1    Title: The roles of Conserved Domains in DEMETER-Mediated Active DNA Demethylation *in*

2    *planta*

3

4    Changqing Zhang[1,2,*], Yu-Hung Hung[1,2,*], Xiang-Qian Zhang[1,2,3,*], Dapeng Zhang[4,5], Wenyan

5    Xiao[4], Lakshminarayan M. Iyer[6], L. Aravind[6], Jin Hoe Huh[7] and Tzung-Fu Hsieh[1,2,*]

6    Affiliations:

7    [1]Department of Plant and Microbial Biology, North Carolina State University, Raleigh, NC

8    27695, USA

9    [2]Plants for Human Health Institute, North Carolina State University, North Carolina Research

10   Campus, Kannapolis, NC 28081, USA

11   [3]Guangdong Engineering Research Center of Grassland Science, College of Forestry and

12   Landscape Architecture, South China Agricultural University, Guangzhou 510642, China.

13   [4]Department of Biology, St. Louis University, St. Louis, MO 63103, USA

14   [5]Program of Bioinformatics and Computational Biology, St. Louis University, St. Louis, MO

15   63103, USA

16   [6]National Center for Biotechnology Information, National Library of Medicine, National

17   Institutes of Health, Bethesda, MD 20894, USA

18   [7]Department of Plant Science, Plant Genomics and Breeding Institute, and Research Institute of

19   Agriculture and Life Sciences, Seoul National University, Seoul 08826, Republic of Korea

20   [*]These authors contribute equally to this study.

21    Correspondence and requests for materials should be addressed to T.-F.H. thsieh3@ncsu.edu

1

**Abstract**

DNA methylation plays critical roles in maintaining genome stability, genomic imprinting, transposon silencing, and development. In Arabidopsis genomic imprinting is established in the central cell by DEMETER (DME)-mediated active DNA demethylation, and is essential for seed viability. DME is a large polypeptide with multiple poorly characterized conserved domains. Here we show that the C-terminal enzymatic core of DME is sufficient to complement *dme* associated developmental defects. When targeted by a native DME promoter, nuclear-localized DME C-terminal region rescues *dme* seed abortion and pollen germination defects, and ameliorates CG hypermethylation phenotype in *dme-2* endosperm. Furthermore, targeted expression of the DME N-terminal region in wild-type central cell induces *dme*-like seed abortion phenotype. Our results support a bipartite organization for DME protein, and suggest that the N-terminal region might have regulatory function such as assisting in DNA binding and enhancing the processivity of active DNA demethylation in heterochromatin targets.

37    Double fertilization during sexual reproduction in flowering-plants is a unique process that

38    underlies the distinctive epigenetic reprogramming of plant gene imprinting. In the ovule, a

39    haploid megaspore undergoes three rounds of mitoses to produce a 7-celled, 8 nuclei embryo sac

40    that consists of egg, central, and accessory cells [1]. During fertilization pollen grain elongates and

41    delivers two sperm nuclei to the female gametophyte to fertilize the egg cell and the central cell,

42    respectively. The fertilized egg cell forms the embryo that marks the beginning of the subsequent

43    generation. Fertilization of the central cell initiates the development of endosperm that

44    accumulates starch, lipids, and storage proteins and serves as a nutrient reservoir for the

45    developing embryo [2, 3]. Endosperm is the major tissue where gene imprinting takes place in plant.

46    Genomic imprinting is the differential expression of the two parental alleles of a gene depending

47    on their parent-of-origin, and is an example of inheritance of differential epigenetic states. In

48    Arabidopsis, MET1-mediated DNA methylation and DME demethylation are two modes of

49    epigenetic regulation critical for imprinted expression of many genes [4, 5, 6, 7, 8]. For example,

50    DEMETER (DME) is required for the expression of *MEA*, *FIS2*, and *FWA* in the central cell and

51    in the endosperm while MET1 is responsible for the silencing of *FIS2* and *FWA* paternal alleles [4,

