

# **METHIONINE ADENOSYLTRANSFERASE AS A USEFUL MOLECULAR SYSTEMATICS TOOL REVEALED BY PHYLOGENETIC AND STRUCTURAL ANALYSES**

Gabino F. Sánchez-Pérez<sup>1</sup>, José M. Bautista<sup>2</sup>, María A. Pajares\*<sup>1</sup>

<sup>1</sup>Instituto de Investigaciones Biomédicas "Alberto Sols" (CSIC-UAM), Arturo Duperier 4, 28029 Madrid (Spain).

<sup>2</sup>Departamento de Bioquímica y Biología Molecular IV, Facultad de Veterinaria, Universidad Complutense de Madrid, 28040 Madrid (Spain).

\*Corresponding author: María A. Pajares, Instituto de Investigaciones Biomédicas "Alberto Sols" (CSIC-UAM), Arturo Duperier 4, 28029 Madrid, Spain. Phone number: 34-91-5854414, FAX number: 34-91-5854401, e-mail:mapajares@iib.uam.es

Evolution of methionine adenosyltransferase

## SUMMARY

Structural and phylogenetic relationships among *Bacteria* and *Eukaryota* were analyzed by examining 292 methionine adenosyltransferase (MAT) amino acid sequences with respect to the crystal structure of this enzyme established for *Escherichia coli* and rat liver. Approximately 30% of MAT residues were found to be identical in all species. Five highly conserved amino acid sequence blocks did not vary in the MAT family. We detected specific structural features that correlated with sequence signatures for several clades, allowing taxonomical identification by sequence analysis. In addition, the number of amino acid residues in the loop connecting  $\alpha$ -strands A2 and A3 served to clearly distinguish sequences between eukaryotes and eubacteria. The molecular phylogeny of MAT genes in eukaryotes can be explained in terms of functional diversification coupled to gene duplication or alternative splicing and adaptation through strong structural constraints. Sequence analyses and intron/exon junction positions among nematodes, arthropods and vertebrates support the traditional Coelomata hypothesis. In vertebrates, the liver MAT I isoenzyme has gradually adapted its sequence towards one providing a more specific liver function. MAT phylogeny also served to cluster the major bacterial groups, demonstrating the superior phylogenetic performance of this ubiquitous, housekeeping gene in reconstructing the evolutionary history of distant relatives.

**Keywords:** Methionine adenosyltransferase, S-adenosylmethionine, evolution, methionine metabolism.

S-adenosylmethionine (SAM) is the main methyl donor in the transmethylation of numerous essential cell constituents (DNA, neurotransmitters, phospholipids, and many small molecules) <sup>1</sup>. After decarboxylation, SAM acts as a propylamine group donor in the biosynthesis of some polyamines (spermine and spermidine) <sup>2</sup>. Its importance is reflected by the fact that this molecule participates in as many reactions as ATP. However, contrary to ATP which is produced in a large number of reactions, SAM synthesis occurs in only one reaction catalyzed by methionine adenosyltransferase (MAT, EC 2.5.1.6). MAT is generally a homotetrameric enzyme that uses methionine and ATP in a reaction dependent on the presence of K<sup>+</sup> and Mg<sup>2+</sup> ions to render SAM, pyrophosphate and inorganic phosphate <sup>3;4</sup>.

To date, many structure/function relationship studies have used either the *Escherichia coli* (c-MAT) or rat liver (rl-MAT) enzyme. These studies have provided a relevant amount of information on key residues of the protein, including cysteines and active-site amino acids <sup>5; 6; 7; 8; 9</sup>. Since the description of the first MAT gene <sup>10</sup>, a substantial number of genes encoding MATs of different origins have been cloned and characterized <sup>11; 12; 13; 14; 15; 16</sup>. The data obtained indicate exceptional conservation of the gene sequence among highly divergent species. At the amino acid level, c- and rl-MATs have been estimated to show 52% identity <sup>17</sup>. The crystal structures of these MATs are the only ones presently available that indicate that conservation also occurs at the structural level, and this is reflected by the essentially identical organization of the domains in the monomer <sup>7; 18</sup>.

The development of molecular phylogenetics has been generally based on small-subunit (SSU) and large-subunit (LSU) ribosomal RNA analysis <sup>19</sup>. However, several recent concerns have challenged the validity of rRNA as a unique phylogenetic marker. These concerns are related to biases in base composition, disparities in evolutionary rates among lineages, position-

dependent substitution patterns, alignment ambiguities among very distant species, etc. Thus, recent efforts have focused on assessing the use of other genes<sup>20; 21; 22; 23</sup> and large combined protein sequence data sets<sup>24; 25</sup> to reconstruct evolutionary relationships among organisms. It has even been suggested that it will be possible to reconstruct a robust universal phylogeny only if a core of conserved markers, not affected by lateral gene transfer, is identified<sup>26</sup>.

Despite a high degree of sequence and structural conservation in MAT, a large number of representative species for which there are available sequences, and vast knowledge on MAT structure and functionality in highly divergent species, this enzyme has not yet been considered as a possible phylogenetic marker. Although attempts have been made to reconstruct partial phylogenies<sup>15; 27; 28</sup> and MAT has been included in studies examining sets of several proteins<sup>29; 30; 31</sup>, no universal phylogenetic evaluation has used MAT as the marker. Hence, the aim of the present study was to assess the performance of MAT in phylogenetic reconstructions using the 292 sequences available to date, and to demonstrate its utility in molecular systematic studies.

## Results and Discussion

### MAT sequence identification and characterization

MAT sequences retrieved by database mining led to the identification of 303 candidate sequences in almost every eukaryote and bacterium, but none in archaea, for which non homologous replacement by a new type of MAT has been recently described<sup>32</sup>. The absence of the MAT gene from the complete genomes of members of the genus *Chlamydia* and the microsporidia *Encephalitozoon cuniculi*, both intracellular parasites<sup>33; 34</sup>, is remarkable. Sequences for *Rickettsia prowazekii* and *R. typhi*, also obligate intracellular parasites, were excluded from the analysis because of the recent detection of stop codons in their MAT gene sequences. This may be interpreted as indicative of a certain degree of genome degeneration<sup>35</sup>. The lack or degeneration of the MAT gene in these species may be explained by functional redundancy (obsolescence) or by the existence of another methyl donating pathway (the host produces the methyl donor to be used by the parasite or another compound replaces SAM).

*Giardia lamblia*, one ancient eukaryote, showed highly divergent sequence and was thus only included in preliminary studies, in which it appeared as the earliest eukaryote with unique structural features. *G. lamblia* has a 41 amino acid insertion in the loop that connects helix 4 with  $\alpha$ -strand A4, the significance of which is unknown.

### Sequence conservation vs. 3D structure constraints

A final alignment of 392 positions was obtained from the 292 MAT amino acid sequences considered. Positional identity was difficult to establish in areas located at the N- and C-terminals and at loops connecting secondary structure elements. These ambiguous positions were therefore excluded, leaving 330 parsimony-informative positions for the final analysis.

MAT protein alignment revealed the presence of 57 amino acids located in identical positions in 100% of MATs analyzed and 61 additional residues that were conserved in 90% of the species studied. This indicates that approximately 30% of MAT residues are identical in all species. Moreover, 49 amino acid positions were identical in 75-90% of the species, 45 in 60-75%, and 39 in 50-60%. Residue conservation was inhomogeneous along the MAT sequence. Further, regions such as the N-terminal, C-terminal and some intermediate regions showed a lower degree of identity, probably due the capacity of these areas to absorb a high variety of substitutions without affecting the overall conformation of the molecule and its function.

Core areas of greatest amino acid conservation were observed along 5 stretches which we denoted blocks I, II, III, IV and V (Figure 1). Block I comprises residues 20-47, including the  $\alpha$ -strand A1 and the  $\alpha$ -helix 1, and one of the methionine binding motifs,  $^{29}$ GHPDK $^{33}$  which is preserved in all MAT sequences<sup>7</sup>. Block II is defined by two separate areas bearing residues 114-122 and 132-143, flanking the flexible loop at the active site of the enzyme. This loop has been recently shown to be involved in controlling the catalytic efficiency of the enzyme<sup>36</sup>. Conservation of these two areas could be due to the need to preserve the correct orientation for the loop, or to the fact that it contains the ATP binding motif  $^{132}$ GAGDQG $^{137}$ <sup>17</sup>. The consensus sequence for ATP binding sites has been defined as GxGDxG plus a lysine located 16-28 residues upstream<sup>37</sup>. However, MAT seems to show no variation in this sequence, which is always GAGDQG, highlighting the significant role of alanine and glutamine residues in this enzyme's ATP binding motif. Block III comprises residues 177-189 that form part of  $\alpha$ -strand A2, including D180 and K182, two of the amino acids involved in catalysis<sup>7;18</sup>. Block IV is the largest and most conserved including 54 amino acids (246-300). This block contains the central loop connecting the N-terminal and central domains, as well as many residues directly involved

in substrate and cation binding<sup>7; 18</sup>. In addition, this block contains a high-glycine stretch (254-281 with 10 Gly), which is fully preserved in all MATs. Finally, block V corresponds to the C-terminal and includes  $\alpha$ -helix 9 (372-389). Besides these blocks, two further reasonably conserved regions were detected: the first comprises  $\alpha$ -strands B1 and B2 with two conserved motifs, residues 55-59 and 70-75, and residue E58, involved in a saline bond with the central loop<sup>7</sup>; and the second area includes residues 303-337, which form part of  $\alpha$ -strands C2, C3, and the end of helix 5. However, it must be emphasized that all the amino acids involved in substrate binding and catalysis occur in the blocks described above, and are fully conserved in the MAT family.

