar residues compatible with an enzymatic function (K.S.M. and E.V.K., unpublished observations). There are no readily detectable eukaryotic orthologs for this protein, which is therefore likely to be an archaea-specific enzyme with a house-keeping function.

Finally, uncharacterized COG1545 consists of genes encoding putative zinc-ribbon-containing proteins that form a stable gene pair with the gene for acetyl-CoA acetyltransferase, a central enzyme of fatty acid biosynthesis (Figure 5e). Both these genes show remarkable paralogous expansion in several archaea, probably as a result of a series of duplications of the gene doublet. It appears likely that proteins from COG1545 form a complex with acetyl-CoA acetyltransferase, with the zinc-ribbon protein regulating and/or stabilizing the enzyme. The predictions depicted in Figure 5c-e and other similar ones ([87]; and K.S.M. and E.V.K., unpublished observations) are not particularly precise, even in terms of the biochemical activity of the respective proteins. Nevertheless, guilt by association implicates each of these proteins in specific biological functions, and the evolutionary conservation of both the proteins themselves and the gene order all but proves that their functions are essential. Thus, these proteins appear to be excellent targets for experimental studies, which have the potential to reveal new facets of central cellular processes in archaea.

Comparative-genomic analysis of prokaryotes and eukaryotes points to lineage-specific expansion (proliferation) of paralogous gene families as a major means by which organisms adapt to their specific environment and lifestyle [93-95]. A number of such expansions are seen in archaea but in most cases we have, at best, only a vague understanding of the associated biology; several examples are given in Figure 6. The expansion of two groups of permeases in *Thermoplasma* and *Sulfolobus* (Figure 6a) clearly reflects the heterotrophic metabolism of the former [96] and the chemorganotrophic lifestyle of the latter [97]. The specific proliferation of ferredoxin in methanogens (Figure 6b) is also easily explained by the role of these proteins in the oxidation-reduction reactions of methanogenesis [98]. The remaining two cases in Figure 6(c,d) are much more enigmatic. The congruent proliferation of the transcription-initiation factors TFIIB and TFIID in *Halobacterium* (Figure 6c) might point to unusual aspects of transcription regulation in this archaeon but the details remain obscure. The proliferation of



**Figure 5**

Prediction of gene functions in archaea by genomic context analysis. **(a)** The superoperon coding for the predicted archaeal exosome (see [88]). **(b)** The partially conserved gene neighborhood coding for the predicted repair system found in archaeal and bacterial thermophiles (see [59] for details). **(c-e)** Predicted operons containing uncharacterized genes in the neighborhood of genes from the following COGs: COG1594, DNA-directed RNA polymerase, subunit M, and transcription elongation factor TFIIIS (RBP9); COG0592, encoding a DNA polymerase sliding clamp subunit (PCNA ortholog); COG1631, ribosomal protein L44E; COG1095, DNA-directed RNA polymerase, subunit E' (RBP7); COG2093, DNA-directed RNA polymerase, subunit E'' (RPE2); COG2004, ribosomal protein S24E; COG1709, transcriptional regulator; COG3425, 3-hydroxy-3-methylglutaryl CoA synthase (PksG); COG0183, acetyl-CoA acetyltransferase (Fad A/PaaJ orthologs). UC, uncharacterized, shown by white arrows. Species abbreviations are as in Table 1. Genes are shown not to scale and are denoted by their respective genes names (some are discussed further in the text); arrows indicate the direction of transcription. A solid line connects genes in a predicted operon. Species that have the same operon organization as the listed species are indicated in parentheses. Orthologous genes are aligned. Genes with similar general functions are shown by the same shading. Broken lines show that genes are in the same predicted operon but are not adjacent. Small arrows indicate the presence of additional functionally related genes in the same predicted operon; these genes are not shown for lack of space.

two subunits of a predicted nucleotidyltransferase in several archaea [99,100] (Figure 6d) is of special interest and might have something to do with thermal adaptation, but the actual functions and even the substrates of these enzymes remain a mystery. Other lineage-specific expansions, such as that of distinct families of predicted ATPases in *Methanocaldococcus* and *Pyrococcus*, or a specific family of RadA(RecA)-like ATPases and the UspA-family of NTP-binding proteins in several archaeal species [39], suggest the existence of unusual pathways, perhaps involved in stress response and signal

transduction, but the actual biology associated with these expansions can only be uncovered experimentally.

Archaeal comparative genomics is a young field and so far, as we have seen, largely predictive. But a few experimental studies have already been instigated as a result of comparative-genomic predictions. The discovery of the archaeal fructose-1,6-bisphosphate aldolase mentioned above [76,77] is a case in point, and several other examples of experimental validation of predictions are given in Table 3. It does not seem to