52    7]. Gene imprinting is essential for reproduction in Arabidopsis, and seeds that inherit a maternal

53    *dme* allele abort due to failure to activate *MEA* and *FIS2*, essential components of the endosperm

54    PRC2 complex required for seed viability, in the central cell [4, 9].

55        *DME* encodes a bifunctional 5mC DNA glycosylase/lyase required for active DNA

56    demethylation in the central cell and the establishment of endosperm gene imprinting in

57    Arabidopsis [5]. Additionally, paralogs of DME, REPRESSOR OF SILENCING 1 (ROS1), DML2,

58    and DML3 are required to counteract the spread of DNA methylation mediated by the RNA-

59    directed DNA methylation (RdDM) machinery into nearby coding genes [10, 11]. The three regions

60    in the C-terminal half of DME protein (the A, Glycosylase, and the B regions, or as the AGB

61    region hereafter) are conserved among the DME/ROS1 DNA glycosylase clade, and are required

62    for DME 5mC excision activity *in vitro*. Thus, the AGB region comprise the minimal catalytic

63    core for the enzymatic function, catalyzing direct excision of 5mC from DNA and initiating

64    active DNA demethylation that influences transcription of nearby genes [5, 9, 12].

65         In Arabidopsis, DME-mediated DNA demethylation is preferentially targeted to small,

66    AT-rich, and nucleosome-poor euchromatic transposons that flank coding genes [13]. Consequently,

67    demethylation in the central cell influences expression of adjacent genes only in the maternal

68    genome, and is a primary mechanism of gene imprinting in plant [5, 13, 14, 15]. In addition to small

69    TEs near coding sequences, DME also targets gene-poor heterochromatin regions for

70    demethylation [13]. The mechanism of DME recruitment to its target sites is not known. Studies in

71    ROS1 have uncovered several players required in the ROS1 demethylation pathway [16, 17, 18].

72    Among them *IDM1* encodes a novel histone acetylase that preferentially acetylates H3K18 and

73    H3K23 *in vitro*, and ROS1 target loci are enriched for H3K18 and K23 acetylation *in vivo* in an

74    IDM1-dependent manner [19]. Thus, IDM1 marks ROS1 target sites by acetylating histone H3 to

75    create a permissible chromatin environment for ROS1 function. The Arabidopsis SSRP1

76    (STRUCTURE SPECIFIC RECOGNITION PROTEIN1), a component of the FACT (facilitates

77    chromatin transcription/transaction) histone chaperone complex, has been shown to regulate

78    DNA demethylation and gene imprinting in Arabidopsis [20]. Linker histone H1 functions in

79    chromatin folding and gene regulation [21, 22, 23, 24], and was shown to interact with DME in a yeast

80    two-hybrid screen and in an *in vitro* pull-down assay [25]. Loss-of-function mutations in *H1* genes

81    affect the imprinted expression of *MEA* and *FWA* in Arabidopsis endosperm, and impair

82    demethylation of their maternal alleles, suggesting that H1 might participate in the DME

83    demethylation process by interaction with DME [25].

84         Computational analysis showed that the DME/ROS1 like DNA glycosylases contain a

85    core with multiple conserved globular domains, and except for the well-characterized

86    glycosylase domain, very little is known about the function of the other domains. Here we show

87    that the C-terminal region of DME necessary for 5-methylcytosine excision activity *in vitro* is

88    sufficient to complement *dme* seed abortion and pollen germination defect, and partially rescue

89    DNA hypermethylation phenotype in endosperm. We present evidence that the region N-

90    terminal to the glycosylase domain can affect endogenous DME activity in a dominant negative

91    manner when ectopically expressed in the nuclei of wild-type central cells. We propose a

92    bipartite structural and functional organization model for the DME/ROS1 family of DNA

93    glycosylases consisting the modular C-terminal AGB region that can substitute for DME's

94    developmental function and the NTD region that might have regulatory functions such as

95    assisting DNA binding and enhancing the processivity of demethylation in heavily methylated

96    genomic regions.