To assess the relationship between evolutionary conservation and surface accessibility, the amino acids were also classified according to the degree of identity among MAT sequences as four categories: fully (100%), highly (75-99%), moderately (50-75%) and poorly conserved (<50%) (Figure 2). Buried residues were often observed among the amino acids within the first and second categories (>75%), whereas exposed amino acids were poorly represented in these categories. This over-representation of buried residues among the most conserved residues could reflect their involvement in catalytic activity, correct folding and the stability of the final structure. The fact that they establish the highest number of interactions among residues in the protein structure and their hydrophobic character, may be indicative of their role in the folding nuclei of the monomeric intermediate, according to the overall folding mechanism that has been established<sup>38</sup>. This area is later involved in the association process that leads to the dimer, the minimum active unit of MAT enzymes<sup>2</sup>.

Surface mapping of the level of evolutionary conservation for each amino acid in the protein structure may help to identify functionally and/or structurally significant regions<sup>39</sup>. For

this purpose, we used available crystallographic data for rl-MAT I/III and color-coded the surface according to the previously defined residue conservation categories (Figure 3A-D). The conservation pattern shows that the preserved blocks defined above occur in the inner channel, where the active site is located, and in the conserved area exposed at both the entrance and walls of this channel. High conservation among residues at the subunit interface is, therefore, consistent with a role for these amino acids in the structure and function of MAT, as reported for other oligomeric enzymes<sup>40</sup>. The study of individual enzyme families reveals how binding and catalysis are optimized in nature through the inclusion of mutations that improve efficiency in cases in which no new function has been acquired<sup>41</sup>. However, the perfect conservation of active site residues in MATs is an exception to this rule, since it indicates the preservation of the catalytic mechanism during evolution with no modification. This would suggest that the special features of the reaction catalyzed by MAT cannot be easily improved.

It is also of interest that several residues outside the subunit interface are also highly conserved. These residues were found to be mainly located in loops connecting secondary structure elements, thus suggesting a key role in preserving the correct orientation between them. The relevance of preserving this orientation probably reflects their essential contribution to final protein folding.

Sequence alignment using conventional algorithms revealed 60% identity between eukaryotic and bacterial MATs. The high degree of sequence conservation indicates severe restrictions for the substitution of certain amino acids. Such restrictions may be determined by two factors. First, the location of certain residues in the active site is important because of their role in catalysis. Mutations affecting these or adjacent residues modify their relative orientation leading to considerably reduced enzyme activity<sup>7</sup>. The second factor is that correct orientation

and positioning of certain secondary structure elements seems to depend strongly on the presence of some amino acids in the connecting loops, thus their substitution may lead to wrongly folded structures with no activity. To date, no such modifications have been identified in any MAT sequence though these loops have not yet been mapped by site-directed mutational analysis. Evolutionary studies could, nevertheless, take advantage of these restrictions in the sequence. Residue changes in areas of low variation that could be functionally absorbed in a certain evolutionary setting, are very unlikely to be repeated or to revert back at a different evolutionary time. Thus, the mutational study of certain areas of the MAT sequence may serve to clarify certain evolutionary relationships due to the low probability of mutational saturation at the protein level.

### **MAT as a phylogenetic marker**

To establish whether MAT could be a useful tool for reconstructing phylogenetic relationships among clades, we used a final data set of 330 positions for distance and parsimony phylogenetic analyses (Figure 4). A general view of the unrooted MAT tree shows the separate grouping of *Eukaryota* and *Eubacteria* with high bootstrap support indicating a common evolutionary origin. Protein structure comparisons showed that sequences belonging to each group can be clearly distinguished according to the number of residues involved in the loop connecting  $\alpha$ -strands A2 and A3. Specifically, this loop was normally 3-4 residues shorter in bacteria, except for *Campylobacter jejuni* and *Deinococcus radiodurans*. This difference has structural implications as shown in Figure 5. A longer loop allows the establishment of 18 favorable interactions between N- and C-terminal domains, including the formation of a salt bridge between D192 and R313. However, in eubacteria, only 2 can be formed among such

interactions including a salt bridge and hence both domains remain more distant. Close inspection of the structural models constructed from the crystallographic data available, shows higher rigidity for the monomer in the eukaryota, due to contact of N- and C-terminal domains<sup>7; 18</sup>.

Members of the MAT family can be grouped into several clades, mostly corresponding to the main taxonomy arrangements. Analysis of MAT phylogeny allows the identification of consensus sequences and specific structural features for the groups (table 2). Additional specific characteristics for each major taxonomic group are detailed below:

### *1. Eukaryotes*

A general view of the eukaryotic MAT phylogeny supports an animal-fungal clade excluding green plants. This is consistent with results obtained using SSU rRNA<sup>42</sup> and protein reconstructions<sup>43; 44; 45</sup>.

#### *1.1 Viridiplantae*

Three to four MAT gene copies were identified in the plants. These show differential expression among tissues and during development<sup>27; 28; 46; 47; 48; 49</sup>. Our analysis also supports (bootstrap value >99%) previous plant MAT gene classification as two types (types I and II)<sup>27</sup>. This topology suggests ancestral duplication, but the absence of MAT genes in conifers prevents us knowing whether duplication occurred before or after the divergence of *Magnoliophyta* from *Coniferophyta*. MAT type I duplications indicate low sequence divergence, thus duplication events are likely to have occurred in a narrow time frame. Consequently, phylogenetic relationships in plants cannot be confidently resolved by sequence similarity analysis. The functional requirements of other MAT isoenzymes may stem from their differential regulation

and/or from new SAM functions observed in plants (e.g., production of the phytohormone ethylene)<sup>50</sup>.

It is generally accepted that the ancestors of vascular plants were similar to green algae (*Chlorophyta*). In our phylogenetic analysis, the green alga *Chlamydomonas reinhardtii* branches earliest from the plant clade with high bootstrap support (100%). Some genes may be non-functional because of extensive divergence (e.g., *Cicer arieticum*) or loss of the C-terminal conservation shown by the remaining MATs (e.g., *Gossypium hirsutum 3*). Some of these MAT copies may no longer be needed as in other gene duplications and thus degenerate to pseudogenes.

The MAT genes corresponding to conifers appeared apart from the remaining plants, indicating earlier divergence for this group. The inclusion of one of the two genes isolated from *Pinus contorta* among the *Leuconostocaceae*, which were initially considered an outgroup for the plants, deserves special mention<sup>28</sup>. Based on the high degree of confidence for this node, the most parsimonious explanation for this would be a horizontal gene transfer event.

### 1.2 Fungi

In this phylogenetic analysis, fungi clustered as a monophyletic group forming the second branch point from that leading to the animals. This topology is well-supported by other studies<sup>51</sup>. Moreover, the MAT tree shows the early separation of the phylum *Basidiomycota* from *Ascomycota*, with high bootstrap support (100%). Only one MAT gene could be identified in most of the fungal species. The exception was *Saccharomyces cerevisiae* which has recently duplicated its genome<sup>52</sup> and has two similar copies of the MAT gene.

### 1.3 Nematoda

Interestingly, *Caenorhabditis elegans* presents 5 MAT genes. One of them, C49F5.1, was separated by an earlier gene duplication explaining its large phylogenetic distance and least similarity to the other *C. elegans* MATs. Genome analysis of the MAT genes yielded the following observations: a) four genes on autosomic chromosome IV, whereas C49F5.1 appears X-linked; b) this last gene presents 4 exons compared to the others showing 6; and c) three MAT genes were identified in *C. briggsae*, one of which showed similar characteristics to C49F5.1. Thus, MAT genes in nematodes probably arose from one single copy that underwent duplication, leading to genes on chromosomes X and IV. Further duplications led to the four copies identified on chromosome IV. Divergences among the five gene sequences could be attributed to a higher mutational rate of the genes on chromosome X, or to the fact that the duplications in chromosome IV have recently taken place. Just as MAT type II in plants, C49F5.1 seems to be the result of a duplication event required for different regulation or new functions of the final reaction product, SAM.

#### *1.4 Arthropoda*

Arthropods differed from the other organisms containing specialized tissues in that only one MAT gene was identified. However, genomic analysis revealed the possibility that this gene may suffer alternative splicing of exon 4 which codifies the flexible loop. Differences among both exons occur in the highly variable region within the loop, whereas conserved regions remain unchanged. These results suggest a role for the flexible loop in the regulation of arthropod MAT that eliminates the need for duplication and tissue specific distribution for each gene, i.e., the most common mechanism seen in MAT from higher organisms.

