# Origin and Evolution of DNA methyltransferases (DNMT) along the tree of life: A multi-genome survey.

**Madhumita Bhattacharyya[1,2], Subhajyoti De[3] and Saikat Chakrabarti[1]\*.**

[1]Structural Biology and Bioinformatics Division, Council of Scientific and Industrial Research-Indian Institute of Chemical Biology, CN 6, Sector V, Kolkata 700091, West Bengal, India

[2]Current address: Chair and Institute of Environmental Medicine (IEM), UNIKA-T, Helmholtz Zentrum Munich and Technical University Munich, Augsburg, Germany

[2]Rutgers Cancer Institute of New Jersey, Rutgers University, New Brunswick, NJ, USA

\* Author for correspondence

Email: saikat@iicb.res.in, saikat273@gmail.com,

1

17

## Abstract:

**Background:** Cytosine methylation is a common DNA modification found in most eukaryotic organisms including plants, animals, and fungi. (Cytosine-5)-DNA methyltransferases (C5-DNA MTases) belong to the DNMT family of enzymes that catalyze the transfer of a methyl group from S-adenosyl methionine (SAM) to cytosine residues of DNA. In mammals, four members of the DNMT family have been reported: DNMT1, DNMT3a, DNMT3b and DNMT3L, but only DNMT1, DNMT3a and DNMT3b possess methyltransferase activity. There have been many reports about the methylation landscape in different organisms yet there is no systematic report of how the enzyme DNA (C5) methyltransferases have evolved in different organisms.

**Result:** DNA methyltransferases are found to be present in all three domains of life. However, significant variability has been observed in length, copy number and sequence identity when compared across kingdoms. Sequence conservation is greatly increased in invertebrates and vertebrates compared to other groups. Similarly, sequence length has been found to be increased while domain lengths remain more or less conserved. Vertebrates are also found to be associated with more conserved DNMT domains. Finally, comparison between single nucleotide polymorphisms (SNPs) prevailing in human populations and evolutionary changes in DNMT vertebrate alignment revealed that most of the SNPs were conserved in vertebrates.

**Conclusion:** The sequences (including the catalytic domain and motifs) and structure of the DNMT enzymes have been evolved greatly from bacteria to vertebrates with a steady increase in complexity and specificity. This study provides a systematic report of the evolution of DNA methyltransferase enzyme across different lineages of tree of life.

39

**Keywords:**

DNA methyltransferase, Tree of Life, Phylogenetic analysis, DNMTs, SNPs

**Background:**

The genomes of eukaryotes are marked with regionally restricted epigenetic information responsible for regulating local activity states. Most widely studied epigenetic modification in humans is cytosine methylation which occurs almost exclusively in the context of CpG dinucleotide. The CpG dinucleotides tend to cluster in regions called CpG islands (Bird A. et. al. 2002). About 60% of human gene promoters are associated with CpG islands and are usually unmethylated in normal cells while some of them (~6%) become methylated in a tissue-specific manner during early development or in differentiated tissues (Brown and Strathdee, 2002). In general, CpG-island methylation is associated with gene silencing, genomic imprinting, X chromosome inactivation in females, histone modification, chromatin remodeling etc. DNA methylation and DNA methylation–associated proteins not only participate in gene transcription regulation in *cis*, but also act in *trans*, being involved in nuclear organization and in the establishment of specific chromosomal territories. Hypermethylated CpGs are needed to protect chromosomal integrity, which is achieved by preventing reactivation of endoparasitic sequences that cause chromosomal instability, translocations and gene disruption (Okano M. et. al. 1999).

DNA methylations in mammalian systems are observed throughout the genome barring the CpG islands (CGIs) (Bird A. 2002, Bird A. P. 1986). In fungi that have genomic 5‑methylcytosine (m5C), only repetitive DNA sequences are methylated (Selker E. U. et. al, 2003). The most frequent pattern observed in invertebrates is 'mosaic methylation', comprising domains of

3

62    heavily methylated DNA interspersed with domains that are methylation free (Bird A. P. et. al.

63    1979, Tweedie S. et. al. 1997). The highest levels of DNA methylation among all eukaryotes

64    have been observed in plants, with up to 50% of cytosine being methylated in some species

65    (Montero L. M. et. al, 1992). In maize, for example, such high levels seem to be due to large

66    numbers of transposons, the degenerate relics of which dominate inter-genic regions and are

67    targeted for methylation (Palmer et. al. 2003, SanMiguel P. et. al. 1996). However, other plants,

68    such as *Arabidopsis thaliana*, display a mosaic DNA methylation pattern that is reminiscent of

69    invertebrate animals.

70    Although DNA methylation appears to be a widespread epigenetic regulatory mechanism,

71    genomes are methylated in diverse ways in different organisms. In animals, DNA methylation

72    occurs mostly symmetrically (both strands) at the cytosines of a CG dinucleotide. DNA

73    methylation in plant genomes can occur symmetrically at cytosines in both CG and CHG (H = A,

74    T, or C) contexts, and also asymmetrically in a CHH context, with the latter directed and

75    maintained by small RNAs (Law and Jacobson, 2010). In the model plant *Arabidopsis thaliana*,

76    levels of cytosine methylation at CG, CHG, and CHH nucleotides are about 24%, 6.7%, and

77    1.7%, respectively (Cokus S. J. et. al. 2008, Lister R. et. al. 2008). However, it is important to

78    bear in mind that the global DNA methylation pattern seen in vertebrates is by no means

79    ubiquitous among eukaryotes. Several well-studied model systems have no recognizable *DNMT*

80    like genes and are devoid of DNA methylation (for example, the yeast *Saccharomyces*

81    *cerevisiae,* fruit fly *Drosophilla melanogaster* and the nematode worm *Caenorhabditis elegans*).

82    Despite the different methylation sequence contexts, cytosine methylation is established and

83    maintained by a family of conserved DNA methyltransferases (Goll and Bestor 2005, Chan S.

84    W. et. al 2005, Cheng and Blumenthal 2008). Not surprisingly, the absence of DNA methylation

85    in some eukaryotes such as yeast, fruit fly, and roundworm is associated with the evolutionary

86    loss of DNA methyltransferase homologs (Goll and Bestor 2005). Different C5-cytosine

87    methyltransferases have been characterized in prokaryotes and eukaryotes. All

88    methyltransferases share a catalytic domain containing 10 conserved small motifs, suggesting a

89    common origin (Posfai J. et al. 1989; Kumar S. et al. 1994). Based on sequence similarity, the

90    eukaryotic methyltransferase have been grouped into different subfamilies at different times,

91    based on different criteria. For example, methylases were classified based of the methylation

92    residue (M4C, M6A, and M5C), methylation activity ("*de novo* methylation" or "maintenance

93    methylation"), and methylation state of substrate DNA (unmethylated / hemimethylated).

94    M6A and M4C methyltransferases are responsible for methylation of adenine residue (at N6) and

95    cytosine residue (at N4), respectively. Both of these enzymes are primarily found in prokaryotes

96    where it functions as restriction endonuclease when DNA is unmethylated and functions as

97    methyltransferase when DNA is hemimethylated. M5C methyltransferases are responsible for

98    methylation in cytosine residue (at N5) position and are found to have a large C-terminal

99    catalytic domain and smaller N- terminal recognition domain.  Until recently only one DNA

100    methyltransferase, *DNMT1*, had been cloned from human and mouse cells. Recently, another

101    group of DNA methyltransferases, *DNMT3A* and *DNMT3B* was isolated by database search.