97

98    **Results**

99    **The DME catalytic core region is sufficient to complement *dme* associated developmental**

100    **defects.** Previous studies have revealed that the C-terminal half of DME comprising the three

101    conserved <u>A</u>, <u>G</u>lycosylase, and <u>B</u> regions (the <u>AGB</u> region, as shown in Supplementary Fig. 1a)

102    are required for *in vitro* 5mC excision activity [5], and deletion of the non-conserved linker

103    between domain A and the glycosylase domain (interdomain 1; ID1) does not affect DME *in*

104    *vitro* enzymatic activity [26, 27]. Thus, the AGB region is thought to be the core catalytic region for

105    DME *in vitro* enzymatic activity. However, it is unknown whether the AGB region alone is

5

106    sufficient for DME function *in vivo*. To determine if the AGB region is functional *in vivo*, we

107    tested if expressing the AGB region in the central cell can complement *dme* seed abortion

108    phenotype. A transgene carrying a 3.1-kb *DME* cDNA that encodes the C-terminal half of DME

109    (DME$^{CTD}$, residue 936-1987) under the control of a native DME promoter was introduced into

110    *DME/dme-2* heterozygous plants by using the floral dipping method [28]. Since DME$^{CTD}$ lacks a

111    nuclear localization signal (data not shown), a classical SV40 nuclear localization signal

112    (PKKKPKV) was introduced in front of the C-terminal fragment (designated as *nDME$^{CTD}$*, see

113    Supplementary Fig. 1b) to ensure proper nuclear localization. We obtained multiple independent

114    transgenic lines and assessed the transgene's ability to complement *dme-2* seed abortion

115    phenotype.

116        The self-pollinated *DME/dme-2* plants produce 50% of normal seed that inherited wild type

117    DME maternal allele, and the other 50% of aborted seed that inherited mutant *dme-2* maternal

118    allele. In self-pollinated transgenic plants that carry a single locus of *nDME$^{CTD}$* or *DME$^{FL}$* (full

119    length DME.2 cDNA, major isoform of DME [29]) transgenes, we observed about 25% aborted

120    seeds among independent transgenic lines, indicating that *nDME$^{CTD}$* and *DME$^{FL}$* complement

121    *dme* seed abortion phenotype (Fig. 1a, b, Supplementary Table 1). In addition, we also

122    transformed *nDME$^{CTD}$* and *DME$^{FL}$* into *dme-2/dme-2* homozygous plants (see Materials and

123    Methods for isolation and characterization of *dme-2/dme-2* homozygous lines in Col-*gl*), both

124    constructs produced T1 transgenic plants that displayed expected 50% seed abortion rate (Fig. 1b,

125    Supplementary Table 1). Seed abortion caused by *dme* mutations is in part due to defects in

126    activating imprinted PRC2 subunit genes required for endosperm development [5, 9, 30, 31, 32]. We

127    use qRT-PCR to check if nDME$^{CTD}$ also restores DME target genes expression in the central cell.

128    Indeed, *FIS2* and *FWA* expression is restored in the complemented lines (Fig. 1c). Thus

129    nDME$^{CTD}$ can substitute for the endogenous DME activity for seed viability, and active DME

130    target genes expression.

131         In addition to maternal effects on seed viability [9], mutations in DME also affect pollen

132    function in Col-0. When *DME/dme-2* heterozygous plants are self-pollinated, only about 20-30%

133    of the viable F1 progeny are heterozygous (Supplementary Table 2), due to decreased *dme* pollen

134    germination rate [33]. To test whether nDME$^{CTD}$ can rescue *dme* pollen phenotype, we pollinated

135    wild type Col-0 with pollen derived from transgenic lines that are homozygous for the *dme-2*

136    allele and carry a single locus of the *nDME$^{CTD}$* transgene (*dme-2/dme-2; nDME$^{CTD}$/~*). If

137    nDME$^{CTD}$ does not complement *dme-2* pollen germination defects, we expect roughly half of the

138    F1 progeny will carry the nDME$^{CTD}$ transgene (hygromycin resistant) because mutant pollen

139    with or without the transgene would germinate with equal frequency. Instead, we observed 65% -

140    90% of the F1 progeny are hygromycin resistant (Table 1), indicating that nDME$^{CTD}$

141    complements *dme-2* pollen germination defect. These results show that expressing the C-

142    terminal half of DME protein in the nucleus is sufficient to rescue *dme* visible phenotypes *in*

143    *planta*.