#### *1.5 Vertebrata*

There are three isoforms of MAT in mammalian tissues that are encoded by two genes. MAT I and MAT III are tetrameric and dimeric forms, respectively, of the same gene product (MAT1A), which is mainly expressed in adult liver<sup>53</sup>. The ubiquitously expressed MAT II isoform is a heterotetramer (2<sub>1</sub>2<sub>2</sub>), whose catalytic subunit ( ) is encoded by a different gene (MAT2A)<sup>2</sup>. The protein isoforms also differ in their regulation and catalytic properties (affinities for methionine)<sup>2; 53</sup>. In every vertebrate reported here, both genes presented around 85% identity, differences always corresponding to the same regions. The MAT tree clearly separates both gene types and locates the duplication event after the divergence from Urochordata, but before the divergence from Teleostomi. This is not surprising, if we consider that many early chordate gene families were formed or expanded by large-scale DNA duplications<sup>54</sup>. Interestingly, liver ontogeny in vertebrates is concomitant with this MAT gene duplication suggesting that tissue-specialized enzyme forms were required for adaptation to the functions of this new organ. MAT I/III has been suggested as a marker for liver development, due to the differential expression of its isoforms in fetal and adult tissue<sup>55</sup>.

A further significant issue is identifying specific isozyme residues to which to ascribe functional roles for these conserved amino acids, such as has been done for other homotetrameric enzymes<sup>56; 57</sup>. This identification uses information from the sequence alignment of vertebrate MAT isozymes and is based on the criterion that the residue is conserved in orthologs and is distinct among paralogs. Surprisingly, not all the residues considered liver-specific occurred in all our vertebrate clades suggesting that the specialization of the liver MAT gene took place gradually. Thus, cysteine 121 is found at the flexible loop of the liver enzyme and, being the target for regulation by nitrosylation of the protein<sup>58</sup>, it mediates the response to oxidative stress. This residue is present in mammalian liver MAT, but not in that of *Gallus gallus*. By comparing

MAT I/III and MAT II among mammals, birds, amphibian and fishes we were able to identify further analogous positions (Table 2). Along with functional studies, the characterization of new MAT sequences from species intermediate between mammals and ancient chordates, may clarify the evolutionary adaptation related to functional diversification of the MAT liver gene.

Finally, sequence comparison and the evolution of intron/exon junction positions among bilateria (data not shown) suggest the presence of a common ancestor for arthropods and vertebrates that branched from nematodes. This observation is in line with previous studies based on 18S rRNA that assign nematodes and arthropods to a common clade (Ecdysozoa) separate from vertebrates<sup>59</sup>. Conversely, our results are more in agreement with a recent study based on more than 100 proteins that supports the traditional Coelomata hypothesis<sup>31</sup> grouping arthropods with vertebrates apart from nematodes. This reinforces the idea that MAT is a better marker than rRNA for resolving the branching order of the main animal lineages.

## 2. Bacteria

Five clades of bacteria were clearly distinguished in our molecular phylogeny of MAT based on the presence of a single gene. These clades are: *Proteobacteria*, *Cyanobacteria*, *CFB group*, *Actinobacteria* and *Firmicutes*. Single representative species of several phyla appear as independent branches of the tree: *Aquificae* (*A. aeolicum*), *Thermotogae* (*T. maritima*), *Deinococcus-Thermus* (*D. radiodurans*), *Chloroflexi* (*C. aurianticus*), *Planctomycetes* (*G. obscuriglobis*), *Fibrobacteres* (*F. succinogenes*), *Spirochaetes* (*B. burgdorferi*) and *Fusobacteria* (*F. nucleatum*). The anaerobic detoxifying bacterium *Dehalococcoides ethenogenes* also appeared as an independent branch confirming its position in a unique phylogenetic group as described by the use of 16S rRNA<sup>60</sup>. There is little resolution in the MAT tree reflecting the

position of these phyla, and hence no consistent phylogenetic relationships with other groups could be clearly established.

### 2.1. Proteobacteria

Proteobacteria have been previously divided into five phylogenetically distinct groups ( , , , and )<sup>19</sup>. In the present phylogenetic analysis, proteobacterial MAT sequences did not consistently form a clade. However, the -, -, - and - groups were found to cluster individually.

#### 2.1.1 $\alpha$ -Proteobacteria

All the -proteobacteria MAT considered grouped with near 100% confidence but branched apart from the rest of proteobacteria because of the extension of several loops only detected in this class (Table 2). The phylogenetic tree presents an early branching point bearing several members of the order *Rickettsiales*, indicating the early separation of this order from the remaining -proteobacteria. The order *Rhizobiales* was also clearly separated from the rest of the -proteobacteria. The most distinctive feature of species of this order was the presence of a 7 amino acid insertion between helix 1 and -strand B1. Other characteristics typical of rhizobial MATs were: a) a flexible loop at least 2 residues longer than in the other -proteobacteria; and b) an insertion between -strand B2 and helix 2. The exception to this rule was shown by *Methylobacterium extorquens* with none of these peculiarities, indicating the very early separation of this family from the remaining rhizobiales.

#### 2.1.2 $\beta$ -Proteobacteria

Based on this phylogenetic tree, MAT sequences from - and -proteobacteria show a common origin, whereas relationships with the rest of the proteobacteria remain obscure, given the <50% bootstrap value for the nodes separating these bacterial classes. Differences between

both groups rely, for example, on the presence of two extra residues in the flexible loop of the  $\alpha$ -proteobacteria. Of interest is the presence of *Acidithiobacillus ferrooxidans* among the  $\alpha$ -proteobacteria, despite its taxonomical classification as belonging to the order *Chromatiales* of the  $\beta$ -proteobacteria. However, the high statistical support for this branch in our analysis and the absence of the typical characteristics of a  $\beta$ -proteobacteria suggest its reclassification within  $\alpha$ -proteobacteria.

### 2.1.3 $\gamma$ -Proteobacteria

Most of the  $\gamma$ -proteobacteria clustered in the tree, showing relationships among the different orders in this class. All the members of the *Pasteurellales* included in the tree grouped together, as occurred for the *Enterobacteria*. MAT sequences in members of the order *Legionellales* slightly diverged from those of the remaining  $\gamma$ -proteobacteria. This leads to the early branching of *Legionellales* from the  $\gamma$ -proteobacteria group.

### 2.2 Cyanobacteria

Based on MAT phylogeny, the *Cyanobacteria* appeared as an independent taxonomic group that was well differentiated by specific insertions (Table 2). Two of the four members of subsection I included in this study (*Prochlorococcus marinus* and *Synechococcus sp.*) carried a 7-residue insertion between helix 4 and  $\alpha$ -strand A4. In contrast, remaining cyanobacteria only have 6 amino acids at this position. In addition, these two species have an extra residue inserted in the loop between helix 2 and  $\alpha$ -strand B3. This could explain why *P. marinus* and *Synechococcus sp.* occur apart from the other members of subsections I, III and IV.

### 2.3 CFB group

The CFB group was included in a homogeneous cluster due to the specific characteristics of this phylum (Table 2). It should be highlighted that phylogenetic reconstruction places

*Chlorobium tepidum* on an early branch from the CFB group, although it presents distinctive features such as a 4-residue insertion in the flexible loop, and the absence of insertions in the loop between helix 4 and  $\alpha$ -strand A4. This position in the tree indicates the existence of a common ancestor for both phyla and their subsequent divergence. The class *Sphingobacteria* reflects this group's peculiarities: the loop between  $\alpha$ -strand A3 and helix 4 does not include the 10-12 residue insertion detected in members of the class *Bacteroidetes*. Moreover, in *Chlorobium* this insertion bears 7 residues.

#### 2.4 Actinobacteria

Members of the group *Actinomycetales* clustered perfectly together, with the order *Bifidobacteriales* remaining apart. The following modifications to the general characteristics were, nevertheless, observed: a) longer insertions in the flexible loop in *B. longum* and *S. fradiae*; b) *Mycobacterium avium*, *M. leprae* and *M. tuberculosis* have a 3-residue insertion in the loop connecting  $\alpha$ -strand B2 and helix 2; c) the insertion in the loop between helix 4 and  $\alpha$ -strand A4 is absent in *Thermobifida* and differs in length between *Streptomyces* (4 residues) and *Mycobacterium*, *Corynebacterium* and *Bifidobacterium* (3-2 residues).

#### 2.5 Firmicutes

2.5.1 *Bacilli/Clostridia*. All *Clostridia* and *Bacilli* clustered together but separately defined each class. Moreover, within *Bacilli*, the orders *Bacillales* and *Lactobacillales* were also separately arranged. The topology observed for the family *Leuconostaceae* is an exception, which can be explained by the early divergence from a common ancestor giving rise to a branch for *Leuconostocaceae* and a later divergence rendering the *Bacilli* and *Clostridia* families. Alternatively, the MAT sequence in *Leuconostocaceae* diverged rapidly.

#### 2.5.2 Mollicutes

This class clustered apart from the other *Firmicutes* (*Clostridia*, *Bacilli*). This might be explained by the presence of a 7-amino acid deletion in the loop connecting helix 2 and  $\beta$ -strand B3. *Mycoplasma pneumoniae* and *M. genitalium* clustered apart from *M. pulmonis* and *U. urealyticum*. These last species presented a 2-amino acid insertion in the loop between helices 8 and 9. In contrast, the separation of *M. pneumoniae* and *M. genitalium* may be attributed to differential structural features: a) a flexible loop two amino acids shorter; b) a residue inserted in the loop between helix 4 and  $\beta$ -strand A4; c) another residue inserted in the loop connecting  $\beta$ -strand C3 and helix 6; and d) the insertion of a residue in the loop between helices 7 and 8. In addition, the positioning of *Mycoplasma* among this group is congruent with results derived from protein fusion trees, in contrast to rRNA trees<sup>26</sup>. This finding is yet another confirmation of the idea that protein sequences more accurately reflect the phylogenetic position of *Mycoplasma*<sup>61</sup>.