102    DNMT1 is the most abundant DNA methyltransferase in mammalian cells, and considered to be

103    the key maintenance methyltransferase in mammals (Okano M. et. al. 1999). The carboxy-

104    terminal catalytic domain of DNMT1 is responsible for transfer of methyl groups from S-

105    adenosyl methionine to cytosines in CpG dinucleotides. The longer N-terminal portion of

106    DNMT1 contains regions responsible for targeting the replication foci. A cysteine-rich Zn-

107    binding motif and a polybromo domain that resemble regions of proteins known to have roles in

108  chromatin modification (Leonhardt H et. al. 1992, Chuang L. S. et. al. 1997) are utilized for

109  discriminating between unmethylated and hemi-methylated DNA,

110  DNMT3 is a family of DNA methyltransferases that could methylate hemimethylated and

111  unmethylated CpG at the same rate. The architecture of DNMT3 enzymes is similar to that of

112  DNMT1, with a regulatory region attached to the catalytic domain. There are three known

113  members of the DNMT3 family: 3A, 3B, and 3L. DNMT3A and DNMT3B can mediate

114  methylation-independent gene repression. DNMT3A can co-localize with

115  heterochromatin protein (HP1) and methyl-CpG-binding protein (MeCBP). They can also

116  interact with DNMT1, which might be a co-operative event during DNA methylation. DNMT3A

117  prefers CpG methylation to CpA, CpT, and CpC methylation, though there appears to be some

118  sequence preference of methylation for DNMT3A and DNMT3B. DNMT3L contains DNA

119  methyltransferase motifs and is required for establishing maternal genomic imprints, despite

120  being catalytically inactive. DNMT3L is expressed during gametogenesis when genomic

121  imprinting takes place. The loss of DNMT3L leads to bi-allelic expression of genes normally not

122  expressed by the maternal allele. DNMT3L interacts with DNMT3A and DNMT3B and co-

123  localizes in the nucleus (Okano M. et. al. 1999).

124  All reported DNMTs have a well conserved C terminal catalytic domain containing an ordered

125  set of sequence motifs which alternates with non-conserved regions. Depending on the criteria

126  used to define the limits of the conserved blocks, up to ten motifs can be identified (Lauster R.

127  et. al. 1989, Posfai J. et. al. 1989, Klimasauskas S. et. al. 1989). In the original analysis of 13

128  M5C-methyltransferases, five motifs were considered highly conserved (I, IV, VI, VIII, and X),

129  and the remaining five moderately conserved. Analysis of 36 sequences resulted in the inclusion

130  of a sixth motif, motif IX, within the highly conserved set (Cheng X. et. al. 1993).

131     Different methyltransferases were reported in different organism and few attempts were made to

132     explore their origin and evolution. Based on these observations, Colot and Rossignol (1999)

133     proposed that methylation has divergent functions in different organisms, consistent with the

134     notion that it is a dynamically evolving mechanism that can be adapted to perform various

135     functions. Another report was published in 2005 where Ponger and Li tried to analyze the

136     variability of methylation systems, with a survey of methytransferases in complete or almost

137     complete eukaryotic genomes, including several species of Protozoa. They also reconstructed a

138     phylogeny of the putative enzymes identified to study the evolutionary history of this function

139     and to classify eukaryotic methyltransferases (Ponger and Li, 2005). Functional and structural

140     conservation of DNMTs in human, mouse and cattle were observed along with similar patterns

141     of transcript abundance for all of the proteins at different stages of early embryo development.

142     Greater degree of structural similarity between human and bovine was observed for all of the

143     DNMT (DNMT1, DNMT3A, DNMT3B, and DNMT3L) than that between human and mouse

144     (Rodriguez-Osorio N. et. al. 2010).

145     Using the information provided by National Center for Biotechnology Information (NCBI)

146     taxonomy, a phylogenetic tree was reported in a science perspective where methylation extent

147     and type of methyltransferases were put together. The report suggested that the  last common

148     ancestor of eukaryotes contained a functional DNA methylation system with secondary

149     expansion and the loss of methylation in some lineages while primitive methylation likely

150     occurred at low to intermediate levels and was targeted to gene bodies and transposable

151     elements, leaving gene promoters unmethylated (Jeltsch A, 2010). Zemach et. al. reported a

152     quantitative estimation of DNA methylation in 17 species and found out that gene body

153     methylation is conserved on the contrary to selective transposon methylation (Zemach A. et. al

7

154     2010). Another publication from the same group reported transposable elements and sex to be

155     the major forces driving the evolution of methylation (Zemach and Zilberman, 2010). Shotgun

156     bisulfite sequencing (BS-seq) was used to compare DNA methylation in eight diverse plant and

157     animals. Different patterns of methylation were detected in flowering plants and animals (Feng

158     S. et. al 2010). In one of the largest phylogenetic analysis till date comprising of 2300 sequences

159     the evolutionary relationship among different DNMT1 and DNMT3 was explored. This study

160     proposed a consensus model of the phylogeny of DNA methyltransferases indicating that

161     DNMT1 and DNMT3A/B enzymes have an independent origin in the prokaryotic DNA

162     methyltransferase sequence space and all were derived from methyltransferases of restriction

163     modification systems (Jurkowska and Jeltsch, 2011).

164     DNA methylation machinery was always a central question of interest and a good number of

165     reports are available about identification and diversification of DNA methylation machinery in

166     higher order organisms and/or having small number of representatives from different

167     taxonomical groups. However, a comprehensive overview of modifications in sequence and

168     structural level over time and taxonomical hierarchy is still lacking. This study attempts to make

169     a detail investigation of evolution of DNA methyltransferase enzymes along all branches of tree

170     of life (TOL) including as many possible organisms from lower (bacteria) to higher (human)

171     order of taxonomy. The objectives of this work are to search and identification of the

172     homologues for DNA methyltransferase in all available genomes followed by a thorough and

173     systematic comparison of the sequences in order to elucidate the evolution of DNMTs along the

174     lineages of tree of life. Finally, naturally occurring variations were compared with SNPs

175     observed in human populations to investigate if the evolutionary selection pressure on structural

176    and functional motifs of DNMT was similar in genetic variations observed in human

177    populations.

178    To reconstruct the tree of life based on DNA methyltransferase using phylogenetic methods,

179    sequences were classified systematically from seven kingdoms. Parameters like gene copy

180    number, sequence length and identity were compared across the kingdoms. Functional motifs

181    and domains were identified and compared across different lineages of the tree of life.

182    Conservation of each site/residue was measured along each kingdom group and finally compared

183    with the available SNP data. All three groups of sequences were found to have a clear

184    evolutionary pattern across kingdoms. Kingdom specific conserved functional motifs and

185    domains were also observed. Moreover, to the best of our knowledge, this is the first large scale

186    systematic report of evolution of DNMT enzymes across different lineages of tree of life

187    including all possible organisms. A clear pattern of evolution was observed where the sequences

188    of the enzyme achieved more complexity and specificity in higher order organisms. Comparison

189    with single nucleotide polymorphism data in human shows that none of the SNPs were

190    overlapping with functional motifs though more than 60% are conserved residues, while more

191    than 80% SNPs occurring in non-conserved residues are tolerated.

192

193    **Results**

194    **Collection and Classification of DNMT sequences**

195    Cytosine (5) methyltransferase domain (Pfam ID: PF00145) was identified from Pfam database

196    (Punta M. et. al. 2010) and was scanned to search similar sequences against different sequence

197    databases like NCBI non redundant database (Coordinator N R 2018) and Uniprot (Leinonen R.

9

198    et. al. 2004). Total 4845 unique sequences were identified in 710 organisms from all three

199    databases.

200    For each type of DNA methyltransferases, (namely DNMT1, DNMT3A, DNMT3B and

201    DNMTL) full length, annotated and experimentally reviewed sequences were collected from

202    Uniprot (Leinonen R. et. al. 2004). 55 DNMT1 unique sequences were identified whereas 17, 26

203    and 5 unique sequences were extracted for DNMT3A, DNMT3B and DNMTL, respectively. For

204    each type of DNMT enzyme a hidden markov model (HMM) profile was created **(See Method)**

205    and each of the collected 4845 sequences was scanned against these signature DNMT profiles.