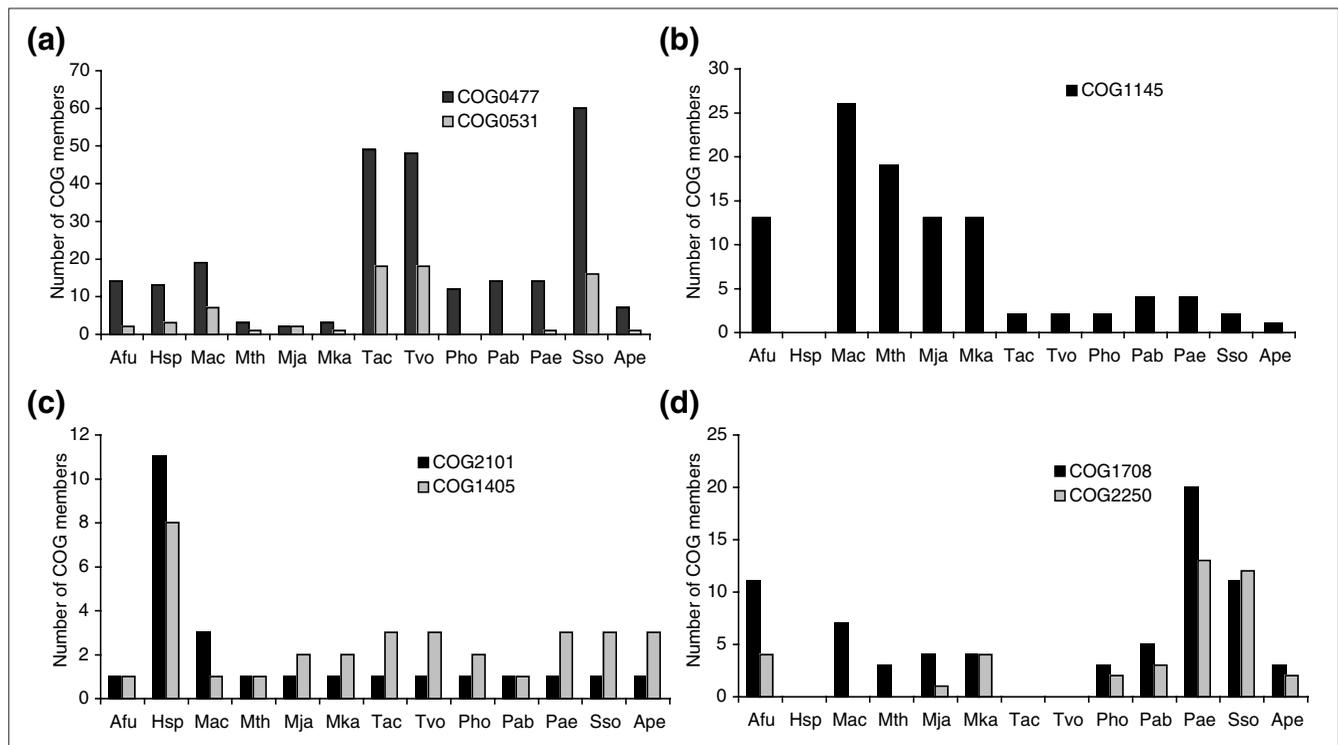
be chance that these examples all involve metabolic enzymes for which the specific reaction could be predicted precisely. Validation is likely to be much more difficult for proteins of other functional groups, such as putative repair enzymes, for which the actual substrates are harder to predict.

For some conserved archaeal proteins, functions cannot be predicted computationally despite considerable effort. Several important discoveries have been made by experimental characterization of such mysterious proteins. The most notable cases include the archaeal Holliday-junction resolvase, which is not related to its functional analog in bacteria [101-103], and DNA polymerase II, a highly conserved euryarchaeal protein that is not found outside this lineage and shows no detectable sequence similarity to any other proteins [104,105]. Additional examples of direct experimental determination of the functions of archaeal proteins that could not be predicted by computational techniques (at least not before the experiment had been reported) are given in Table 3.

Especially notable is the story of the Alba protein, a DNA-binding component of chromatin in Crenarchaeota [106,107]. As noted above, crenarchaea lack histones and in these

organisms Alba appears to be the main chromatin protein, in a striking case of non-orthologous gene displacement. But orthologs of Alba are also present in thermophilic Euryarchaeota and in some eukaryotic lineages, where its functions remain to be elucidated. The most remarkable discovery regarding Alba is the regulation of its interaction with DNA and with the chromatin-associated protein deacetylase Sir2 via lysine acetylation and deacetylation [106,108]. In eukaryotes, regulation of chromatin dynamics via acetylation and deacetylation occurs through histone tails [109]. Thus, a special case of non-orthologous gene displacement seems to have taken place whereby the regulation mechanism is conserved but the actual substrates are different in archaea and eukaryotes. To add an extra twist to the story, *Thermoplasma* lacks both histones and Alba but has the bacterial DNA-binding protein HU, pointing to three distinct solutions to the problem of chromatin organization in archaea [107].

The last subject we have to briefly touch upon is structural genomics of the archaea. The ultimate goal of the structural genomics enterprise is determining the three-dimensional structure for all proteins, or at least for all sufficiently different proteins encoded in the genomes of diverse life forms [110]. This goal is far from being reached, and targets for



**Figure 6**  
Lineage-specific expansions of paralogous gene families in archaea. The vertical axis shows the number of members of the indicated COGs. (a) COG0477, permeases of the major facilitator superfamily; COG0531, amino-acid transporters. (b) COG1145, ferredoxin. (c) COG2101, TATA-box binding protein (TBP), a component of transcription initiation factors TFIID and TFIIB; COG1405, Brf1 subunit of transcription-initiation factor TFIIB and transcription-initiation factor TFIIB. (d) COG1708, 'minimal' nucleotidyltransferase catalytic subunit; COG2250, 'minimal' nucleotidyltransferase accessory subunit. Species abbreviations are as in Table 1.

structural determination have been prioritized by different researchers on the basis of different principles, from nearly random choice to relatively elaborate strategies, including the use of the COG database [111-115]. The development of structural genomics so far has been a mixture of success, when informative and interesting structures have been solved, and mild disappointment in cases when the structure determination did not seem to shed any light on a protein's function. Structural genomics could be particularly important in the case of archaea, for which a miniscule number of structures had been solved prior to the launch of structural genomic initiatives, and in which proteins often show low similarity to bacterial or eukaryotic homologs, making homology modeling difficult.

Notable developments that illustrate both the benefits and the pitfalls of structural genomics, are the concerted effort on 'structural proteomics' of *Methanothermobacter thermoautotrophicus* [116] and a similar project on *M. jannaschii* [117]. The elucidation of the structure of the *M. jannaschii* protein MJ0577 [117] is an excellent case for the power of structural genomics. Analysis of this structure and accompanying biochemical experiments revealed a distinct nucleotide-binding domain that is distantly related to the catalytic domains of class I aminoacyl-tRNA synthetases and belongs to the so-called HUP fold of nucleotide-binding domains [118]. Together with comprehensive sequence analysis, the determination of this structure provided the structural, functional and evolutionary context for the UspA protein family, which is specifically expanded in archaea [39]. The exact function(s) of these proteins remains unknown but, in this case, structural genomics ensured a substantial functional insight. On several other occasions, however, determination of the structures of archaeal proteins has failed to provide clear functional clues; these remain structures in search of a function.