## Conclusions

The present study is a first attempt at using the housekeeping MAT gene as a marker in eukarya and bacterial systematics, as an alternative to rRNA and other protein reference markers. This work is the result of integrating data yielded by intensive data mining, robust aligning of MAT sequences and structural-functional analyses. Through the detection of fully conserved regions in the MAT protein of all species, regions varying even among close relatives and characteristic structural features missing from certain taxonomic groups, this new tool enabled us to resolve phylogenetic relationships between close and distant relatives with high bootstrap support.

## Materials and Methods

### MAT sequences

MAT amino acid sequences were deduced from DNA sequence data available from complete or nearly complete publicly available genomes by conducting a TBLASTN search using rl-MAT as probe. For further references to amino acid positions rl-MAT is used as the consensus sequence. Candidate sequences were identified as MAT when they met the following criteria: a) a length of 370-414 amino acids; b) an N-terminal sequence containing the motif <sup>21</sup>FTSESVxEGHPDK<sup>33</sup>; and c) a C-terminal including the motif <sup>379</sup>GHFGxxxxxWE<sup>389</sup>. MAT sequences were obtained from GenBank, The Institute for Genomic Research (TIGR), DOE Joint Genome Institute, The Sanger Institute, University of Oklahoma, Baylor College of Medicine, Genome Sequencing Center at Washington University Medical School, Columbia Genome Center, Whitehead/MIT Genome Center, Genoscope, DNA Data Bank of Japan, Université Catholique de Louvain and European Bioinformatics Institute. A complete list of the sequences and their sources is provided in table 1.

### Alignment and sequence analysis

Sequences were aligned using the program BIOEDIT <sup>62</sup>. The alignment was manually refined to correct for large inserts and ambiguities. Final alignments contained 292 sequences 370-414 amino acids long from which 392 positions were selected for phylogenetic reconstruction. A total number of 330 informative positions were taken into account. Identity percentages were calculated using MULTALIN <sup>63</sup>. For additional information visit our database home page at <http://www.iib.uam.es/MAT>.

## **Evolutionary analysis**

MAT's evolutionary tree was established from the 292 sequence alignment of the 303 sequences available using neighbor-joining and parsimony methods. MAT sequences for the following microorganisms were discarded in the final analysis due to ambiguity: *Amoeba proteus*, *Campylobacter jejuni*, *Escherichia coli II*, *Giardia lamblia*, *Mesorhizobium loti II*, *Nostoc sp. PCC 7120. II*, *Rickettsia prowazekii*, *Rickettsia typhi*, *Treponema pallidum*, *Leptospira interrogans* and *Tropheryma whipplei*. Ambiguity may arise when there is a loss in the phylogenetic signal, obsolescence, insufficient taxon sampling or high sequence divergence. The phylogenetic analysis was performed using the PHYLIP package v.3.5<sup>64</sup>. Distances were calculated using Dayhoff in PROTDIST followed by the use of NEIGHBOR or PROTPARS for tree reconstruction. Phylogenetic trees were drawn using the MEGA2<sup>65</sup> and ATV programs<sup>66</sup>. Statistical support for the groups in the tree was evaluated by bootstrap analysis of 500 iterations using SEQBOOT.

## **Amino acid solvent accessibility**

Coordinates for rl-MAT crystal structure<sup>7</sup> were used to calculate the solvent accessibility (SA) of each individual residue in the dimeric structure using the SwissPdbViewer<sup>67</sup>. Amino acids with a SA <10% are regarded as buried, whereas residues with SA >10% are considered exposed<sup>68</sup>.

## Footnotes

## Acknowledgments

This work has been supported by grants of Fondo de Investigación Sanitaria of the Instituto de Salud Carlos III (01/1077 and RCMN C03/08) and MCYT (BMC-2002-00243) (to M.A.P.), and MCYT (PM99-0049-C02-01) (to J.M.B.).

**Abbreviations:** MAT, ATP:L-methionine adenosyltransferase; c-MAT, *E. coli* methionine adenosyltransferase; rl-MAT, rat liver methionine adenosyltransferase; SAM, S-adenosylmethionine; L-*cis*AMB, L-2-amino-4-methoxy-*cis*-but-3-enoic acid.

.

## References

1. Cantoni, G. L. (1975). Biological methylation: selected aspects. *Annu Rev Biochem* 44, 435-51.
2. Mato, J. M., Alvarez, L., Ortiz, P. & Pajares, M. A. (1997). S-adenosylmethionine synthesis: molecular mechanisms and clinical implications. *Pharmacol Ther* 73, 265-80.
3. Cantoni, G. L. (1953). S-adenosylmethionine: a new intermediate formed enzymatically from L-methionine and adenosinetriphosphate. *J Biol Chem* 204, 403-416.
4. McQueney, M. S. & Markham, G. D. (1995). Investigation of monovalent cation activation of S-adenosylmethionine synthetase using mutagenesis and uranyl inhibition. *J Biol Chem* 270, 18277-84.
5. Mingorance, J., Alvarez, L., Sanchez-Gongora, E., Mato, J. M. & Pajares, M. A. (1996). Site-directed mutagenesis of rat liver S-adenosylmethionine synthetase. Identification of a cysteine residue critical for the oligomeric state. *Biochem J* 315 ( Pt 3), 761-6.
6. Taylor, J. C. & Markham, G. D. (1999). The bifunctional active site of s-adenosylmethionine synthetase. Roles of the active site aspartates. *J Biol Chem* 274, 32909-14.
7. Gonzalez, B., Pajares, M. A., Hermoso, J. A., Alvarez, L., Garrido, F., Sufrin, J. R. & Sanz-Aparicio, J. (2000). The crystal structure of tetrameric methionine adenosyltransferase from rat liver reveals the methionine-binding site. *J Mol Biol* 300, 363-75.
8. Taylor, J. C. & Markham, G. D. (2000). The bifunctional active site of S-adenosylmethionine synthetase. Roles of the basic residues. *J Biol Chem* 275, 4060-5.
9. Sanchez-Perez, G. F., Gasset, M., Calvete, J. J. & Pajares, M. A. (2003). Role of an Intrasubunit Disulfide in the Association State of the Cytosolic Homo-oligomer Methionine Adenosyltransferase. *J Biol Chem* 278, 7285-7293.
10. Markham, G. D., DeParasis, J. & Gatmaitan, J. (1984). The sequence of metK, the structural gene for S-adenosylmethionine synthetase in Escherichia coli. *J Biol Chem* 259, 14505-7.
11. Thomas, D. & Surdin-Kerjan, Y. (1987). SAM1, the structural gene for one of the S-adenosylmethionine synthetases in Saccharomyces cerevisiae. Sequence and expression. *J Biol Chem* 262, 16704-9.
12. Larsson, J. & Rasmuson-Lestander, A. (1994). Molecular cloning of the S-adenosylmethionine synthetase gene in Drosophila melanogaster. *FEBS Lett* 342, 329-33.
13. Mautino, M. R., Barra, J. L. & Rosa, A. L. (1996). eth-1, the Neurospora crassa locus encoding S-adenosylmethionine synthetase: molecular cloning, sequence analysis and in vivo overexpression. *Genetics* 142, 789-800.
14. Yocum, R. R., Perkins, J. B., Howitt, C. L. & Pero, J. (1996). Cloning and characterization of the metE gene encoding S-adenosylmethionine synthetase from Bacillus subtilis. *J Bacteriol* 178, 4604-10.
15. Chiang, P. K., Chamberlin, M. E., Nicholson, D., Soubes, S., Su, X., Subramanian, G., Lanar, D. E., Prigge, S. T., Scovill, J. P., Miller, L. H. & Chou, J. Y. (1999). Molecular characterization of Plasmodium falciparum S-adenosylmethionine synthetase. *Biochem J* 344 Pt 2, 571-6.
16. Reguera, R. M., Balana-Fouce, R., Perez-Pertejo, Y., Fernandez, F. J., Garcia-Estrada, C., Cubria, J. C., Ordonez, C. & Ordonez, D. (2002). Cloning expression and