206    Finally each sequence was classified to a particular DNMT group based on its alignment

207    matching with the respective DNMT HMM profile. From the initial 4845 sequences, 1783

208    sequence were found to be DNMT1, whereas only 67 and 172 sequences were identified to be

209    DNMT3A and DNMT3B, respectively. No DNMTL sequences were identified. The protocol for

210    sequence identification is described as a flow chart in **Figure S1**.

211    In this study 16S rRNA based tree of life (Yarza P. et. al. 2008) was used as a reference where

212    all organisms were grouped into three super kingdoms namely archea, bacteria and eukarya.

213    However, for broader understanding, all the DNMT containing organisms were grouped into

214    seven kingdoms such as archea, bacteria, algae, fungi, plants, invertebrates and vertebrates. Each

215    genus was classified up to phylum level and very large phyla like ascomycota/basidiomycota or

216    angiospermata were classified up to class level. The whole classification is presented in **Tables**

217    **S1, S2 and S3.** Most prevalent bacterial organisms were from phylum proteobacteria and

218    actinobacteria. Arthropods were found to be most prevalent DNMT1 containing phyla. Most

219    prevalent vertebrates were rodent and primate mammals; other mammals were also present in the

220    classification. Most prevalent DNMT3A and DNMT3B containing organisms were mammals.

10

221   DNMT3B is also present in plants and invertebrates. Presence of both of these enzymes only in

222   higher order organisms suggests that they probably have originated much later than DNMT1.

223

224   **Phylogeny of DNMTs along tree of life**

225   In order to trace the evolutionary history of DNMT enzymes along different branches of tree of

226   life the 710 organisms were classified into seven major kingdoms as mentioned above. No

227   archea, algae and fungi were found to possess sequences of DNMT3A and DNMT3B. Only 15

228   plants were detected with DNMT3B sequences but no DNMT3A sequences.   Among 292

229   bacteria only 7 bacteria contain DNMT3A sequences and 32 genuses contain DNMT3B

230   sequences. Similar distribution was observed in invertebrates where only one organism was

231   detected with DNMT3A sequence and nine organisms were detected with DNMT3B sequences.

232   On the contrary, most of the vertebrates (75%) were detected with sequences of all three DNMT

233   enzymes. Distribution of the organisms is presented in **Table 1.** A matrix was also created to

234   represent the number of unique organisms having any one, two or all of the three enzymes

235   **(Figure S2)**. In 2006 Ciccarelli proposed one of the most extensive tree of life which spans

236   across 191 genomes (Ciccarelli F. D. et. al. 2006). For detection of evolutionary lineage, 16

237   protein families were considered to build this huge tree. The presence of DNMT enzymes were

238   mapped on the same tree. Only 134 organisms were common between both the dataset **(Figure**

239   **1)**.

240

241   **Table 1: Number of Different DNMT containing organisms in various phylogenetic groups.**

242

| Kingdom level of organisms | Total number of organisms | Number of organisms containing DNMT1 | Number of organisms containing DNMT3A | Number of organisms containing DNMT3B |
|---|---|---|---|---|
| Bacteria | 292 | 281 | 7 | 32 |
| Archea | 30 | 30 | 0 | 0 |
| Algae | 5 | 5 | 0 | 0 |
| Fungi | 46 | 46 | 0 | 0 |
| Invertebrates | 25 | 23 | 1 | 9 |
| Plants | 26 | 26 | 0 | 15 |
| Vertebrates | 40 | 37 | 31 | 31 |

243

244    All three DNMT sequences were subjected to multiple sequence alignment (MSA) followed by a

245    phylogenetic tree creation. Total 448 organisms were found to contain 1773 sequences of

246    DNMT1 sequences. Longest sequence from each organism was selected for the MSA. DNMT1

247    sequences were identified in all kingdoms of tree of life. In the tree, higher organism like

248    vertebrate and invertebrates were clustered together, whereas in plants, fungi and algal species

249    DNMT sequences are more closely related. DNMT1 in certain archea groups were more closely

250    related to bacteria than to themselves **(Figure 2).** On the contrary DNMT3A sequences were

251    identified only in vertebrates and bacteria. DNMT3B sequences were detected in bacteria, plants,

252    invertebrate and vertebrates. Two important observations from this analysis were (i) unlike r-

12

253 RNA based tree of life few of the bacteria and archea are closely related, (ii) eukaryotic

254 kingdoms are well clustered among themselves and distantly related from one another.

255

256 **Variation in DNMT sequences among different organisms**

257 Collection and phylogenetic analysis of the sequences were followed up by analysis of variation

258 across seven major kingdoms of life. The variation was investigated using parameters like gene

259 copy number, length of protein, sequence identity etc. Many organisms were identified with

260 multiple copies of DNMT sequences. Highest copy number of DNMT1 genes (97) was found in

261 *Helicobacter sp.,* a proteobacteria. Many other bacteria (101) were found to possess 3 or more

262 copies of DNMT genes. *Clostridium sp.* was observed to contain as high as 28 copies of DNMT1

263 and 12 copies of DNMT3B. Most of the algal and fungal species were found to contain two to

264 four copies of DNMT1 and DNMT3B. On the contrary, in higher order organism like mammals

265 one or two copies of genes were identified. Zebra fish *Danio rerio* was observed to have highest

266 copy number of all three enzymes having six copies of DNMT1, seven copies of DNMT3A and

267 twelve copies DNMT3B. Few plants like *Oryza sativa* and *Vitis vinifera* were observed with

268 DNMT1 copy number as high as 17. The distribution of copy numbers different kingdoms are

269 presented as a box whisker plot in **Figure 3A**.

270 Many organisms from the same group were observed to have DNMT enzymes of different

271 lengths. The mean length of whole enzyme increases with the evolutionary hierarchy. Average

272 length of DNMT enzymes in bacteria is about 424 amino acids whereas algal species contain

273 very long (>= 2000) DNMT1 sequences. Average length of vertebrate DNMT1 sequences is

274 about 1452 amino acids while the same of DNMT3A and DNMT3B are about 834 and 970

275 amino acids. The distribution of sequence lengths are plotted in **Figure 3B.**

276   From each DNMT1 full-length sequence, Cytosine specific (C5) DNA methyltransferase domain

277   was extracted for length comparison. Interestingly, average length of the DNMT domain in

278   DNMT1 sequences also increases with evolutionary hierarchy while the same in DNMT3A

279   follows the opposite trend. Average lengths of DNMT domain from DNMT3B sequences are

280   very similar in different kingdoms. The distributions of domain lengths are presented in **Figure**

281   **3C.** Another interesting observation was that the variability in the length of whole enzyme is

282   higher than the length of the domain.

283   From all the sequences the DNA methylase domains were extracted and aligned pairwise in all-

284   to-all combinations in order to identify the overall conservation of DNA methylase domain

285   across each group. Though functionally all of them are responsible for methyl transfer to C5

286   position yet the sequence identities were low in most of the groups especially in bacteria, algae

287   and fungi. Interestingly, archea DNMT1 are relatively more conserved (51%) similar to that of

288   plants and invertebrates.  However, vertebrate DNMT1 sequences were found to be very highly

289   conserved (80%) compared to other groups. Interestingly, both the bacterial and vertebrate

290   DNMT3A sequences were also found to be quite highly conserved (67%) compared to

291   DNMT3B where bacterial sequences possess significantly lower conservation pattern compared

292   higher organisms **(Figure 4)**.

293

294   **Variation of DNA methyltransferase motifs and associated domains**

295   DNA methyltransferase domain is known to have ten small motifs. All of the motifs play

296   significant role in DNA binding, Adomet binding and catalytic activity. Motif IV (PCQ) is the

297   catalytic motif. All these motifs were reported to be present in vertebrate DNMT1, DNMT3A

298   and DNMT3B. As overall sequence identities were found to be low in the DNMT sequences,

14

299    conservation pattern of individual motifs was examined in each of the enzyme class across each

300    kingdom **(Figure S3)**.