### What's around the corner?

The first sequenced archaeal genome was a veritable *terra incognita*. Six years after that sequence appeared, the archaeal genomescape looks quite different. The principal landmarks have been mapped and now, when a new archaeal genome is released, we largely know what to expect from it. Computational approaches to comparative genomics, combining in-depth sequence and structure comparison with genome context analysis, have led to the reconstruction of the central functional systems of archaeal cells. But these approaches have also produced numerous isolated predictions of biochemical activities of archaeal proteins that remain to be fitted into a general picture, and this can be done only through 'wet' experiments, although new genome sequences will substantially help by enriching the genomic context. A shrinking but still notable set of archaeal genes includes those that encode highly conserved proteins without any clue to function; solving these mysteries has the potential

to bring out truly new biology. Furthermore, in this article we have not even touched upon important aspects of archaeal genomics, such as the in-depth studies of the translation system, which have revealed several highly unusual, remarkable mechanisms and enzymatic systems [63,119] or the identification of regulatory sites in DNA and patterns of transcription regulation [120,121]. The latter avenue of research is still in its infancy but will certainly grow in scale once more archaeal genomes, and in particular closely related ones, are sequenced.

Because of the lack of established model systems for archaeal experimental biology and the resulting difficulty with large-scale experimentation, clues from genome comparison are even more crucial for archaeal functional genomics than they are in the case of bacteria or eukaryotes. So far, the input of comparative genomics into actual experiments has been less prominent than we would hope. Simply put, it is not often that experimenters rush to test predictions produced by *in silico* genome comparison and, furthermore, it is even rarer that targets for functional characterization are carefully prioritized on the basis of how unusual and fundamental the predictions are. As discussed above, however, the few cases when such tests have been performed are encouraging. It is our hope that the future belongs to a much tighter integration of comparative, structural and functional genomics.

Beyond functional studies, archaeal genomics is fundamental to our understanding of two critical transitions in the evolution of life. The first is the primary split between the bacterial and archaeo-eukaryotic lineages, which might have involved the origin of the DNA-replication machinery and of the large, double-stranded DNA genomes themselves [22,23], and the second is the origin of eukaryotes [122]. With regard to the latter problem, archaea are a particularly valuable source of information because, on many occasions, they seem to have retained primitive traits while eukaryotes have undergone major changes. A characteristic example is the small DNA polymerase subunit, which has all the hallmarks of an active phosphatase in archaea, but not in eukaryotes, in which the phosphatase activity is predicted to be inactivated [123]. Indubitably, archaea resemble the common ancestor of the archaeo-eukaryotic line of descent more closely than eukaryotes do, so archaeal genomics is our best chance to reconstruct this critical intermediate in the evolution of life. We are confident that comparative archaeogenomics has a bright future, with major progress in both the functional and the evolutionary avenues of research expected within the next few years.

### Additional data file

The list of genes in the reconstructed gene set of the last common ancestor of archaea is available with the complete version of this article, online.

## Acknowledgements

We thank Boris Mirkin for producing the data used for Figure 3 and Stephen Bell, Michael Galperin, Dieter Söll and Yuri Wolf for useful discussions.