144

145    **nDME$^{CTD}$ partially rescue *dme-2* CG hypermethylation phenotype in the endosperm.**  In

146    Arabidopsis seed viability depends on the DME activity in the central cell to activate the

147    MEDEA/FIS2/MSI1/FIE PRC2 complex required for endosperm development. In addition,

148    DME is required to demethylate multiple maternally (*MEGs*) or paternally expressed imprinted

149    genes (*PEGs*) to establish their parent-of-origin specific expression patterns in the endosperm [13,

150    15]. Thus, in *dme* mutant endosperm, discrete genomic loci targeted by DME for demethylation

151    are hypermethylated [13]. Since nDME$^{CTD}$ complements *dme* seed abortion, and activates DME

152   target gene expression (Fig. 1), we assumed it does so by demethylating the central cell genome

153   and activating PRC2 genes essential for seed development.  To test this hypothesis, and to

154   examine the extend of nDME$^{CTD}$ demethylation activity *in vivo*, we manually isolated *nDME$^{CTD}$*-

155   complemented   endosperm   (*dme-2/dme-2;nDME$^{CTD}$/nDME$^{CTD}$*),   determined   the   DNA

156   methylation profile by whole genome bisulfite sequencing, and compared the complemented

157   methylomes to those of wild-type and *dme-2* endosperm. Methylomes from three independent

158   lines were generated and compared with that of *dme-2* endosperm. We observed although the

159   differentially methylated regions (DMRs) between each independent lines do not completely

160   overlap, the DMRs unique to each line are also demethylated in other lines (Supplementary Fig.

161   2, 3), suggesting that the number of overlapped DMRs was underestimated due to the cutoff used

162   in defining the DMRs, similar to what's observed in a recent study [34]. We therefore used the

163   combined reads from three independent lines for the subsequent analyses so that all comparisons

164   are confined to the same cutoff criteria (see Materials and Methods). As expected, several DME

165   regulated *MEGs* and *PEGs* are demethylated compared to *dme-2* endosperm, indicating that

166   nDME$^{CTD}$ is correctly recruited to these loci for demethylation (Fig. 2a). We focused our analysis

167   on   previously   determined   differentially   methylated   sites   between   *dme-2*   and   wild-type

168   endosperm (*dme* hyper-DMRs, the DME canonical targets) [13, 15]. Overall, the CG methylation

169   levels in these canonical DME target sites are reduced in the complemented endosperm,

170   indicating that nDME$^{CTD}$ is directed to these endogenous DME target sites for demethylation.

171   However, compared to wt endosperm, these *dme* hyper-DMRs are demethylated to a lesser

172   degree by the nDME$^{CTD}$ (Fig. 2b). Thus nDME$^{CTD}$ only partially rescues the *dme* CG

173   hypermethylation phenotype in the endosperm. The DMRs of *dme* relative to wild-type

174   endosperm or to *nDME$^{CTD}$*-complemented endosperm partially overlap (Supplemental Fig. 4).

175    However, among the DMRs unique to nDME$^{\text{CTD}}$, we also observed decreased CG methylation in

176    WT endosperm compared to *dme*, indicating that they are also demethylated by the endogenous

177    DME. Similarly, among the DMRs unique to wt endosperm, these regions are also demethylated

178    by the nDME$^{\text{CTD}}$. Thus nDME$^{\text{CTD}}$ appears to partially demethylate the majority of the loci

179    targeted by the endogenous DME. These observations also suggest that intact full-length DME

180    protein is required for robust and complete demethylation *in vivo*.