301    All the DNMT sequences from each organism were subjected to a MSA followed by a motif

302    scanning protocol. Though not all the motifs are present in all kingdoms and in all enzyme

303    groups, yet six of them are found to be present in every kingdom with little or no mutation.

304    Motifs I, which is responsible for Adomet (the methyl doner) binding, and motif IV, which is the

305    catalytic motif, were found to be present in all organism groups. Motif VIII of DNMT3A and

306    DNMT3B was not present in vertebrate sequences. All motifs except motif I and Motif IV are

307    absent in DNMT3B sequences of plants. A new motif (**R**x**R**) was identified in our analysis in the

308    DNMT1 and DNMT3A and DNMT3B sequences in all organism groups except plant DNMT3B

309    sequences. Though the motif residues were conserved yet some changes were observed across

310    the organisms. The glutamine (Q) in the catalytic motif IV has been replaced by asparagine (N)

311    in many of the DNMT3A and DNMT3B sequences. Also in case of motif I (**FxGxG**) the last

312    Glycine have been observed to have higher rate of mutation in DNMT3A and DNMT3B. Also

313    motif I in vertebrate DNMT1 sequences (**FSGCG**) is changed into **FDGIA** in DNMT3A and

314    DNMT3B sequences. Motif X in DNMT1 sequences (GN) is also replaced by SN in many of the

315    DNMT3A and DNMT3B sequences. The logo plots of motifs are presented in **Figure S3**

316    whereas an estimation of the diversity within the motifs is shown in **Figure S4**.

317    All DNMT sequences were scanned against Pfam database in order to map annotated protein

318    domains in the sequences. About 275 other domains (including 61 domains with no known

319    function) were found to be present in DNMT1 sequences but only 37 of them were present in

320    more than 5 sequences **(Figure 5A)**. In DNMT3A sequences only one domain (PWWP) was

321    found to be present in almost all sequences. Though 33 different domains other than DNA

322     methylase domain were mapped on DNMT3B sequences yet only six domains were mapped

323     onto more than 10 DNMT3B sequence.

324     All the domains were mapped onto DNMT1, DNMT3A and DNMT3B phylogenetic trees

325     mentioned above. Almost all organism groups were observed to contain unique domain

326     organization. Bacterial and archeal DNMT1 sequences were found to be comparatively shorter

327     and methyltransferase domains were present at the N-terminal. Many DUF (Domain with

328     Unknown Functions) were also mapped onto different bacterial DNMT1 sequences.  Algal and

329     fungal sequences are comparatively longer are mapped with single or double BAH domain and

330     DNMT1-RFD domain. Fungal DNMT1 contained long stretch of sequences with no mapped

331     domains. Interestingly, both vertebrate and invertebrate sequences were found to have multiple

332     other domains. Many invertebrate and all vertebrate sequences were mapped with DMAP1

333     binding, DNMT1-RFD, zf-CXXC, two consecutive BAH domain and DNA methylase domain

334     from N to C terminal **(Figure 5A)**. All DNMT3A sequences were mapped with only two

335     domains PWWP at the N terminal and DNA methylase at the C-terminal **(Figure 5B)**. Most of

336     the DNMT3B bacterial sequences were mapped with only DNA methylase domain whereas

337     UBA domain was mapped in many plant sequences. All invertebrate and vertebrate sequences

338     were mapped with PWWP domain. Another interesting observation was that in both DNMT3A

339     and DNMT3B there is a large insertion inside the DNA methylase domain **(Figure 5C)**.

340

341     **Single nucleotide polymorphisms in DNMTs**

342     A list of reported single nucleotide polymorphisms (SNPs) for all three human DNMT1,

343     DNMT3A and DNMT3B were compiled from ensemble transcripts sequences (**Table S4**).

344     Mutation in DNMTs could be very crucial for different diseases. In order to investigate whether

16

345    the functionally important SNP positions are conserved in vertebrates, all SNPs were mapped

346    over vertebrate specific alignment of DNMT1, DNMT3A and DNMT3B. In all three cases

347    percentage of total SNP positions that are not conserved was found to be connected with

348    deleterious effect, whereas most of the SNPs which are deleterious in effect were found to be

349    more conserved **(Figure 6).** The distribution of SIFT score, which indicates the deleterious effect

350    of a SNP and the conservation index (measured by AL2CO scores) was compared together to

351    identify both deleterious and conserved SNPs **(Figure 7)**. Mutations map was created for all

352    three enzymes. It was observed that most of the conserved and deleterious SNPs were result of

353    mutation replacement of polar uncharged residues with hydrophobic residue or vice versa. On

354    the contrary, SNPs without any deleterious effect were resulted from replacement with similar

355    residues **(Figure 7).**

356

## Discussion:

358    In this study, evolutionary studies of three DNA (C5) methyl transferase enzymes namely

359    DNMT1, DNMT3A and DNMT3B are performed. This is one of the few large scale study

360    considering about five hundred organisms and about two thousand sequences. This was also the

361    first attempt made towards exploring the landscape of evolutionary history of an extremely

362    critical and important enzyme in different classes of organism. This study shows the presence of

363    cytosine specific methyltransferase in many primitive organisms like bacteria and archea, which

364    contradicts the common idea that cytosine specific methylation, is a regulatory mechanism

365    present only in higher order organisms. Rather, it has been found to be a global phenomenon

366    with conspicuous absence in few organisms like yeast, worm and fruit fly.

17

367      DNA methylation was found to play important roles in the biology of bacteria: phenomena such

368      as timing of DNA replication, partitioning nascent chromosomes to daughter cells, repair of

369      DNA, and timing of transposition and conjugal transfer of plasmids are sensitive to the

370      methylation states of specific DNA regions. All of these above mentioned events use the hemi-

371      methylated state of newly replicated DNA as a signal. In the case of DNA replication, the protein

372      SeqA binds preferentially to hemi-methylated DNA target sites (GATC sequence) clustered in

373      the origin of replication (*oriC*) and sequesters the origin from replication initiation. In addition,

374      SeqA also transiently blocks synthesis of the DnaA protein, which is necessary for replication

375      initiation, by binding to hemi-methylated GATC sites in the *DnaA* promoter. In DNA repair, the

376      methyl-directed mismatch repair protein MutH recognizes hemi-methylated DNA sites and cuts

377      the non-methylated daughter DNA strand, ensuring that the methylated parental strand will be

378      used as the template for repair-associated DNA synthesis. DNA methyltransferases in bacteria

379      were best understood in the context of restriction–modification (R–M) systems, which act as

380      bacterial immune systems against incoming DNA including phages. But several orphan

381      methyltransferases, which were not associated with any restriction enzyme, have also been

382      characterized and may protect against parasitism by R–M systems. Interestingly, in a report by

383      Sandip Krishna in 2012 this orphan methyltransferases were found to be more conserved than R-

384      M methyltransferase. In our study, we have identified nearly 1400 enzymes across 450 genomes

385      of bacteria.

386      Apart from bacteria, methylation mediated regulation was poorly understood in algae and fungi.

387      A firsthand report of algal and fungal DNMTs along with their structural features and

388      conservation pattern is provided by this study. In plants, methylation of cytosine bases was found

389      in all sequence contexts: the symmetric CG and CHG contexts (where H is A, T, or C) and the

390    asymmetric CHH context. Specific enzymes were reported that establish and successively

391    maintain methylation patterns during DNA replication. It was suggested that methylation occurs

392    predominantly at repeats and transposons (more than 90% are methylated), but approximately

393    the 20% of genes also exhibit a certain degree of methylation. Overall, the levels of methylation

394    in the *Arabidopsis thaliana* genome at CG, CHG, and CHH are about 24%, 6.7%, and 1.7%,

395    respectively, but methylation within the genes is primarily restricted to CG sites and was

396    predominantly observed in the transcribed coding region or the so called gene body (Kokus S. J,

397    et. al 2008 and Lister R. et. al 2008). It was also reported that modestly expressed genes are more

398    likely to be methylated within gene body, while genes expressed at high and low levels are

399    usually less methylated. In our study it was found that DNMTs of plants acquire a certain

400    conserved sequence with all the conserved motifs. It was also reported that Plants sequence were

401    more closely related to the Algal and Fungal sequences than which corresponds to the protein

402    based or RNA based tree of life (See Figure 2).