## References

- Woese CR, Fox GE: **Phylogenetic structure of the prokaryotic domain: the primary kingdoms.** *Proc Natl Acad Sci USA* 1977, **74**:5088-5090.
- Fox GE, Stackebrandt E, Hespell RB, Gibson J, Maniloff J, Dyer TA, Wolfe RS, Balch WE, Tanner RS, Magrum LJ, et al.: **The phylogeny of prokaryotes.** *Science* 1980, **209**:457-463.
- Woese CR, Kandler O, Wheelis ML: **Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya.** *Proc Natl Acad Sci USA* 1990, **87**:4576-4579.
- Woese CR, Gupta R: **Are archaeobacteria merely derived 'prokaryotes'?** *Nature* 1981, **289**:95-96.
- Mayr E: **Two empires or three?** *Proc Natl Acad Sci USA* 1998, **95**:9720-9723.
- Woese CR: **Default taxonomy: Ernst Mayr's view of the microbial world.** *Proc Natl Acad Sci USA* 1998, **95**:11043-11046.
- Gupta RS: **Life's third domain (Archaea): an established fact or an endangered paradigm?** *Theor Popul Biol* 1998, **54**:91-104.
- Kushner DJ: **Lysis and dissolution of cells and envelopes of an extremely halophilic bacterium.** *J Bacteriol* 1964, **87**:1147-1156.
- Langworthy TA, Smith PF, Mayberry WR: **Lipids of *Thermoplasma acidophilum*.** *J Bacteriol* 1972, **112**:1193-1200.
- Brock TD, Brock KM, Belly RT, Weiss RL: ***Sulfolobus*: a new genus of sulfur-oxidizing bacteria living at low pH and high temperature.** *Arch Mikrobiol* 1972, **84**:54-68.
- Stetter KO: **Extremophiles and their adaptation to hot environments.** *FEBS Lett* 1999, **452**:22-25.
- Seegerer AH, Burggraf S, Fiala G, Huber G, Huber R, Pley U, Stetter KO: **Life in hot springs and hydrothermal vents.** *Orig Life Evol Biosph* 1993, **23**:77-90.
- DeLong EF: **A phylogenetic perspective on hyperthermophilic microorganisms.** *Methods Enzymol* 2001, **330**:3-11.
- DeLong EF, Pace NR: **Environmental diversity of bacteria and archaea.** *Syst Biol* 2001, **50**:470-478.
- Hanford MJ, Peebles TL: **Archaeal tetraether lipids: unique structures and applications.** *Appl Biochem Biotechnol* 2002, **97**:45-62.
- Engelhardt H, Peters J: **Structural research on surface layers: a focus on stability, surface layer homology domains, and surface layer-cell wall interactions.** *J Struct Biol* 1998, **124**:276-302.
- Edgell DR, Doolittle WF: **Archaea and the origin(s) of DNA replication proteins.** *Cell* 1997, **89**:995-998.
- Sandman K, Pereira SL, Reeve JN: **Diversity of prokaryotic chromosomal proteins and the origin of the nucleosome.** *Cell Mol Life Sci* 1998, **54**:1350-1364.
- Sandman K, Bailey KA, Pereira SL, Soares D, Li WT, Reeve JN: **Archaeal histones and nucleosomes.** *Methods Enzymol* 2001, **334**:116-129.
- Lecompte O, Ripp R, Thierry JC, Moras D, Poch O: **Comparative analysis of ribosomal proteins in complete genomes: an example of reductive evolution at the domain scale.** *Nucleic Acids Res* 2002, **30**:5382-5390.
- Forterre P, Brochier C, Philippe H: **Evolution of the Archaea.** *Theor Popul Biol* 2002, **61**:409-422.
- Leipe DD, Aravind L, Koonin EV: **Did DNA replication evolve twice independently?** *Nucleic Acids Res* 1999, **27**:3389-3401.
- Forterre P: **The origin of DNA genomes and DNA replication proteins.** *Curr Opin Microbiol* 2002, **5**:525-532.
- Koonin EV, Mushegian AR, Galperin MY, Walker DR: **Comparison of archaeal and bacterial genomes: computer analysis of protein sequences predicts novel functions and suggests a chimeric origin for the archaea.** *Mol Microbiol* 1997, **25**:619-637.
- Jain R, Rivera MC, Lake JA: **Horizontal gene transfer among genomes: the complexity hypothesis.** *Proc Natl Acad Sci USA* 1999, **96**:3801-3806.
- Koonin EV, Galperin MY: *Sequence - Evolution - Function. Computational Approaches in Comparative Genomics.* New York: Kluwer Academic; 2002.
- Bult CJ, White O, Olsen GJ, Zhou L, Fleischmann RD, Sutton GG, Blake JA, FitzGerald LM, Clayton RA, Gocayne JD, et al.: **Complete genome sequence of the methanogenic archaeon, *Methanococcus jannaschii*.** *Science* 1996, **273**:1058-1073.
- Pace NR: **A molecular view of microbial diversity and the biosphere.** *Science* 1997, **276**:734-740.
- Huber H, Hohn MJ, Rachel R, Fuchs T, Wimmer VC, Stetter KO: **A new phylum of Archaea represented by a nanosized hyperthermophilic symbiont.** *Nature* 2002, **417**:63-67.
- Huber H, Hohn MJ, Stetter KO, Rachel R: **The phylum Nanoarchaeota: present knowledge and future perspectives of a unique form of life.** *Res Microbiol* 2003, **154**:165-171.
- Tatusov RL, Koonin EV, Lipman DJ: **A genomic perspective on protein families.** *Science* 1997, **278**:631-637.
- Fitch WM: **Distinguishing homologous from analogous proteins.** *Syst Zool* 1970, **19**:99-113.
- Fitch WM: **Homology: a personal view on some of the problems.** *Trends Genet* 2000, **16**:227-231.
- Sonnhammer EL, Koonin EV: **Orthology, paralogy and proposed classification for paralog subtypes.** *Trends Genet* 2002, **18**:619-620.
- Gaasterland T, Ragan MA: **Microbial genescapes: phyletic and functional patterns of ORF distribution among prokaryotes.** *Microb Comp Genomics* 1998, **3**:199-217.
- Pellegrini M, Marcotte EM, Thompson MJ, Eisenberg D, Yeates TO: **Assigning protein functions by comparative genome analysis: protein phylogenetic profiles.** *Proc Natl Acad Sci USA* 1999, **96**:4285-4288.
- Tatusov RL, Natale DA, Garkavtsev IV, Tatusova TA, Shankavaram UT, Rao BS, Kiryutin B, Galperin MY, Fedorova ND, Koonin EV: **The COG database: new developments in phylogenetic classification of proteins from complete genomes.** *Nucleic Acids Res* 2001, **29**:22-28.
- Prokaryotic COGs project phyletic pattern search** [<http://www.ncbi.nlm.nih.gov/COG/new/release/phylox.cgi>]
- Makarova KS, Aravind L, Galperin MY, Grishin NV, Tatusov RL, Wolf YI, Koonin EV: **Comparative genomics of the Archaea (Euryarchaeota): evolution of conserved protein families, the stable core, and the variable shell.** *Genome Res* 1999, **9**:608-628.
- Koonin EV, Mushegian AR, Bork P: **Non-orthologous gene displacement.** *Trends Genet* 1996, **12**:334-336.
- Doolittle WF, Logsdon JM, Jr.: **Archaeal genomics: do archaea have a mixed heritage?** *Curr Biol* 1998, **8**:R209-R211.
- Doolittle WF: **Phylogenetic classification and the universal tree.** *Science* 1999, **284**:2124-2129.
- Doolittle WF: **Lateral genomics.** *Trends Cell Biol* 1999, **9**:M5-M8.
- Koonin EV, Makarova KS, Aravind L: **Horizontal gene transfer in prokaryotes - quantification and classification.** *Annu Rev Microbiol* 2001, **55**:709-42.
- Slesarev AI, Mezhevaya KV, Makarova KS, Polushin NN, Shcherbinina OV, Shakhova VV, Belova GI, Aravind L, Natale DA, Rogozin IB, et al.: **The complete genome of hyperthermophile *Methanopyrus kandleri* AV19 and monophyly of archaeal methanogens.** *Proc Natl Acad Sci USA* 2002, **99**:4644-4649.
- Aravind L, Tatusov RL, Wolf YI, Walker DR, Koonin EV: **Evidence for massive gene exchange between archaeal and bacterial hyperthermophiles.** *Trends Genet* 1998, **14**:442-444.
- Nelson KE, Clayton RA, Gill SR, Gwinn ML, Dodson RJ, Haft DH, Hickey EK, Peterson JD, Nelson WC, Ketchum KA, et al.: **Evidence for lateral gene transfer between Archaea and bacteria from genome sequence of *Thermotoga maritima*.** *Nature* 1999, **399**:323-329.
- Brown JR: **Ancient horizontal gene transfer.** *Nat Rev Genet* 2003, **4**:121-132.
- Kyrpides NC, Olsen GJ: **Archaeal and bacterial hyperthermophiles: horizontal gene exchange or common ancestry?** *Trends Genet* 1999, **15**:298-299.
- Aravind L, Tatusov RL, Wolf YI, Walker DR, Koonin EV: **Reply. Archaeal and bacterial hyperthermophiles: horizontal gene exchange or common ancestry?** *Trends Genet* 1999, **15**:299-300.
- Bao Q, Tian Y, Li W, Xu Z, Xuan Z, Hu S, Dong W, Yang J, Chen Y, Xue Y, et al.: **A complete sequence of the *T. tengcongensis* genome.** *Genome Res* 2002, **12**:689-700.
- Thermoanaerobacter tengcongensis* proteins** [<http://www.ncbi.nlm.nih.gov/sutils/taxik.cgi?gi=237>]
- Nesbo CL, Boucher Y, Doolittle WF: **Defining the core of non-transferable prokaryotic genes: the euryarchaeal core.** *J Mol Evol* 2001, **53**:340-350.