- characterization of methionine adenosyltransferase in *Leishmania infantum* promastigotes. *J Biol Chem* 277, 3158-67.
17. Horikawa, S., Ishikawa, M., Ozasa, H. & Tsukada, K. (1989). Isolation of a cDNA encoding the rat liver S-adenosylmethionine synthetase. *Eur J Biochem* 184, 497-501.
  18. Takusagawa, F., Kamitori, S. & Markham, G. D. (1996). Structure and function of S-adenosylmethionine synthetase: crystal structures of S-adenosylmethionine synthetase with ADP, BrADP, and PPi at 28 angstroms resolution. *Biochemistry* 35, 2586-96.
  19. Woese, C. R. (1987). Bacterial evolution. *Microbiol Rev* 51, 221-71.
  20. Gupta, R. S. & Golding, G. B. (1993). Evolution of HSP70 gene and its implications regarding relationships between archaeobacteria, eubacteria, and eukaryotes. *J Mol Evol* 37, 573-82.
  21. Gupta, R. S., Aitken, K., Falah, M. & Singh, B. (1994). Cloning of *Giardia lamblia* heat shock protein HSP70 homologs: implications regarding origin of eukaryotic cells and of endoplasmic reticulum. *Proc Natl Acad Sci U S A* 91, 2895-9.
  22. Eisen, J. A. (1995). The RecA protein as a model molecule for molecular systematic studies of bacteria: comparison of trees of RecAs and 16S rRNAs from the same species. *J Mol Evol* 41, 1105-23.
  23. Baldauf, S. L., Palmer, J. D. & Doolittle, W. F. (1996). The root of the universal tree and the origin of eukaryotes based on elongation factor phylogeny. *Proc Natl Acad Sci U S A* 93, 7749-54.
  24. Brown, J. R., Douady, C. J., Italia, M. J., Marshall, W. E. & Stanhope, M. J. (2001). Universal trees based on large combined protein sequence data sets. *Nat Genet* 28, 281-5.
  25. Baldauf, S. L., Roger, A. J., Wenk-Siefert, I. & Doolittle, W. F. (2000). A kingdom-level phylogeny of eukaryotes based on combined protein data. *Science* 290, 972-7.
  26. Brochier, C., Bapteste, E., Moreira, D. & Philippe, H. (2002). Eubacterial phylogeny based on translational apparatus proteins. *Trends Genet* 18, 1-5.
  27. Schroder, G., Eichel, J., Breinig, S. & Schroder, J. (1997). Three differentially expressed S-adenosylmethionine synthetases from *Catharanthus roseus*: molecular and functional characterization. *Plant Mol Biol* 33, 211-22.
  28. Lindroth, A. M., Saarikoski, P., Flygh, G., Clapham, D., Gronroos, R., Thelander, M., Ronne, H. & von Arnold, S. (2001). Two S-adenosylmethionine synthetase-encoding genes differentially expressed during adventitious root development in *Pinus contorta*. *Plant Mol Biol* 46, 335-46.
  29. Feng, D. F., Cho, G. & Doolittle, R. F. (1997). Determining divergence times with a protein clock: update and reevaluation. *Proc Natl Acad Sci U S A* 94, 13028-33.
  30. Mushegian, A. R., Garey, J. R., Martin, J. & Liu, L. X. (1998). Large-scale taxonomic profiling of eukaryotic model organisms: a comparison of orthologous proteins encoded by the human, fly, nematode, and yeast genomes. *Genome Res* 8, 590-8.
  31. Blair, J. E., Ikeo, K., Gojobori, T. & Hedges, S. B. (2002). The evolutionary position of nematodes. *BMC Evol Biol* 2, 7.
  32. Graham, D. E., Bock, C. L., Schalk-Hihi, C., Lu, Z. J. & Markham, G. D. (2000). Identification of a highly diverged class of S-adenosylmethionine synthetases in the archaea. *J Biol Chem* 275, 4055-9.
  33. Wicher, V., Baughn, R. E., Fuentealba, C., Shaddock, J. A., Abbruscato, F. & Wicher, K. (1991). Enteric infection with an obligate intracellular parasite, *Encephalitozoon cuniculi*, in an experimental model. *Infect Immun* 59, 2225-31.

34. Zomorodipour, A. & Andersson, S. G. (1999). Obligate intracellular parasites: *Rickettsia prowazekii* and *Chlamydia trachomatis*. *FEBS Lett* 452, 11-5.
35. Andersson, J. O. & Andersson, S. G. (1999). Genome degradation is an ongoing process in *Rickettsia*. *Mol Biol Evol* 16, 1178-91.
36. Taylor, J. C., Takusagawa, F. & Markham, G. D. (2002). The active site loop of S-adenosylmethionine synthetase modulates catalytic efficiency. *Biochemistry* 41, 9358-69.
37. Wierenga, R. K., Terpstra, P. & Hol, W. G. (1986). Prediction of the occurrence of the ADP-binding beta alpha beta-fold in proteins, using an amino acid sequence fingerprint. *J Mol Biol* 187, 101-7.
38. Gasset, M., Alfonso, C., Neira, J. L., Rivas, G. & Pajares, M. A. (2002). Equilibrium unfolding studies of the rat liver methionine adenosyltransferase III, a dimeric enzyme with intersubunit active sites. *Biochem J* 361, 307-15.
39. Glaser, F., Pupko, T., Paz, I., Bell, R. E., Bechor-Shental, D., Martz, E. & Ben-Tal, N. (2003). ConSurf: identification of functional regions in proteins by surface-mapping of phylogenetic information. *Bioinformatics* 19, 163-4.
40. Notaro, R., Afolayan, A. & Luzzatto, L. (2000). Human mutations in glucose 6-phosphate dehydrogenase reflect evolutionary history. *Faseb J* 14, 485-94.
41. Todd, A. E., Orengo, C. A. & Thornton, J. M. (2002). Plasticity of enzyme active sites. *Trends Biochem Sci* 27, 419-26.
42. Van de Peer, Y., Baldauf, S. L., Doolittle, W. F. & Meyer, A. (2000). An updated and comprehensive rRNA phylogeny of (crown) eukaryotes based on rate-calibrated evolutionary distances. *J Mol Evol* 51, 565-76.
43. Stechmann, A. & Cavalier-Smith, T. (2002). Rooting the eukaryote tree by using a derived gene fusion. *Science* 297, 89-91.
44. Roger, A. J., Sandblom, O., Doolittle, W. F. & Philippe, H. (1999). An evaluation of elongation factor 1 alpha as a phylogenetic marker for eukaryotes. *Mol Biol Evol* 16, 218-33.
45. Baldauf, S. L. (1999). A Search for the Origins of Animals and Fungi: Comparing and Combining Molecular Data. *Am Nat* 154, S178-S188.
46. Espartero, J., Pintor-Toro, J. A. & Pardo, J. M. (1994). Differential accumulation of S-adenosylmethionine synthetase transcripts in response to salt stress. *Plant Mol Biol* 25, 217-27.
47. Vander Mijnsbrugge, K., Van Montagu, M., Inze, D. & Boerjan, W. (1996). Tissue-specific expression conferred by the S-adenosyl-L-methionine synthetase promoter of *Arabidopsis thaliana* in transgenic poplar. *Plant Cell Physiol* 37, 1108-15.
48. Lee, J. H., Chae, H. S., Hwang, B., Hahn, K. W., Kang, B. G. & Kim, W. T. (1997). Structure and expression of two cDNAs encoding S-adenosyl-L-methionine synthetase of rice (*Oryza sativa* L.). *Biochim Biophys Acta* 1354, 13-8.
49. Gomez-Gomez, L. & Carrasco, P. (1998). Differential expression of the S-adenosyl-L-methionine synthase genes during pea development. *Plant Physiol* 117, 397-405.
50. Kende, H. (1993). Ethylene biosynthesis. *Annu Rev Plant Physiol Plant Mol Biol* 44, 283-307.
51. Baldauf, S. L. & Palmer, J. D. (1993). Animals and fungi are each other's closest relatives: congruent evidence from multiple proteins. *Proc Natl Acad Sci U S A* 90, 11558-62.

52. Wolfe, K. H. & Shields, D. C. (1997). Molecular evidence for an ancient duplication of the entire yeast genome. *Nature* 387, 708-13.
53. Lu, S. C., Gukovsky, I., Lugea, A., Reyes, C. N., Huang, Z. Z., Chen, L., Mato, J. M., Bottiglieri, T. & Pandol, S. J. (2003). Role of S-adenosylmethionine in two experimental models of pancreatitis. *Faseb J* 17, 56-8.
54. McLysaght, A., Hokamp, K. & Wolfe, K. H. (2002). Extensive genomic duplication during early chordate evolution. *Nat Genet* 31, 200-4.
55. Gil, B., Casado, M., Pajares, M. A., Bosca, L., Mato, J. M., Martin-Sanz, P. & Alvarez, L. (1996). Differential expression pattern of S-adenosylmethionine synthetase isoenzymes during rat liver development. *Hepatology* 24, 876-81.
56. Chang, S. H. & Kemp, R. G. (2002). Role of Ser530, Arg292, and His662 in the allosteric behavior of rabbit muscle phosphofructokinase. *Biochem Biophys Res Commun* 290, 670-5.
57. Pezza, J. A., Choi, K. H., Berardini, T. Z., Beernink, P. T., Allen, K. N. & Tolan, D. R. (2003). Spatial Clustering of Isozyme-specific Residues Reveals Unlikely Determinants of Isozyme Specificity in Fructose-1,6-bisphosphate Aldolase. *J Biol Chem* 278, 17307-13.
58. Avila, M. A., Mingorance, J., Martinez-Chantar, M. L., Casado, M., Martin-Sanz, P., Bosca, L. & Mato, J. M. (1997). Regulation of rat liver S-adenosylmethionine synthetase during septic shock: role of nitric oxide. *Hepatology* 25, 391-6.
59. Aguinaldo, A. M., Turbeville, J. M., Linford, L. S., Rivera, M. C., Garey, J. R., Raff, R. A. & Lake, J. A. (1997). Evidence for a clade of nematodes, arthropods and other moulting animals. *Nature* 387, 489-93.
60. Hendrickson, E. R., Payne, J. A., Young, R. M., Starr, M. G., Perry, M. P., Fahnestock, S., Ellis, D. E. & Ebersole, R. C. (2002). Molecular analysis of Dehalococcoides 16S ribosomal DNA from chloroethene-contaminated sites throughout North America and Europe. *Appl Environ Microbiol* 68, 485-95.
61. Kamla, V., Henrich, B. & Hadding, U. (1996). Phylogeny based on elongation factor Tu reflects the phenotypic features of mycoplasmas better than that based on 16S rRNA. *Gene* 171, 83-7.
62. Hall, T. A. (1999). BioEdit: a user-friendly biological alignment editor and analysis program for Windows 95/98/NT. *Nucl Acids Symp* 41, 95-98.
63. Corpet, F. (1988). Multiple sequence alignment with hierarchical clustering. *Nucleic Acids Res* 16, 10881-90.
64. Felsenstein, J. (1993). PHYLIP (phylogeny inference package) version 3.5c. *University of Washington, Seattle*.
65. Kumar, S., Tamura, K., Jakobsen, I. B. & Nei, M. (2001). MEGA2: molecular evolutionary genetics analysis software. *Bioinformatics* 17, 1244-5.
66. Zmasek, C. M. & Eddy, S. R. (2001). ATV: display and manipulation of annotated phylogenetic trees. *Bioinformatics* 17, 383-4.
67. Guex, N. & Peitsch, M. C. (1997). SWISS-MODEL and the Swiss-PdbViewer: an environment for comparative protein modeling. *Electrophoresis* 18, 2714-23.
68. Goldman, N., Thorne, J. L. & Jones, D. T. (1998). Assessing the impact of secondary structure and solvent accessibility on protein evolution. *Genetics* 149, 445-58.