403    However, the sequence and structure of the enzyme has been evolved greatly from bacteria to

404    vertebrates. Great variability has been observed in length, copy number and sequence identity of

405    DNA (C5) methyltransferase domain when compared across kingdoms. It has been observed that

406    the sequence conservation is greatly increased in invertebrates and vertebrates in comparison to

407    other groups. Sequence length has been found to increase while domain lengths remain more or

408    less conserved with the evolution indication association of other functional domains. By domain

409    analysis it was found that vertebrates not only have a conserved methylation domain but also

410    have other conserved domain which aid in the methyltransferase function. No such signature

411    association of domains was observed. As a whole this study reports a history of evolution of

412    DNA methyltransferase enzyme across different leaves of tree of life. There is a future scope to

19

413    trace the origin of this enzyme if a phylogenetic analysis is performed including DNMTs and

414    other methyltransferase enzymes like RNMT, DAM and MGMTs.

415    Comparison between naturally occurring SNPs and residues conserved in vertebrate alignment

416    revealed that most of the tolerated SNPs are conserved in vertebrates. It was also observed that

417    mutations caused by residues with similar chemical property were more tolerated than the other

418    way round. In future the comparison of SNPs can be extended to alignment of all organisms to

419    see similar selection pressure exists throughout all branches of tree of life.

420

## Methods:

421    **Methods:**

422    **Sequence collection**

423    Sequences of DNA (C5) methyltransferase domain (PF00145) were extracted from Pfam and

424    subjected to homology search using BLAST (Altschul S. F. et. al. 1990) against NR (Coordinator

425    N. R. 2018) and Uniprot (Leinonen R. et. al. 2004) sequence datasets. Homologues were

426    identified based on a threshold E-value of $<= 10^{-5}$, query coverage $>= 50\%$ and subject coverage

427    $>= 50\%$.

428    All reviewed full length sequences of DNMT1, DNMT3A, DNMT3B and DNMTL were

429    collected from Uniprot (Leinonen R. et. al. 2004). HMM-profile was made using MAFFT and

430    HMMER for each of the four enzymes. Each sequence of DNMT was scanned against these

431    profiles using *Hmmersearch* (Johnson L. S. et al. 2010)). Again the sequences were identified

432    using a threshold of E-value $<= 10^{-5}$, query coverage $>= 50\%$ and subject coverage$>= 50\%$.

433    Unique sets of homologous sequences were identified using an in house alignment scoring

434    method.

435

**Obtaining the tree of life**

437 Here, two trees of life have been used as references to detect the evolutionary relationship among

438 DNMTs. First one was the 16S r-RNA based tree of life according to which evolution occurred

439 in three different branches of life namely archea, bacteria and eukarya.

440 The second one was proposed by Ciccarelli in 2006 (Ciccarelli F. D. et. al. 2006). This tree was

441 constructed using 236 protein families spanning over 191 genomes. With further investigation it

442 was found that the tree has 112 unique genuses among which 71 are bacteria, 23 are eukaryotes

443 and 18 are archea. Keeping this in mind, all the organism found to containing DNMTs were

444 classified into seven groups; along with archea and bacteria, eukarya was classified into 5 more

445 groups such as algae, fungi, plants, invertebrates and vertebrates.

446

**Phylogenetic tree construction**

448 A common protocol was followed for construction of all phylogenetic trees in this study.

449 Sequence set for tree construction was identified and redundancy was removed using CD-HIT at

450 100% (Li and Godzik 2006). Multiple sequence alignments (MSA) were created using MAFFT

451 5.1 (Kattoh K. et. al. 2002). Phylogenetic tree was constructed by RAXML-HPC 7.0.4

452 (Stamatakis A. 2006) package using maximum likelihood method (Bootstraping value was set as

453 100). In all seven groups DNA methylase domain was marked and extracted for tree

454 construction. Tree images were created using ITOL (Interactive Tree of Life) server (Letunic and

455 Borc, 2007).

456

**Calculation of sequence identity**

458     Sequences from each 7 kingdom was subjected to all-to-all pairwise alignment using NEEDLE

459     tool from EMBOSS package (Rice P. et. al. 2000) and identity matrices were created using in-

460     house programs. The numerical matrices were converted to colour matrices by using Matrix2png

461     1.0.6 tool (Pavlidis and Noble 2003).

462

463     **Scanning and identification of motifs**

464     Sequence conservation was measured in all the kingdom level MSAs by AL2CO software

465     package (Pei and Grishin 2001). Scores for all columns were sorted and subjected to statistical

466     analysis. All columns with positive scores were identified as conserved column which then

467     divided into three classes. The columns in the top quartile of the distribution were considered as

468     highly conserved, the columns belonging to the second and third quartiles were termed as

469     moderately conserved and the columns in last quartile were marked as conserved columns. All

470     signature motifs were identified in conserved columns of the alignment and converted to logo

471     plot using WEBLOGO 3.3 (Crooks G.E. et. al. 2004). The conservation score of each motif was

472     calculated by an in-house Perl program.

473

474     **Scanning and identification of domains**

475     All the sequences were scanned against Pfam database (Punta M. et. al. 2012) using *hmmsearch*

476     by HMMER 2.1 (Johnson L. S. et al. 2010). Domains were identified with the threshold of E-

477     value $10^{-5}$ and subject coverage of >=50%. Position and combination of domain in each sequence

478     were identified using in-house perl program. The domains were mapped onto the phylogenetic

479     trees using ITOL server v2 (Letunic and Bork, 2006,).

480

481     **Collection of SNP data**

22

482    Single nucleotide polymorphisms (SNPs) for all three human DNMT1, DNMT3A and DNMT3B

483    were obtained from the Ensembl Genome Browser database (Hunt S. E. et al., 2018). This

484    represented the pooled list of SNPs obtained from the dbSNP, ClinVar and 1000 Genomes

485    Project that were mapped onto the Ensembl transcripts of DNMT1, DNMT3A and DNMT3B.

486

487    **Calculation of SIFT score**

488    SIFT (Sorting Intolerant From Tolerant) uses sequence homology to predict whether an amino

489    acid substitution will affect protein function and hence, potentially alter phenotype. For

490    calculation of  score a query protein is searched against a protein database to obtain homologous

491    protein sequences. Sequences with appropriate sequence diversity are chosen. The chosen

492    sequences are aligned, and for a particular position, SIFT looks at the composition of amino

493    acids and computes the score. A SIFT score is a normalized probability of observing the new

494    amino acid at that position, and ranges from 0 to 1 (NG and Henikoff, 2001, NG and Henikoff

495    2002, Sim et. al. 2012).

496

497    **List of abbreviations**

498    DNMT: DNA methyl transferases

499    RNMT: RNA methyl transferase

500    DAM: Deoxy Adenosine Methylase

501    MGMT: Methylgualine (O6) DNA methyltransferase

502    SNP**:** Single Nucleotide Polymorphism

503    SIFT: Sorting Intolerant From Tolerant

504    iTOL: Interactive Tree of Life

505    **Declarations**

506    **Ethics approval and consent to participate:**

507    Not applicable

508    **Consent for publication**

509    **Availability of Data**

510    Not applicable

511    **Competing Interest**

512    The authors declare that they have no competing interests.

513

514    **Funding**

517

518    **Authors' contributions**

519    MB and SC have designed computational experiments. SD has provided the experimental data.

520    MB has performed experiments. MB, SD and SC has analyzed the results. MB, SD and SC has

521    written the manuscript.