54. Deppenmeier U, Johann A, Hartsch T, Merkl R, Schmitz RA, Martinez-Arias R, Henne A, Wiezer A, Baumer S, Jacobi C, et al.: **The genome of *Methanosarcina mazei*: evidence for lateral gene transfer between bacteria and archaea.** *J Mol Microbiol Biotechnol* 2002, **4**:453-461.
55. Galagan JE, Nusbaum C, Roy A, Endrizzi MG, Macdonald P, FitzHugh W, Calvo S, Engels R, Smirnov S, Atnoor D, et al.: **The genome of *M. acetivorans* reveals extensive metabolic and physiological diversity.** *Genome Res* 2002, **12**:532-542.
56. Forterre P: **A hot story from comparative genomics: reverse gyrase is the only hyperthermophile-specific protein.** *Trends Genet* 2002, **18**:236-237.
57. Forterre P, Bergerat A, Lopez-Garcia P: **The unique DNA topology and DNA topoisomerases of hyperthermophilic archaea.** *FEMS Microbiol Rev* 1996, **18**:237-248.
58. Makarova KS, Wolf YI, Koonin EV: **Potential genomic determinants of hyperthermophily.** *Trends Genet* 2003, **19**:172-176.
59. Makarova KS, Aravind L, Grishin NV, Rogozin IB, Koonin EV: **A DNA repair system specific for thermophilic archaea and bacteria predicted by genomic context analysis.** *Nucleic Acids Res* 2002, **30**:482-496.
60. White RH: **Biosynthesis of the methanogenic cofactors.** *Vitam Horm* 2001, **61**:299-337.
61. Galperin MY, Koonin EV: **Who's your neighbor? New computational approaches for functional genomics.** *Nat Biotechnol* 2000, **18**:609-613.
62. Ibba M, Morgan S, Curnow AW, Pridmore DR, Vothknecht UC, Gardner W, Lin W, Woese CR, Soll D: **A euryarchaeal lysyl-tRNA synthetase: resemblance to class I synthetases.** *Science* 1997, **278**:1119-1122.
63. Praetorius-Ibba M, Ibba M: **Aminoacyl-tRNA synthesis in archaea: different but not unique.** *Mol Microbiol* 2003, **48**:631-637.
64. Myllykallio H, Lipowski G, Leduc D, Filee J, Forterre P, Liebl U: **An alternative flavin-dependent mechanism for thymidylate synthesis.** *Science* 2002, **297**:105-107.
65. Wolf YI, Rogozin IB, Grishin NV, Koonin EV: **Genome trees and the tree of life.** *Trends Genet* 2002, **18**:472-479.
66. Wolf YI, Rogozin IB, Grishin NV, Tatusov RL, Koonin EV: **Genome trees constructed using five different approaches suggest new major bacterial clades.** *BMC Evol Biol* 2001, **1**:8.
67. Clarke GD, Beiko RG, Ragan MA, Charlebois RL: **Inferring genome trees by using a filter to eliminate phylogenetically discordant sequences and a distance matrix based on mean normalized BLASTP scores.** *J Bacteriol* 2002, **184**:2072-2080.
68. Korbelt JO, Snel B, Huynen MA, Bork P: **SHOT: a web server for the construction of genome phylogenies.** *Trends Genet* 2002, **18**:158-162.
69. Matte-Tailliez O, Brochier C, Forterre P, Philippe H: **Archaeal phylogeny based on ribosomal proteins.** *Mol Biol Evol* 2002, **19**:631-639.
70. Snel B, Bork P, Huynen MA: **Genomes in flux: the evolution of archaeal and proteobacterial gene content.** *Genome Res* 2002, **12**:17-25.
71. Mirkin BG, Fenner TI, Galperin MY, Koonin EV: **Algorithms for computing parsimonious evolutionary scenarios for genome evolution, the last universal common ancestor and dominance of horizontal gene transfer in the evolution of prokaryotes.** *BMC Evol Biol* 2003, **3**:2.
72. Wilson CA, Kreychman J, Gerstein M: **Assessing annotation transfer for genomics: quantifying the relations between protein sequence, structure and function through traditional and probabilistic scores.** *J Mol Biol* 2000, **297**:233-249.
73. Luo Y, Wasserfallen A: **Gene transfer systems and their applications in Archaea.** *Syst Appl Microbiol* 2001, **24**:15-25.
74. Liu L, Komori K, Ishino S, Bocquier AA, Cann IK, Kohda D, Ishino Y: **The archaeal DNA primase: biochemical characterization of the p41-p46 complex from *Pyrococcus furiosus*.** *J Biol Chem* 2001, **276**:45484-45490.
75. Bocquier AA, Liu L, Cann IK, Komori K, Kohda D, Ishino Y: **Archaeal primase: bridging the gap between RNA and DNA polymerases.** *Curr Biol* 2001, **11**:452-456.
76. Galperin MY, Aravind L, Koonin EV: **Aldolases of the Dhna family: a possible solution to the problem of pentose and hexose biosynthesis in archaea.** *FEMS Microbiol Lett* 2000, **183**:259-264.
77. Siebers B, Brinkmann H, Dorr C, Tjaden B, Lilie H, van der Oost J, Verhees CH: **Archaeal fructose-1,6-bisphosphate aldolases constitute a new family of archaeal type class I aldolase.** *J Biol Chem* 2001, **276**:28710-28718.