## Figure legends

**Figure 1.** Evolution of the MAT protein. This diagram shows the analysis of 292 different MAT sequences. Rat liver MAT (lower row) was used as the reference sequence. The lettering for rat liver MAT indicates amino acid solvent accessibility (SA): blue, buried (SA<10%); black, exposed (SA>10%). Immediately above the rat liver MAT sequence, we show the identity consensus sequence for the 292 sequences analyzed. At each position, the most frequent amino acid is coded according to its frequency (blue: 100%; green: 90-99%; red: 75-90%; yellow: 60-75%). The conserved blocks described in the text are boxed and identified with Roman numerals. The locations of  $\alpha$ -helices and  $\beta$ -strands as found in the rat liver MAT 3D structure are provided above the identity consensus sequence. The structure of the first 16 residues (lowercase) in rat liver MAT is not known.

**Figure 2.** Conservation and 3-dimensional structure. Each bar represents the distribution of MAT amino acids among the different conservation categories for exposed, total, and buried residues, respectively.

**Figure 3.** Diagram showing MAT sequence conservation in 3-dimensional structure. The conservation pattern is color-coded on the molecular surface of rat liver MAT: dark violet indicates maximal conservation (100-90% identity), white indicates an average conservation level (90-50%) and dark turquoise indicates maximal variability (<50% identity). **(A)** Dimer viewed from the active site entrance. **(B)** Side view of the dimer after a 90° rotation to the right showing the exposed surface. **(C)** A view of the monomer from the monomer-monomer interface. **(D)** Side view of the monomer after a 180° rotation showing the exposed surface.

**Figure 4.** Phylogenetic tree derived from MAT sequences. Numbers on the branches show the percentage occurrence of nodes in 500 bootstrap replicates in the neighbor-joining and maximum parsimony analyses. Bootstrap values are indicated only when greater than 70%. For simplicity, the branches within each major group were collapsed.

**Figure 5.** Differences between rat liver MAT (**A**) (PDB ID code: 1QM4) and *E. coli* MAT (**B**) (PDB ID code: 1fug) for the atomic interactions between loops connecting A2 and A3 strands in the N-terminal domain (yellow ribbons and red backbone and side chains) and loops connecting helix 5, C2 and C3 strands in the C-terminal domain (blue ribbons and cyan backbone and side chains).

## Tables

**Table 1.** List of MAT sequences used in the present analysis.

<b><u>Specie</u></b>	<b><u>Accession N° / Genome center</u></b>
<i>Acanthamoeba castellanii</i>	<b><u>927487</u></b>
<i>Acidithiobacillus ferrooxidans</i>	TIGR
<i>Acinetobacter ADP-1</i>	Genoscope
<i>Actinidia chinensis 1</i>	<b><u>726030</u></b>
<i>Actinidia chinensis 2</i>	<b><u>726028</u></b>
<i>Actinobacillus actinomycetemcomitans</i>	Univ. Oklahoma
<i>Agrobacterium tumefaciens</i>	<b><u>17934278</u></b>
<i>Ajellomyces capsulatus</i>	Univ. Washington
<i>Anaplasma phagocytophila</i>	TIGR
<i>Anopheles gambiae</i>	<b><u>21291484</u></b>
<i>Aquifex aeolicus</i>	<b><u>7387875</u></b>
<i>Arabidopsis thaliana 1</i>	<b><u>15217781</u></b>
<i>Arabidopsis thaliana 2</i>	<b><u>15234354</u></b>
<i>Arabidopsis thaliana 3</i>	<b><u>15229033</u></b>
<i>Arabidopsis thaliana 4</i>	<b><u>15228048</u></b>
<i>Ascaris lumbricoides</i>	EBI Parasites
<i>Ascobolus immersus</i>	<b><u>836960</u></b>
<i>Aspergillus fumigatus</i>	TIGR
<i>Aspergillus nidulans</i>	Whitehead/MIT
<i>Azotobacter vinelandii</i>	<b><u>23102724</u></b>
<i>Bacillus anthracis</i>	<b><u>21402812</u></b>
<i>Bacillus cereus 1</i>	TIGR
<i>Bacillus cereus 2</i>	<b><u>29898393</u></b>
<i>Bacillus halodurans</i>	<b><u>20138752</u></b>
<i>Bacillus stearothermophilus</i>	Univ. Oklahoma
<i>Bacillus subtilis</i>	<b><u>7434008</u></b>
<i>Bacteroides forsythus</i>	TIGR
<i>Bacteroides fragilis</i>	Sanger Institute
<i>Bacteroides thetaiotaomicron</i>	<b><u>29340533</u></b>
<i>Bifidobacterium longum</i>	<b><u>23327085</u></b>
<i>Bombyx mori</i>	EST NCBI
<i>Bordetella avium</i>	Sanger Institute
<i>Bordetella bronchiseptica</i>	Sanger Institute
<i>Bordetella pertussis</i>	Sanger Institute
<i>Borrelia burgdorferi</i>	<b><u>7434004</u></b>
<i>Botrytis cinerea</i>	Genoscope
<i>Bradyrhizobium japonicum</i>	<b><u>27354222</u></b>
<i>Brassica juncea 1</i>	<b><u>10443981</u></b>
<i>Brassica juncea 2</i>	<b><u>14600070</u></b>
<i>Brassica juncea 3</i>	<b><u>14600072</u></b>
<i>Brucella melitensis</i>	<b><u>17088253</u></b>

<i>Brucella suis</i>	<b><u>23349048</u></b>
<i>Buchnera aphidicola</i>	<b><u>11386917</u></b>
<i>Burkholderia cepacia</i>	Sanger Institute
<i>Burkholderia fungorum</i>	<b><u>22985056</u></b>
<i>Burkholderia mallei</i>	TIGR
<i>Burkholderia pseudomallei</i>	Sanger Institute
<i>Caenorhabditis briggsae 1</i>	Sanger Institute
<i>Caenorhabditis briggsae 2</i>	Sanger Institute
<i>Caenorhabditis briggsae 3</i>	Sanger Institute
<i>Caenorhabditis elegans 1</i>	<b><u>17538494</u></b>
<i>Caenorhabditis elegans 2</i>	<b><u>1753849</u></b>
<i>Caenorhabditis elegans 3</i>	<b><u>21106027</u></b>
<i>Caenorhabditis elegans 4</i>	<b><u>7509275</u></b>
<i>Caenorhabditis elegans 5</i>	<b><u>17551082</u></b>
<i>Camellia sinensis</i>	<b><u>7594741</u></b>
<i>Campylobacter jejuni</i>	<b><u>11258525</u></b>
<i>Candida albicans</i>	<b><u>7271000</u></b>
<i>Carboxydotherrmus hydrogenoformans</i>	TIGR
<i>Carica papaya</i>	<b><u>22774026</u></b>
<i>Catharanthus roseus 1</i>	<b><u>1655576</u></b>
<i>Catharanthus roseus 2</i>	<b><u>1655578</u></b>
<i>Catharanthus roseus 3</i>	<b><u>1655580</u></b>
<i>Caulobacter crescentus</i>	<b><u>16124306</u></b>
<i>Chlamydomonas reinhardtii</i>	DOE Joint Institute
<i>Chlorobium tepidum</i>	<b><u>21646663</u></b>
<i>Chloroflexus aurantiacus</i>	<b><u>22971187</u></b>
<i>Cicer arietinum</i>	<b><u>1808591</u></b>
<i>Ciona intestinalis</i>	<b><u>23586111</u></b>
<i>Clavibacter michiganensis</i>	Sanger Institute
<i>Clostridium acetobutylicum</i>	<b><u>15896110</u></b>
<i>Clostridium botulinum</i>	Sanger Institute
<i>Clostridium perfringens</i>	<b><u>18311159</u></b>
<i>Clostridium tetani</i>	<b><u>28202517</u></b>
<i>Coccidioides posadasii</i>	TIGR
<i>Colwellia psychroerythraea</i>	TIGR
<i>Corynebacterium diphtheriae</i>	TIGR
<i>Corynebacterium efficiens</i>	<b><u>23493563</u></b>
<i>Corynebacterium glutamicum</i>	<b><u>19552815</u></b>
<i>Coxiella burnetii</i>	TIGR
<i>Cryptococcus neoformans</i>	TIGR
<i>Cryptosporidium parvum</i>	TIGR
<i>Cytophaga hutchinsonii</i>	<b><u>23137174</u></b>
<i>Danio rerio 1</i>	<b><u>28278852</u></b>
<i>Danio rerio 2</i>	TIGR
<i>Dechloromonas aromatica</i>	DOE Joint Institute
<i>Dehalococcoides ethenogenes</i>	TIGR
<i>Deinococcus radiodurans</i>	<b><u>15805667</u></b>
<i>Dendrobium crumenatum</i>	<b><u>16226050</u></b>