522    **Acknowledgement**

527

24

528

529

530

# References:

532      Altschul, S.F., Gish, W., Miller, W., Myers, E.W. & Lipman, D.J. (1990) Basic local alignment search tool. J. Mol.

533      Biol. 215:403-410 .

534      Bird, A. P. (2002). DNA methylation patterns and epigenetic memory. Genes Dev. 16, 6–21.

535      Bird, A. P. (1986). CpG-rich islands and the function of DNA methylation. Nature 321, 209–213.

536      Bird, A. P., Taggart, M. H. & Smith, B. A. (1979). Methylated and unmethylated DNA compartments in the sea

537      urchin genome. Cell 17, 889–901 .

538      Brown, R. & Strathdee, G. Epigenomics and epigenetic therapy of cancer. (2002). Trends Mol. Med. 8 (Suppl.),

539      S43–S48 .

540      Chan, S.W., Henderson, I.R., & Jacobsen, S. E. (2005). Gardening the genome: DNA methylation in Arabidopsis

541      thaliana. Nat Rev Genet 6:351–360.

542      Cheng, X. & Blumenthal, R.M. (2008). Mammalian DNA methyltransferases: A structural perspective. Structure

543      16:341–350.

544      Cheng., X., Kumar, S., Posfai, J., Pflugrath, J.W. & Roberts, R.J. (1993). Crystal structure of the HhaI DNA

545      methyltransferase complexed with S-adenosyl-L-methionine Cell 74: 299-307.

546      Chuang, L.S. et al. (1997). Human DNA-(cytosine-5) methyltransferase-PCNA complex as a target for p21WAF1.

547      Science 277, 1996– 2000.

548      Ciccarelli, F.D., Doerks, T., von Mering, C., Creevey, C.J., Sne,l B., & Bork, P. (2006). Toward automatic

549      reconstruction of a highly resolved tree of life. Science. 311(5765):1283-7

550    Cokus, S.J., Feng, S., Zhang, X., Chen, Z., Merriman, B., Haudenschild, C.D., Pradhan, S., Nelson, S.F., Pellegrini,

551    M., & Jacobsen, S.E. (2008). Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA methylation

552    patterning. Nature 452:215–219.

553    Colot, V., & J. L. Rossignol. (1999). Eukaryotic DNA methylation as an evolutionary device. Bioessays 21:402–

554    411.

555    Coordinators, N. R. (2018) Database resources of the national center for biotechnology information. Nucleic Acids

556    Res, 46 (Database issue), D8.

557    Crooks, G.E., Hon, G., Chandonia, J.M., Brenner, S. E. (2004) WebLogo: A Sequence Logo Generator Genome

558    Res. 2004 Jun; 14(6): 1188–1190.

559    Feng, S., Cokus, S.J., Zhang, X., Chen, P.Y., Bostick, M., Goll, M.G., Hetzel, J., Jain, J., Strauss, S.H., Halpern,

560    M.E., Ukomadu, C., Sadler, K.C., Pradhan, S., Pellegrini, M., & Jacobsen, S. E. (2010) Conservation and divergence

561    of methylation patterning in plants and animals. Proc Natl Acad Sci U S A. 107(19):8689-94.

562    Goll, M.G, & Bestor, T.H. (2005) Eukaryotic cytosine methyltransferases. Annu Rev Biochem 74:481–514.

563    Hunt, S.E., McLaren, W., Gil, L., Thormann, A., Schuilenburg, H., Sheppard, D., Parton, A., Armean, I.M.,

564    Trevanion, S.J., Flicek, P., Cunningham, F. Ensembl variation resources. (2018). Database (Oxford). Jan 1.

565    Jeltsch, A. (2010) Molecular biology. Phylogeny of methylomes. Science. 328(5980):837-8.

566    Johnson, L.S., Eddy, S.R., Portugaly, E. (2010) Hidden Markov model speed heuristic and iterative HMM search

567    procedure. BMC Bioinformatics. Aug 18;11:431.

568    Jurkowski, T.P., & Jeltsch, A. (2011) On the evolutionary origin of eukaryotic DNA methyltransferases and Dnmt2.

569    PLoS One. 6(11):e28104.

570    Katoh, K., Misawa, K., Kuma, K., & Miyataa, T. (2002) MAFFT: a novel method for rapid multiple sequence

571    alignment based on fast Fourier transform Nucleic Acids Res. 30(14): 3059–3066.

572    Klimasauskas, S., Timinskas, A., Menkevicius, S., Butkiene, D., Butkus, V. and Janulaitis, A.A. (1989) The

573    sequence specificity domain of cytosine-C5 methylases. Nucleic Acids Res. 17: 9823-9832.

574    Kumar, S., Tamura, K., & Nei, M. (2004). MEGA3: integrated software for molecular volutionary genetics analysis

575    and sequence alignment. Brief Bioinform. 5:150–163.

576    Lauster, R., Trautner, T.A. and Noyer-Weidner, M. (1989) Cytosine-specific type II DNA methyltransferases. A

577    conserved enzyme core with variable target-recognizing domains. J. Mol. Biol. 206: 305-312.

578    Law, J.A., & Jacobsen, S.E. (2010) Establishing, maintaining and modifying DNA methylation atterns in plants and

579    animals. Nat Rev Genet 11:204–220.

580    Leinonen, R., Diez, F.G., Binns, D., Fleischmann, W., Lopez, R., & Apweiler, R. (2004) UniProt archive.

581    Bioinformatics. 20(17):3236-7.

582    Leonhardt, H., Page, A.W., Weier, H.U. & Bestor, T.H. (1992) A targeting sequence directs DNA methyltransferase

583    to sites of DNA replication in mammalian nuclei. Cell 71, 865–873.

584    Letunic, I., & Bork, P. (2007) Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and

585    annotation. Bioinformatics. 23:127–128.

586    Li, W.,  & Godzik, A. (2006) Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide

587    sequences. Bioinformatics. 22(13):1658-9.

588    Lister, R., O'Malley, R.C., Tonti-Filippini, J., Gregory, B.D., Berry, C.C., Millar, A.H., Ecker, J.R. (2008) Highly

589    integrated single-base resolution maps of the epigenome in Arabidopsis. Cell 133:523–536.

590    Montero, L. M. et al. (1992) The distribution of 5‑methylcytosine in the nuclear genome of plants. Nucleic Acids

591    Res. 20, 3207–3210.

592    Ng, P,C., & Henikoff, S. (2002) Accounting for human polymorphisms predicted to affect protein function. Genome

593    Res. 12(3):436-46.

594    Ng, P.C., & Henikoff, S. (2001) Predicting deleterious amino acid substitutions. Genome Res. 11(5):863-74.

27

595   Okano, M., Bell, D. W., Haber, D. A. & Li, E. (1999) DNA methyltransferases Dnmt3a and Dnmt3b are essential
596   for de novo methylation and mammalian development. Cell 99, 247–257.

597   Palmer, L. E. et al. (2003) Maize genome sequencing by methylation filtration. Science 302, 2115–2117.

598   Pavlidis, P. & Noble, W.S. (2003) Matrix2png: A Utility for Visualizing Matrix Data. Bioinformatics 19: 295-296.

599   Pei, J., & Grishin N.V. (2001) AL2CO: calculation of positional conservation in a protein sequence alignment.
600   Bioinformatics. 17(8):700-12.

601   Ponger, L.,  & Li, W.H. (2005) Evolutionary diversification of DNA methyltransferases in eukaryotic genomes. Mol
602   Biol Evol. 22(4):1119-28.