78. Fabrega C, Farrow MA, Mukhopadhyay B, de Crecy-Lagard V, Ortiz AR, Schimmel P: **An aminoacyl tRNA synthetase whose sequence fits into neither of the two known classes.** *Nature* 2001, **411**:110-114.
79. Stathopoulos C, Li T, Longman R, Vothknecht UC, Becker HD, Ibba M, Soll D: **One polypeptide with two aminoacyl-tRNA synthetase activities.** *Science* 2000, **287**:479-482.
80. Iyer LM, Aravind L, Bork P, Hofmann K, Mushegian AR, Zhulin IB, Koonin EV: **Quod erat demonstrandum? The mystery of experimental validation of apparently erroneous computational analyses of protein sequences.** *Genome Biol* 2001, **2**:research0051.1-0051.11
81. Kamtekar S, Kennedy WD, Wang J, Stathopoulos C, Soll D, Steitz TA: **The structural basis of cysteine aminoacylation of tRNA<sup>Pro</sup> by prolyl-tRNA synthetases.** *Proc Natl Acad Sci USA* 2003, **100**:1673-1678.
82. Xie G, Forst C, Bonner C, Jensen RA: **Significance of two distinct types of tryptophan synthase beta chain in Bacteria, Archaea and higher plants.** *Genome Biol* 2002, **3**:research0004.1-0004.13
83. Panek H, O'Brian MR: **A whole genome view of prokaryotic haem biosynthesis.** *Microbiology* 2002, **148**:2273-2282.
84. Huynen M, Snel B, Lathe W, Bork P: **Exploitation of gene context.** *Curr Opin Struct Biol* 2000, **10**:366-370.
85. Huynen M, Snel B, Lathe W 3rd, Bork P: **Predicting protein function by genomic context: quantitative evaluation and qualitative inferences.** *Genome Res* 2000, **10**:1204-1210.
86. Aravind L: **Guilt by association: contextual information in genome analysis.** *Genome Res* 2000, **10**:1074-1077.
87. Wolf YI, Rogozin IB, Kondrashov AS, Koonin EV: **Genome alignment, evolution of prokaryotic genome organization, and prediction of gene function using genomic context.** *Genome Res* 2001, **11**:356-372.
88. Koonin EV, Wolf YI, Aravind L: **Prediction of the archaeal exosome and its connections with the proteasome and the translation and transcription machineries by a comparative-genomic approach.** *Genome Res* 2001, **11**:240-252.
89. Decker CJ: **The exosome: a versatile RNA processing machine.** *Curr Biol* 1998, **8**:R238-R240.
90. van Hoof A, Parker R: **The exosome: a proteasome for RNA?** *Cell* 1999, **99**:347-350.
91. Mitchell P, Petfalski E, Shevchenko A, Mann M, Tollervey D: **The exosome: a conserved eukaryotic RNA processing complex containing multiple 3'→5' exoribonucleases.** *Cell* 1997, **91**:457-466.
92. Rogozin IB, Makarova KS, Murvai J, Czabarka E, Wolf YI, Tatusov RL, Szekely LA, Koonin EV: **Connected gene neighborhoods in prokaryotic genomes.** *Nucleic Acids Res* 2002, **30**:2212-2223.
93. Lespinet O, Wolf YI, Koonin EV, Aravind L: **The role of lineage-specific gene family expansion in the evolution of eukaryotes.** *Genome Res* 2002, **12**:1048-1059.
94. Jordan IK, Makarova KS, Spouge JL, Wolf YI, Koonin EV: **Lineage-specific gene expansions in bacterial and archaeal genomes.** *Genome Res* 2001, **11**:555-565.
95. Remm M, Storm CE, Sonnhammer EL: **Automatic clustering of orthologs and in-paralogs from pairwise species comparisons.** *J Mol Biol* 2001, **314**:1041-1052.
96. Ruepp A, Graml W, Santos-Martinez ML, Koretke KK, Volker C, Mewes HW, Frishman D, Stocker S, Lupas AN, Baumeister W: **The genome sequence of the thermoacidophilic scavenger *Thermoplasma acidophilum*.** *Nature* 2000, **407**:508-513.
97. She Q, Singh RK, Confalonieri F, Zivanovic Y, Allard G, Awayez MJ, Chan-Weiher CC, Clausen IG, Curtis BA, De Moors A, et al.: **The complete genome of the crenarchaeon *Sulfolobus solfataricus* P2.** *Proc Natl Acad Sci USA* 2001, **98**:7835-7840.
98. Deppenmeier U: **The unique biochemistry of methanogenesis.** *Prog Nucleic Acid Res Mol Biol* 2002, **71**:223-283.
99. Aravind L, Koonin EV: **DNA polymerase beta-like nucleotidyltransferase superfamily: identification of three new families, classification and evolutionary history.** *Nucleic Acids Res* 1999, **27**:1609-1618.
100. Lehmann C, Lim K, Chalamasetty VR, Krajewski W, Melamud E, Galkin A, Howard A, Kelman Z, Reddy PT, Murzin AG, Herzberg O: **The HI0073/HI0074 protein pair from *Haemophilus influenzae* is a member of a new nucleotidyltransferase family: structure, sequence analyses, and solution studies.** *Proteins* 2003, **50**:249-260.