<i>Desulfovibrio desulfuricans</i>	<b><u>23475965</u></b>
<i>Desulfovibrio vulgaris</i>	TIGR
<i>Dianthus caryophyllus</i>	<b><u>7434012</u></b>
<i>Dichelobacter nodosus</i>	TIGR
<i>Dictyostelium discoideum</i>	Sanger Institute
<i>Drosophila melanogaster</i>	<b><u>7296263</u></b>
<i>Drosophila pseudoobscura</i>	Baylor
<i>Ehrlichia chaffensis</i>	TIGR
<i>Ehrlichia ruminantium</i>	Sanger Institute
<i>Elaeagnus umbellata 1</i>	<b><u>13540316</u></b>
<i>Elaeagnus umbellata 2</i>	<b><u>13540318</u></b>
<i>Enterococcus faecalis</i>	<b><u>29342835</u></b>
<i>Escherichia coli (metK)</i>	<b><u>1708999</u></b>
<i>Escherichia coli (metX)</i>	<b><u>26250367</u></b>
<i>Fibrobacter succinogenes</i>	TIGR
<i>Fusarium sporotrichioides</i>	Univ. Oklahoma
<i>Fusobacterium nucleatum</i>	<b><u>19703697</u></b>
<i>Gallus gallus</i>	TIGR
<i>Gemmata obscuriglobus</i>	TIGR
<i>Giardia lamblia</i>	<b><u>29247850</u></b>
<i>Glycine max 1</i>	TIGR
<i>Glycine max 2</i>	TIGR
<i>Glycine max 3</i>	TIGR
<i>Gossypium hirsutum 1</i>	TIGR
<i>Gossypium hirsutum 2</i>	TIGR
<i>Gossypium hirsutum 3</i>	TIGR
<i>Haemophilus influenzae</i>	<b><u>1170942</u></b>
<i>Haemophilus somnus</i>	<b><u>23467639</u></b>
<i>Helicobacter pylori</i>	<b><u>3024119</u></b>
<i>Helicobacter pylori J99</i>	<b><u>6685665</u></b>
<i>Heliobacillus mobilis</i>	<b><u>27262362</u></b>
<i>Homo sapiens 1</i>	<b><u>4557737</u></b>
<i>Homo sapiens 2</i>	<b><u>284394</u></b>
<i>Hordeum vulgare 1</i>	<b><u>7434000</u></b>
<i>Hordeum vulgare 2</i>	TIGR
<i>Hordeum vulgare 3</i>	TIGR
<i>Ictalurus punctatus</i>	TIGR
<i>Klebsiella pneumoniae</i>	Univ. Washington
<i>Lactobacillus gasseri</i>	DOE Joint Institute
<i>Lactobacillus plantarum</i>	<b><u>28378057</u></b>
<i>Lactococcus lactis</i>	<b><u>13878576</u></b>
<i>Lactuca sativa 1</i>	TIGR
<i>Lactuca sativa 2</i>	TIGR
<i>Legionella pneumophila</i>	Columbia
<i>Leishmania infantum</i>	<b><u>20387266</u></b>
<i>Leptospira interrogans</i>	<b><u>24215333</u></b>
<i>Leuconostoc mesenteroides</i>	<b><u>23023832</u></b>
<i>Listeria innocua</i>	<b><u>1680041</u></b>

<i>Listeria monocytogenes</i>	<b><u>16411100</u></b>
<i>Litchi chinensis</i>	<b><u>30142157</u></b>
<i>Lotus japonicus 1</i>	<b><u>21907982</u></b>
<i>Lotus japonicus 2</i>	TIGR
<i>Lycopersicon esculentum 1</i>	<b><u>1084406</u></b>
<i>Lycopersicon esculentum 2</i>	<b><u>481566</u></b>
<i>Lycopersicon esculentum 3</i>	<b><u>1084408</u></b>
<i>Lycopersicon esculentum 4</i>	TIGR
<i>Magnaporthe grisea</i>	TIGR
<i>Magnetococcus sp. MC-1</i>	<b><u>23000957</u></b>
<i>Magnetospirillum magnetotacticum</i>	<b><u>23015399</u></b>
<i>Mannheimia haemolytica</i>	Baylor
<i>Medicago truncatula 1</i>	TIGR
<i>Medicago truncatula 2</i>	TIGR
<i>Medicago truncatula 3</i>	TIGR
<i>Medicago truncatula 4</i>	TIGR
<i>Mesembryanthemum crystallinum</i>	<b><u>1724104</u></b>
<i>Mesorhizobium loti 1</i>	<b><u>20803994</u></b>
<i>Mesorhizobium loti 2</i>	<b><u>13475107</u></b>
<i>Methylobacterium extorquens</i>	Univ. Washington
<i>Methylococcus capsulatus</i>	TIGR
<i>Microbulbifer degradans</i>	<b><u>23027148</u></b>
<i>Moraxella catarrhalis</i>	DDJB Japan
<i>Mus musculus 1</i>	<b><u>476917</u></b>
<i>Mus musculus 2</i>	<b><u>13097429</u></b>
<i>Musa acuminata</i>	<b><u>2305014</u></b>
<i>Mycobacterium avium</i>	TIGR
<i>Mycobacterium bovis</i>	Sanger Institute
<i>Mycobacterium leprae</i>	<b><u>15214074</u></b>
<i>Mycobacterium marinum</i>	Sanger Institute
<i>Mycobacterium smegmatis</i>	TIGR
<i>Mycobacterium tuberculosis</i>	<b><u>3915763</u></b>
<i>Mycoplasma genitalium</i>	<b><u>1346527</u></b>
<i>Mycoplasma penetrans</i>	<b><u>26553547</u></b>
<i>Mycoplasma pneumoniae</i>	<b><u>2500686</u></b>
<i>Mycoplasma pulmonis</i>	<b><u>15829173</u></b>
<i>Myxococcus xanthus</i>	<b><u>27804841</u></b>
<i>Neisseria gonorrhoeae</i>	Univ. Oklahoma
<i>Neisseria meningitidis</i>	<b><u>11258516</u></b>
<i>Neorickettsia sennetsu</i>	TIGR
<i>Neurospora crassa</i>	<b><u>2133316</u></b>
<i>Nicotiana tabacum</i>	<b><u>7230379</u></b>
<i>Nitrosomonas europaea</i>	<b><u>22954639</u></b>
<i>Nostoc sp. PCC 7120 1</i>	<b><u>17231616</u></b>
<i>Nostoc sp. PCC 7120 2</i>	<b><u>17132339</u></b>
<i>Novosphingobium aromaticivorans</i>	<b><u>23110658</u></b>
<i>Oceanobacillus iheyensis</i>	<b><u>22777999</u></b>
<i>Opococcus canis</i>	<b><u>23038264</u></b>