603   Posfai, J., Bhagwat, A. S.,. Posfai, G., & Roberts, R.J. (1989). Predictive motifs derived from tytosine
604   methyltransferases. Nucleic Acids Res. 17:2421–2435.

605   Punta, M., Coggill, P.C., Eberhardt , R.Y., Mistry, J., Tate, J., Boursnell, C., Pang, N., Forslund, K., Ceric, G.,
606   Clements, J., Heger, A., Holm, L., Sonnhammer, E.L., Eddy, S.R, Bateman, A., & Finn, R. D. (2012) The Pfam
607   protein families database. Nucleic Acids Res. 40(Database issue):D290-301.

608   Rice, P.,  Longden, I., & Bleasby, A (2000). "EMBOSS: The European Molecular Biology Open Software Suite".
609   Trends in Genetics. 16(6): 276–277.

610   Rodriguez-Osorio, N., Wang, H., Rupinski, J., Bridges, S.M., Memili, E. (2010) Comparative functional genomics
611   of mammalian DNA methyltransferases. Reprod Biomed Online. 20(2):243-55.

612   SanMiguel, P. et al. (1996) Nested retrotransposons in the intergenic regions of the maize genome. Science 274,
613   765–768.

614   Selker, E. U. et al. (2003). The methylated component of the Neurospora crassa genome. Nature 422, 93–897.

615   Seshasayee, A.S., Singh, P., & Krishna, S. (2012) Context-dependent conservation of DNA methyltransferases in
616   bacteria. Nucleic Acids Res. 40(15):7066-73.

617    Sim, N.L., Kumar, P., Hu, J., Henikoff, S., Schneider, G., Ng, P.C. (2012). SIFT web server: predicting effects of

618    amino acid substitutions on proteins. Nucleic Acids Res. W452–W457.

619    Stamatakis, A. (2006) RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa

620    and mixed models. Bioinformatics. 22(21):2688-90.

621    Tweedie, S., Charlton, J., Clark, V. & Bird, A. (1997) Methylation of genomes and genes at the invertebrate–

622    vertebrate boundary. Mol. Cell Biol. 17, 1469–1475.

623    Yarza, P., Richter, M., Peplies, J., Euzeby, J., Amann, R., Schleifer, K.H., Ludwig, W., Glöckner, F.O., & Rosselló-

624    Móra, R. (2008) The All-Species Living Tree project: a 16S rRNA-based phylogenetic tree of all sequenced type

625    strains. Syst Appl Microbiol. 31(4):241-50.

626    Zemach, A., McDaniel, I.E., Silva, P., & Zilberman, D. (2010) Genome-wide evolutionary analysis of eukaryotic

627    DNA methylation. Science. 328(5980):916-9.

628    Zemach, A., & Zilberman, D. (2010) Evolution of eukaryotic DNA methylation and the pursuit of safer sex. Curr

629    Biol. 20(17):R780-5.

630

631

632

633

634

635

636

637 **Figure Legends**

638 **Figure 1: DNMTs mapped on protein based tree of life (TOL).** Presence of each of the 3 DNMT enzymes
639 (indicated by 3 different colors) is marked beside each leaf of the tree of life (Ciccarelli F. D. et. al. 2006). The tree
640 and image was obtained from ITOL server.

641 **Figure 2: Phylogenetic trees for all DNMT1, DNMT3A and DNMT3B sequences.** One full length representative
642 sequence from each organism was used to construct the tree for each enzyme. The source organisms/leaves in the
643 tree are color coded according to their kingdom.

644 **Figure 3: Variation in DNMT enzymes in (A) copy number (B) Sequence length and (C) Domain Length.** The
645 distributions are plotted as box whisker plot for each enzyme from each of the seven kingdoms. Boxes represent
646 second and $3^{rd}$ quartile of the distribution while whiskers represent the standard deviation. Horizontal line across the
647 box denotes the median and a small square represent the mean.

648

649 **Figure 4: Sequence identity of DNMT enzyme across groups.** In the left panel all-to-all sequence identity for
650 each sequence in each groups for DNMT1, DNMT3A and DNMT3B are presented as color matrix and in the right
651 panel box distributions are presented for same matrices. Percentage sequence identity is presented as a red/green
652 color matrix. Boxes in the box plot represent second and $3^{rd}$ quartile of the distribution while whiskers represent the
653 standard deviation. Horizontal line across the box denotes the median and a small square represent the mean. In the
654 corresponding table mean and median values are presented along with maximum and minimum values of the
655 distribution.

656 **Figure 5: Domain architecture of DNMT1, DNMT3A and DNMT3B sequences across tree of life.** In the left
657 panel, numbers of sequences with other associated domains are presented. In the right panel the global tree of
658 DNMT1, DNMT3A and DNMT3B are enriched with the combination of domains present in each sequence. The
659 colors on the organisms are according to their kingdom level classification. Different domains are also marked by
660 different colors.

30

661  **Figure 6: Distribution of different types of SNPs in all three DNMT enzymes.**  Both deleterious and tolerated

662  SNPs are mapped and are plotted against with the sequence conservation of those particular alignment columns.

663  Each of the 4 combinations is marked with different color as mentioned in X axis.

664  **Figure 7: SNP conservation and mutation map for DNMT1, DNMT3A and DNMT3B.** (A-C) Count of different

665  types of SNPs based on the combination of SIFT score and Al2Co score. Each combination is indicated by a colour.

666  (D-F) Distribution of SIFT score and Al2CO score. Background color indicates different combination range of both

667  score as indicated in panel A. (G-I) Raw numbers of a particular mutation was marked into a 20 X 20 amino acid

668  matrix to create a mutation map for each of the combination range of SIFT score and Al2CO score. Order of the

669  combination is same with panel A.

670

671  # Supplementary Figure Legends

672  **Figure S1: Identification and classification of DNMT sequences.** The identification and classification protocol

673  along with result is described as a flow chart.

674  **Figure S2: Distribution of enzyme in different organism group.** Number of organisms containing unique

675  combinations of enzymes is represented in a color matrix.

676  **Figure S3: Logo Plots of Motif I to X in DNMT1, DNMT3A and DNMT3B sequences across each kingdom.**

677  **Figure S4: Conservation Score of motifs in DNMT1, DNMT3A and DNMT3B.** Conservation score was

678  calculated for each signature motifs from multiple sequence alignment of each kingdom for three enzymes

679  separately. The conservation score per residue is presented as a bar diagram. Full length of the grey bar indicates

680  100 % conservation.

681  **Supplementary Table S1:** Number of DNMT1 sequences in Phylum/Class/Family/Order

682  **Supplementary Table S2:** Number of DNMT3A sequences in Phylum/Class/Family/Order

683  **Supplementary Table S3:** Number of DNMT3B sequences in Phylum/Class/Family/Order