101. Komori K, Sakae S, Shinagawa H, Morikawa K, Ishino Y: **A Holliday junction resolvase from *Pyrococcus furiosus*: functional similarity to *Escherichia coli* RuvC provides evidence for conserved mechanism of homologous recombination in Bacteria, Eukarya, and Archaea.** *Proc Natl Acad Sci USA* 1999, **96**:8873-8878.
102. Aravind L, Makarova KS, Koonin EV: **Survey and summary: Holliday junction resolvases and related nucleases: identification of new families, phyletic distribution and evolutionary trajectories.** *Nucleic Acids Res* 2000, **28**:3417-3432.
103. Daiyasu H, Komori K, Sakae S, Ishino Y, Toh H: **Hjc resolvase is a distantly related member of the type II restriction endonuclease family.** *Nucleic Acids Res* 2000, **28**:4540-4543.
104. Ishino Y, Komori K, Cann IK, Koga Y: **A novel DNA polymerase family found in Archaea.** *J Bacteriol* 1998, **180**:2232-2236.
105. Ishino Y, Ishino S: **DNA polymerases from euryarchaeota.** *Methods Enzymol* 2001, **334**:249-260.
106. Bell SD, Botting CH, Wardleworth BN, Jackson SP, White MF: **The interaction of Alba, a conserved archaeal chromatin protein, with Sir2 and its regulation by acetylation.** *Science* 2002, **296**:148-151.
107. White MF, Bell SD: **Holding it together: chromatin in the Archaea.** *Trends Genet* 2002, **18**:621-626.
108. Wardleworth BN, Russell RJ, Bell SD, Taylor GL, White MF: **Structure of Alba: an archaeal chromatin protein modulated by acetylation.** *EMBO J* 2002, **21**:4654-4662.
109. Kurdistani SK, Grunstein M: **Histone acetylation and deacetylation in yeast.** *Nat Rev Mol Cell Biol* 2003, **4**:276-284.
110. Vitkup D, Melamud E, Moutl J, Sander C: **Completeness in structural genomics.** *Nat Struct Biol* 2001, **8**:559-566.
111. Elofsson A, Sonnhammer EL: **A comparison of sequence and structure protein domain families as a basis for structural genomics.** *Bioinformatics* 1999, **15**:480-500.
112. Gaasterland T: **Structural genomics: bioinformatics in the driver's seat.** *Nat Biotechnol* 1998, **16**:625-627.
113. Brenner SE: **Target selection for structural genomics.** *Nat Struct Biol* 2000, **7 Suppl**:967-969.
114. Gerstein M: **Integrative database analysis in structural genomics.** *Nat Struct Biol* 2000, **7 Suppl**:960-963.
115. Koonin EV, Wolf YI, Aravind L: **Protein fold recognition using sequence profiles and its application in structural genomics.** *Adv Protein Chem* 2000, **54**:245-275.
116. Christendat D, Yee A, Dharamsi A, Kluger Y, Savchenko A, Cort JR, Booth V, Mackereth CD, Saridakis V, Ekiel I, et al.: **Structural proteomics of an archaeon.** *Nat Struct Biol* 2000, **7**:903-909.
117. Zarembinski TI, Hung LW, Mueller-Dieckmann HJ, Kim KK, Yokota H, Kim R, Kim SH: **Structure-based assignment of the biochemical function of a hypothetical protein: a test case of structural genomics.** *Proc Natl Acad Sci USA* 1998, **95**:15189-15193.
118. Aravind L, Anantharaman V, Koonin EV: **Monophyly of class I aminoacyl tRNA synthetase, USPA, ETFP, photolyase, and PP-ATPase nucleotide-binding domains: implications for protein evolution in the RNA.** *Proteins* 2002, **48**:1-14.
119. Woese CR: **Translation: in retrospect and prospect.** *RNA* 2001, **7**:1055-1067.
120. Gelfand MS, Koonin EV, Mironov AA: **Prediction of transcription regulatory sites in Archaea by a comparative genomic approach.** *Nucleic Acids Res* 2000, **28**:695-705.
121. Rodionov DA, Mironov AA, Gelfand MS: **Conservation of the biotin regulon and the BirA regulatory signal in eubacteria and archaea.** *Genome Res* 2002, **12**:1507-1516.
122. Dacks JB, Doolittle WF: **Reconstructing/deconstructing the earliest eukaryotes: how comparative genomics can help.** *Cell* 2001, **107**:419-425.
123. Aravind L, Koonin EV: **Phosphoesterase domains associated with DNA polymerases of diverse origins.** *Nucleic Acids Res* 1998, **26**:3746-3752.
124. Klenk HP, Clayton RA, Tomb JF, White O, Nelson KE, Ketchum KA, Dodson RJ, Gwinn M, Hickey EK, Peterson JD, et al.: **The complete genome sequence of the hyperthermophilic, sulphate-reducing archaeon *Archaeoglobus fulgidus*.** *Nature* 1997, **390**:364-370.
125. Ng WV, Kennedy SP, Mahairas GG, Berquist B, Pan M, Shukla HD, Lasky SR, Baliga NS, Thorsson V, Sbrogna J, et al.: **Genome sequence of *Halobacterium* species NRC-1.** *Proc Natl Acad Sci USA* 2000, **97**:12176-12181.