181        We next examined the methylome of *dme-2* endosperm complemented by the full length

182    *DME.2 cDNA* (designated as *DME$^{FL}$*). Unexpectedly, based on the number of DMRs between

183    *dme* and *DME$^{FL}$*-complemented endosperm and the level of CG methylation within the DMRs

184    (Fig. 2c), DME$^{\text{FL}}$ appears to be less active compared to endogenous DME, or to nDME$^{\text{CTD}}$, albeit

185    it being able to complement *dme* seed abortion (Fig. 1b) [9, 35]. Since the *DME$^{FL}$* transgene only

186    differs from *nDME$^{CTD}$* by the N-terminal region, reduced activity of DME$^{\text{FL}}$ compared to

187    DME$^{\text{CTD}}$ cannot be attributed to the lack of introns or 3' flanking sequences that might be needed

188    for robust DME protein production. Indeed, we found both transgenes are expressed at

189    comparable levels in *DME$^{FL}$*- and *nDME$^{CTD}$*-complemented lines used in the methylome study

190    (Supplemental Fig. 5), indicating lower activity of DME$^{\text{FL}}$ compared to nDME$^{\text{CTD}}$ is not due to

191    their differential transcript abundance. Nevertheless, comparison of CG methylation levels in

192    DMR regions unique to DME$^{\text{FL}}$, nDME$^{\text{CTD}}$, or endogenous DME also reveals that unique DMR

193    regions are more or less hypomethylated in WT or in complemented endosperm relative to *dme*

194    endosperm. Thus the methylome difference between wt, *DME$^{FL}$*-, and *nDME$^{CTD}$*-complemented

195    endosperm appears to be more in the degree of demethylation, rather than in targeting specificity.

196

197 **Function of the N-terminal region in DME-mediated active DNA demethylation.** The *dme-2*

198 allele is caused by an activation-tagging T-DNA insertion in the middle of the A region

199 (Supplementary Fig. 1a)[9]. We found that in floral buds of *dme-2/dme-2* plants, the endogenous

200 *DME* transcripts downstream of T-DNA insertion site is greatly reduced compared to wild-type

201 Col-0 plants, but the level of DME transcripts upstream of the T-DNA insertion site is relatively

202 high (Supplementary Fig. 6). We suspected these transcripts could produce truncated form of

203 DME proteins that might interfere with the DME$^{FL}$ transgene activity. To test this hypothesis, we

204 transformed wild-type Col-0 plants with an engineered GFP-tagged DME NTD (using the

205 genomic DNA sequence upstream of T-DNA insertion site, encoding residues 1-1022, designated

206 as *DME$^{NTD}$-GFP*) transgene mimicking the *dme-2* T-DNA insertion (Supplementary Fig. 1B).

207 Clear GFP signals are observed in the central cell nuclei of transgenic lines (data not shown). We

208 also observed about one third of transgenic lines showing apparent *dme-2* like seed abortion

209 phenotype, with abortion rates ranging from 10% to ~ 40% (Supplementary Table 3, 4) in the T1

210 plants, suggesting that expression of DME$^{NTD}$ has a dominant negative effect on endogenous

211 DME protein.

212     To minimize the possibility and the degree of transgene induced sense co-suppression, we

213 reverse translated DME$^{NTD}$ protein sequence into cDNA sequence using the human codon usage

214 table. As a result, the re-engineered "humanized" version of NTD (mDME$^{NTD}$) codes for the

215 identical protein sequence but with no significant nucleotide sequence similarity to the original

216 cDNA sequence to induce co-suppression (Supplementary Table 5). In addition, a GFP tag was

217 added to the C-terminus (mDME$^{NTD}$-GFP) to monitor its expression (Fig. 3a). We generated 28

218 independent transgenic lines, and among them 16 lines showed seed abortion rate of 5% - 52%

219 (Supplementary Table 3, 6). The aborted seeds resemble *dme* mutant seeds with abnormal

10

220   endosperm, arrested embryo, and shriveled brown seeds (Fig. 3b, c). We selected four lines with

221   high, medium, or no seed abortion rate (Fig. 3d), and assessed the endogenous DME transcript

222   abundance. As shown in Fig. 3e, among lines with different seed abortion rate, the endogenous

223   DME mRNA abundance is similar to that of the vector control line, indicating the severity of

224   seed abortion phenotype is not due to interference of endogenous *DME* transcripts. Furthermore,

225   the rate of seed abortion is positively correlated with the levels of *mDME^{NTD}-GFP* mRNA (Fig.