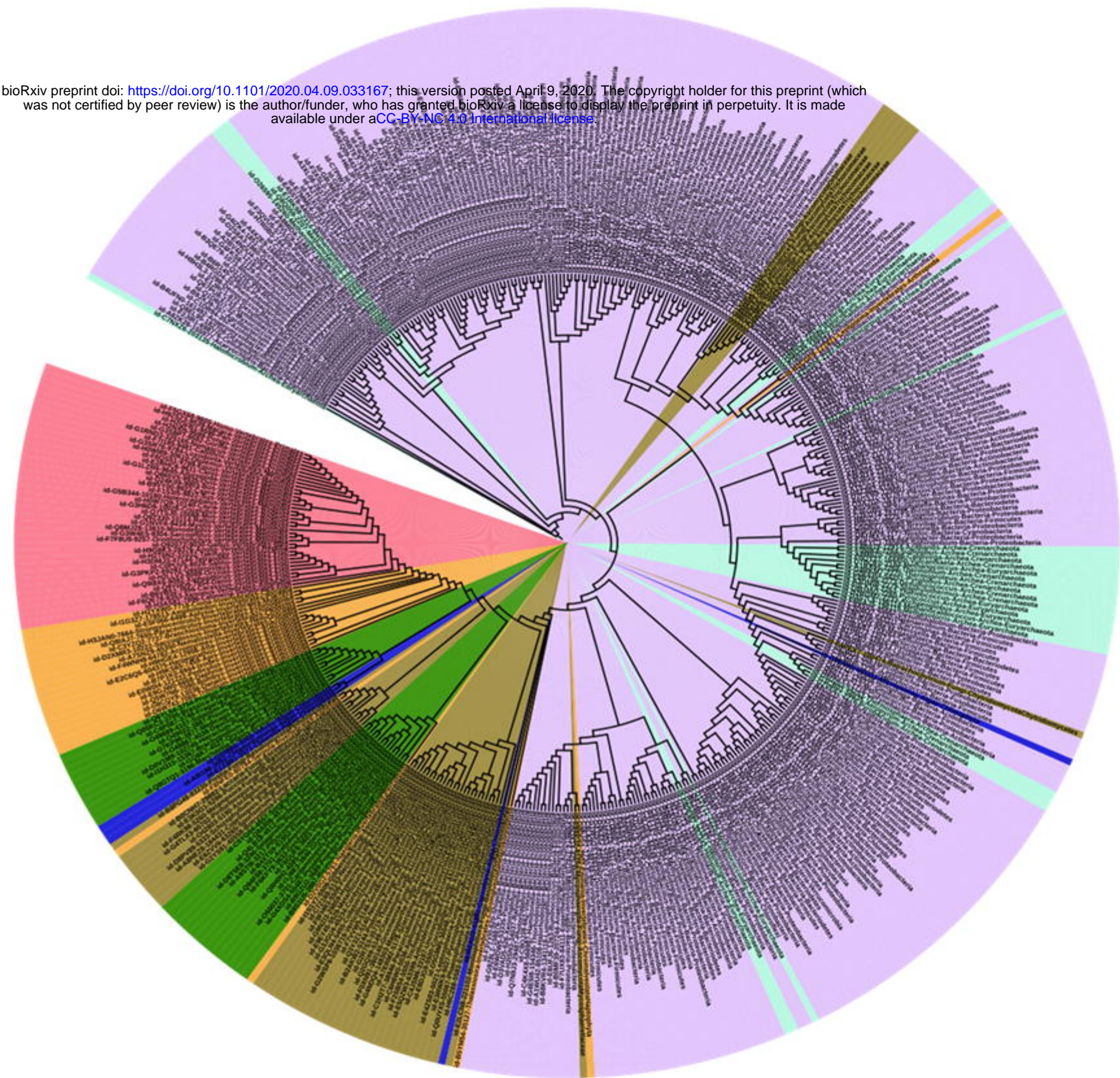684  **Supplementary Table S4:** Details of SNPs in DNMT1, DNMT3A and DNMT3B
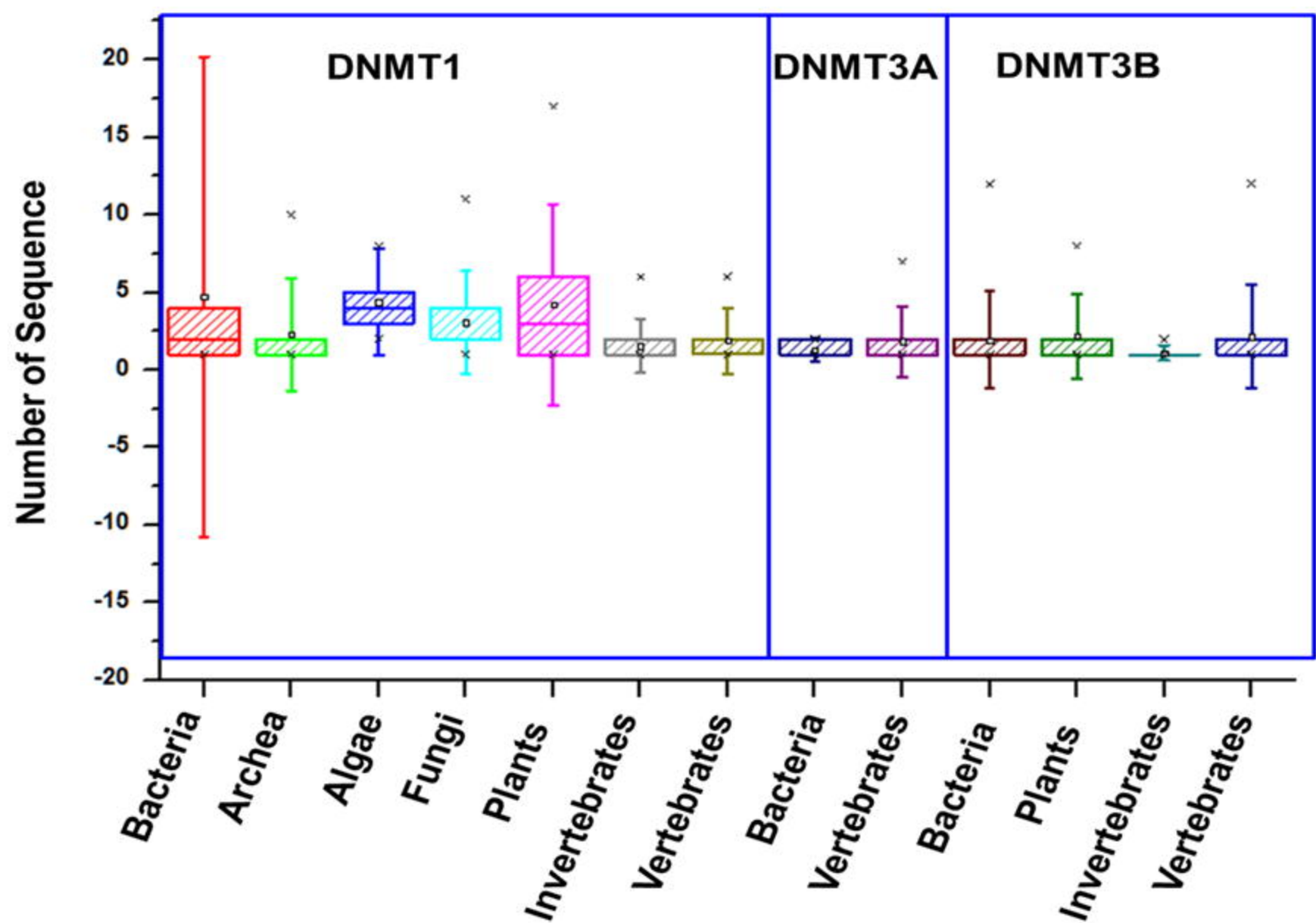
685

# DNMT1

# DNMT3A

# DNMT3B

| | | | | | | |
|---|---|---|---|---|---|---|
| Archea | Bacteria | Algae | Fungi | Plants | Invertebrates | Vertebrates |

# A. *Copy number*



# B. *Sequence Length*

# C. *Domain Length*

| | Bacteria---- | Archea---- | Algae---- | Fungi---- | Plants---- | Invertebrates | -Vertebrates- | ----Bacteria---- | -Vertebrates- | ----Bacteria---- | ----Plants---- | Invertebrates | -Vertebrates- |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Maximum | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Minimum | 10 | 14 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 12 | 10 | 14 | 11 |
| Mean | 37.9 | 51.1 | 32.6 | 39.9 | 51.4 | 46.9 | 80.9 | 65.6 | 65.6 | 41.9 | 79.4 | 62.5 | 81.7 |
| Median | 37.8 | 50.3 | 29.9 | 33.8 | 44.7 | 39.6 | 93.9 | 83.1 | 83.1 | 38.5 | 83.5 | 60.4 | 85.2 |