<i>Oryza sativa</i> 1	<b><u>450549</u></b>
<i>Oryza sativa</i> 2	<b><u>1778821</u></b>
<i>Oryza sativa</i> 3	<b><u>8468037</u></b>
<i>Pasteurella multocida</i>	<b><u>13431697</u></b>
<i>Pectobacterium carotovorum</i>	Sanger Institute
<i>Pectobacterium chrysanthemi</i>	TIGR
<i>Petunia x hybrida</i> 1	<b><u>1084428</u></b>
<i>Petunia x hybrida</i> 2	<b><u>5726594</u></b>
<i>Phanerochaete chrysosporium</i>	DOE Joint Institute
<i>Phaseolus lunatus</i>	<b><u>18157331</u></b>
<i>Photorhabdus asymbiotica</i>	Sanger Institute
<i>Phytophthora infestans</i>	<b><u>23394401</u></b>
<i>Pinus banksiana</i>	<b><u>1033190</u></b>
<i>Pinus contorta</i> 1	<b><u>10441429</u></b>
<i>Pinus contorta</i> 2	<b><u>10441431</u></b>
<i>Pinus contorta</i> 3	TIGR
<i>Pinus contorta</i> 4	TIGR
<i>Pinus contorta</i> 5	TIGR
<i>Pisum sativum</i>	<b><u>2129889</u></b>
<i>Plasmodium berghei</i>	Sanger Institute
<i>Plasmodium chabaudi</i>	Sanger Institute
<i>Plasmodium falciparum</i>	<b><u>10129955</u></b>
<i>Plasmodium knowlesi</i>	Sanger Institute
<i>Plasmodium yoelii</i>	<b><u>23482440</u></b>
<i>Populus deltoides</i>	<b><u>497900</u></b>
<i>Porphyromonas gingivalis</i>	TIGR
<i>Prevotella intermedia</i>	TIGR
<i>Prochlorococcus marinus</i> 1	<b><u>23132136</u></b>
<i>Prochlorococcus marinus</i> 2	<b><u>23122176</u></b>
<i>Pseudomonas aeruginosa</i>	<b><u>15595743</u></b>
<i>Pseudomonas fluorescens</i>	<b><u>23060543</u></b>
<i>Pseudomonas putida</i>	<b><u>26991645</u></b>
<i>Pseudomonas syringae</i>	<b><u>23471420</u></b>
<i>Psychrobacter sp.</i> 273-4	DOE Joint Institute
<i>Ralstonia eutropha</i>	DOE Joint Institute
<i>Ralstonia metallidurans</i>	<b><u>22979451</u></b>
<i>Ralstonia solanacearum</i>	<b><u>17544853</u></b>
<i>Rattus norvegicus</i> 1	<b><u>92483</u></b>
<i>Rattus norvegicus</i> 2	<b><u>19705457</u></b>
<i>Rhizobium leguminosarum</i>	Sanger Institute
<i>Rhodobacter sphaeroides</i>	<b><u>22957691</u></b>
<i>Rhodopseudomonas palustris</i>	<b><u>22961217</u></b>
<i>Rhodospirillum rubrum</i>	<b><u>22965557</u></b>
<i>Rickettsia prowazekii</i>	<b><u>7387879</u></b>
<i>Rickettsia typhi</i>	<b><u>11133588</u></b>
<i>Saccharomyces cerevisiae</i> 1	<b><u>1346525</u></b>
<i>Saccharomyces cerevisiae</i> 2	<b><u>83324</u></b>
<i>Salmonella typhi</i>	<b><u>16766301</u></b>

<i>Schistosoma japonicum</i>	EBI Parasites
<i>Schizosaccharomyces pombe</i>	<b><u>7493352</u></b>
<i>Secale cereale</i>	TIGR
<i>Serratia marcescens</i>	Sanger Institute
<i>Shewanella oneidensis</i>	<b><u>24372516</u></b>
<i>Shigella flexneri</i>	<b><u>24114197</u></b>
<i>Silicibacter pomeroyi</i>	TIGR
<i>Sinorhizobium meliloti</i>	<b><u>15964164</u></b>
<i>Solanum tuberosum 1</i>	TIGR
<i>Solanum tuberosum 2</i>	TIGR
<i>Solanum tuberosum 3</i>	TIGR
<i>Sorghum bicolor</i>	TIGR
<i>Spiroplasma kunkelii</i>	Univ. Oklahoma
<i>Staphylococcus aureus</i>	<b><u>1709003</u></b>
<i>Staphylococcus epidermidis</i>	<b><u>9624212</u></b>
<i>Streptococcus agalactiae</i>	<b><u>29611810</u></b>
<i>Streptococcus equi</i>	Sanger Institute
<i>Streptococcus gordonii</i>	TIGR
<i>Streptococcus mitis</i>	TIGR
<i>Streptococcus mutans</i>	<b><u>29611808</u></b>
<i>Streptococcus pneumoniae</i>	<b><u>15902715</u></b>
<i>Streptococcus pyogenes</i>	<b><u>19746332</u></b>
<i>Streptococcus sobrinus</i>	TIGR
<i>Streptococcus suis</i>	Sanger Institute
<i>Streptococcus thermophilus</i>	Univ. C. Louvain
<i>Streptococcus uberis</i>	Sanger Institute
<i>Streptomyces avermitilis</i>	<b><u>29833416</u></b>
<i>Streptomyces coelicolor</i>	<b><u>21219978</u></b>
<i>Streptomyces fradiae</i>	<b><u>15554326</u></b>
<i>Streptomyces pristinaespiralis</i>	<b><u>2294502</u></b>
<i>Streptomyces spectabilis</i>	<b><u>7387884</u></b>
<i>Suaeda maritima</i>	<b><u>11992267</u></b>
<i>Synechococcus sp.</i>	<b><u>23134392</u></b>
<i>Synechocystis sp.</i>	<b><u>7434007</u></b>
<i>Takifugu rubripes 1</i>	DOE Joint Institute
<i>Takifugu rubripes 2</i>	DOE Joint Institute
<i>Tetraodon nigroviridis</i>	Genoscope
<i>Thalassiosira pseudonana</i>	DOE Joint Institute
<i>Theileria annulata</i>	Sanger Institute
<i>Theileria parva</i>	TIGR
<i>Thermobifida fusca</i>	<b><u>23017995</u></b>
<i>Thermosynechococcus elongatus</i>	<b><u>29611806</u></b>
<i>Thermotoga maritima</i>	<b><u>7387883</u></b>
<i>Treponema pallidum</i>	<b><u>15639781</u></b>
<i>Trichodesmium erythraeum</i>	<b><u>23040652</u></b>
<i>Triticum aestivum 1</i>	TIGR
<i>Triticum aestivum 2</i>	TIGR
<i>Triticum aestivum 3</i>	TIGR

<i>Triticum aestivum 4</i>	TIGR
<i>Tropheryma whipplei</i>	<b><u>28572562</u></b>
<i>Trypanosoma brucei</i>	TIGR
<i>Trypanosoma cruzi</i>	Sanger Institute
<i>Ureaplasma urealyticum</i>	<b><u>13357974</u></b>
<i>Vibrio cholerae</i>	<b><u>9654898</u></b>
<i>Vibrio parahaemolyticus</i>	<b><u>28899380</u></b>
<i>Vibrio vulnificus</i>	<b><u>27364907</u></b>
<i>Wigglesworthia brevipalpis</i>	<b><u>24324056</u></b>
<i>Wolbachia sp.</i>	TIGR
<i>X-bacteria</i>	Direct submission
<i>Xanthomonas axonopodis</i>	<b><u>21241583</u></b>
<i>Xanthomonas campestris</i>	<b><u>21230235</u></b>
<i>Xenopus laevis 1</i>	TIGR
<i>Xenopus laevis 2</i>	<b><u>27882050</u></b>
<i>Xenopus tropicalis</i>	Sanger Institute
<i>Xylella fastidiosa</i>	<b><u>15836994</u></b>
<i>Yersinia enterocolitica</i>	Sanger Institute
<i>Yersinia pestis</i>	<b><u>16121235</u></b>
<i>Zea mays 1</i>	TIGR
<i>Zea mays 2</i>	TIGR

**Table 2.** Specific structural features shown by each group on the MAT phylogenetic tree.

<b>Eukaryota</b>	The loop connecting $\alpha$ -strands A2-A3 has 10 residues		
	<b>Plants</b>	Presence of the sequences $^{93}\text{FXSXDVLXAD}^{103}$ , $^{216}\text{NDEIA}^{220}$ , and $^{331}\text{VFVD}^{334}$	
		Insertion of 3-4 residues in the loop connecting helices 7 and 8, between positions 369-370	
	<b>Vertebrates</b>	<b>Fishes</b>	Present specific residues for the MAT I enzyme, C69, C377 and V262
		<b>Birds/reptiles</b>	Present additional specific residues such as R82 and D167, besides those shown in fishes
<b>Mammals</b>		Present a typical cysteine in the MAT I flexible loop, C121 besides the residues specific for fishes, birds and reptiles	
<b>Eubacteria</b>	Deletion of 3-4 residues in the loop connecting $\alpha$ -strands A2-A3		
	<b>Proteobacteria</b>	Deletion of 2 residues in the loop connecting helix 2 and $\alpha$ -strand B3	
		Insertion of 4 residues between helix 4 and $\alpha$ -strand A4	
		Insertion of 4-5 residues between helices 8-9	
	<b>Cyanobacteria</b>	Insertion of 4 residues in the loop between $\alpha$ -strand B3 and helix 4	
		Insertion of 6-7 residues in the loop connecting helix 4 and $\alpha$ -strand A4	
		Insertion of 7 residues in the loop between helices 7 and 8	
		Insertion of 2 residues in the loop connecting helices 8 and 9	
	<b>CFB group</b>	Insertion of 1-2 residues in the loop connecting helix 3 and $\alpha$ -strand A2	
		Insertion of 2 residues in the loop connecting $\alpha$ -strand C3 and helix 6	
		Insertion of 9-11 residues in the loop connecting helix 4 and $\alpha$ -strand A4	
		Insertion of 21-22 residues in the loop connecting helices 8 and 9	
	<b>Actinobacteria</b>	Insertion of 8 residues in the flexible loop	
		Insertion of 2-4 residues in the loop connecting helix 4 and $\alpha$ -strand A4	
	<b>Firmicutes</b>	<b>Mollicutes</b>	Deletion of 7 residues in the loop connecting helix 2 and $\alpha$ -strand B3