126. Smith DR, Doucette-Stamm LA, Deloughery C, Lee H, Dubois J, Aldredge T, Bashirzadeh R, Blakely D, Cook R, Gilbert K, et al.: **Complete genome sequence of *Methanobacterium thermoautotrophicum* deltaH: functional analysis and comparative genomics.** *J Bacteriol* 1997, **179**:7135-7155.
127. Kawarabayasi Y, Sawada M, Horikawa H, Haikawa Y, Hino Y, Yamamoto S, Sekine M, Baba S, Kosugi H, Hosoyama A, et al.: **Complete sequence and gene organization of the genome of a hyper-thermophilic archaeobacterium, *Pyrococcus horikoshii* OT3.** *DNA Res* 1998, **5**:55-76.
128. Cohen GN, Barbe V, Flament D, Galperin M, Heilig R, Lecompte O, Poch O, Prieur D, Querellou J, Ripp R, et al.: **An integrated analysis of the genome of the hyperthermophilic archaeon *Pyrococcus abyssi*.** *Mol Microbiol* 2003, **47**:1495-1512.
129. Robb FT, Maeder DL, Brown JR, DiRuggiero J, Stump MD, Yeh RK, Weiss RB, Dunn DM: **Genomic sequence of hyperthermophile, *Pyrococcus furiosus*: implications for physiology and enzymology.** *Methods Enzymol* 2001, **330**:134-157.
130. Kawashima T, Amano N, Koike H, Makino S, Higuchi S, Kawashima-Ohya Y, Watanabe K, Yamazaki M, Kanehori K, Kawamoto T, et al.: **Archaeal adaptation to higher temperatures revealed by genomic sequence of *Thermoplasma volcanium*.** *Proc Natl Acad Sci USA* 2000, **97**:14257-14262.
131. Fitz-Gibbon ST, Ladner H, Kim UJ, Stetter KO, Simon MI, Miller JH: **Genome sequence of the hyperthermophilic crenarchaeon *Pyrobaculum aerophilum*.** *Proc Natl Acad Sci USA* 2002, **99**:984-989.
132. Kawarabayasi Y, Hino Y, Horikawa H, Yamazaki S, Haikawa Y, Jin-no K, Takahashi M, Sekine M, Baba S, Ankai A, et al.: **Complete genome sequence of an aerobic hyper-thermophilic crenarchaeon, *Aeropyrum pernix* K1.** *DNA Res* 1999, **6**:83-101.
133. Kawarabayasi Y, Hino Y, Horikawa H, Jin-no K, Takahashi M, Sekine M, Baba S, Ankai A, Kosugi H, Hosoyama A, et al.: **Complete genome sequence of an aerobic thermoacidophilic crenarchaeon, *Sulfolobus tokodaii* strain 7.** *DNA Res* 2001, **8**:123-140.
134. Tersteegen A, Hedderich R: ***Methanobacterium thermoautotrophicum* encodes two multisubunit membrane-bound [NiFe] hydrogenases. Transcription of the operons and sequence analysis of the deduced proteins.** *Eur J Biochem* 1999, **264**:930-943.
135. Natale DA, Shankavaram UT, Galperin MY, Wolf YI, Aravind L, Koonin EV: **Towards understanding the first genome sequence of a crenarchaeon by genome annotation using clusters of orthologous groups of proteins (COGs).** *Genome Biol* 2000, **1**:research0009.1-0009.19.
136. Aravind L, Walker DR, Koonin EV: **Conserved domains in DNA repair proteins and evolution of repair systems.** *Nucleic Acids Res* 1999, **27**:1223-1242.
137. Cope GA, Suh GS, Aravind L, Schwarz SE, Zipursky SL, Koonin EV, Deshaies RJ: **Role of predicted metalloprotease motif of Jab1/Csn5 in cleavage of Nedd8 from Cul1.** *Science* 2002, **298**:608-611.
138. Verma R, Aravind L, Oania R, McDonald WH, Yates JR, 3rd, Koonin EV, Deshaies RJ: **Role of Rpn11 metalloprotease in deubiquitination and degradation by the 26S proteasome.** *Science* 2002, **298**:611-615.
139. Aravind L, Koonin EV: **DNA polymerase beta-like nucleotidyltransferase superfamily: identification of three new families, classification and evolutionary history.** *Nucleic Acids Res* 1999, **27**:1609-1618.
140. Daugherty M, Vonstein V, Overbeek R, Osterman A: **Archaeal shikimate kinase, a new member of the GHMP-kinase family.** *J Bacteriol* 2001, **183**:292-300.
141. van der Oost J, Huynen MA, Verhees CH: **Molecular characterization of phosphoglycerate mutase in archaea.** *FEMS Microbiol Lett* 2002, **212**:111-120.
142. Rashid N, Imanaka H, Kanai T, Fukui T, Atomi H, Imanaka T: **A novel candidate for the true fructose-1,6-bisphosphatase in archaea.** *J Biol Chem* 2002, **277**:30649-30655.
143. Constantinesco F, Forterre P, Elie C: **NurA, a novel 5'-3' nuclease gene linked to rad50 and mre11 homologs of thermophilic Archaea.** *EMBO Rep* 2002, **3**:537-542.
144. Graham DE, Bock CL, Schalk-Hihi C, Lu ZJ, Markham GD: **Identification of a highly diverged class of S-adenosylmethionine synthetases in the archaea.** *J Biol Chem* 2000, **275**:4055-4059.
145. Graham DE, Xu H, White RH: ***Methanococcus jannaschii* uses a pyruvoyl-dependent arginine decarboxylase in polyamine biosynthesis.** *J Biol Chem* 2002, **277**:23500-23507.

comment

reviews

reports

deposited research

refereed research

interactions

information