226   3f), suggesting the degree of seed abortion is likely due to the levels of transgene expression. We

227   next tested whether expression of *nDME^{CTD}* or *DME^{FL}* in WT Col-0 can also induce seed

228   abortion phenotype. For each construct, more than 25 independent transgenic lines were

229   examined and none resulted in any seed abortion phenotype (Supplementary Table 3). Thus the

230   dominant negative effect appears to be specific to the DME NTD region.

231

232   **Evolutionary history and late acquisition of the N-terminal region of DME-like proteins.**

233   We show the C-terminal half of DME is sufficient to complement *dme* mutant developmental

234   phenotypes, and can be recruited to most of the DME target loci. Thus the DME^{CTD} most likely

235   contains intrinsic targeting information. To gain insights from the evolution of the conserved

236   domains in DME, we conducted sequence searches of the NR database with various homologs as

237   query. The core of the DME-like proteins, as previously reported [36], comprises the catalytic

238   glycosylase domain of the HhH (helix-hairpin-helix) modules followed by the FCL ([Fe4S4]

239   cluster loop) motif and a divergent version of an RRM (RNA Recognition Motif) fold domain

240   (Fig. 4). The DNA glycoslase and FCL domains span the A and G regions, whereas the RRM

241   fold domain corresponds to the B region of angiosperm DME homologs.  A diversity of domains

242   associate with the basic DME core can be found across various clades. Land plants and

11

243    charophytes (Streptophyta) possess a permuted divergent version of the umethylated CpG

244    recognizing CXXC domain (containing only one of two structural repeats of the classical CXXC

245    domain) between the FCL and RRM domains. By contrast, one or more copies of the CXXC

246    domain can be found in chlorophyte and stramenopile algae at distinct positions.  Some algal

247    DME homologs (from Chlorophyte and stramenopile) also possess other chromatin-modification

248    reader (Tudor and PHD domains), DNA binding (AT-hook motif), and the DnaJ domain which

249    interacts with the chaperone Hsp70 [36, 37]. These accessory domains suggest a potential role for

250    regulating the associated DNA glycosylase activity according to the DNA methylation (via

251    CXXC) or chromatin status (via PHD, Tudor) of the cell in which they are expressed.

252        The N-terminal half of the DME consists of a large portion of unstructured, low complexity

253    sequences (residues 346-947), a stretch of basic amino acid-rich direct repeats (residues 291-

254    345), and a 120 amino-acid N-terminal domain (DemeN) of unknown function (residues 1-

255    120)(see Supplementary Fig. 7 for sequence alignment). The DemeN domain and charged

256    repeats are restricted to the angiosperm lineage and appears to be a late acquisition during land

257    plant evolution.

258        In summary, the evolutionary history of the DME domains can be summarized as follows:

259    bacterial versions of the HhH-FCL pair from a cyanobacterial source fused to an RRM-fold

260    domain and further acquired an insert in the glycosylase domain to give the ancestral form in the

261    plant lineage. This was likely then transferred to the stramenopiles from a secondary chlorophyte

262    endosymbiont of this lineage. Finally, at the base of the streptophyte radiation, DME acquired a

263    permuted CXXC, and later the DemeN domain and the associated charged repeats were acquired

264    in the angiosperm lineage, possibly to facilitate and ensure a robust and thorough demethylation.

265

**Discussion**

266

267    We show for the first time that the core conserved region of the DME protein containing the

268    DNA-glycosylase, FCL, divergent and permuted CXXC and divergent RRM domains is

269    sufficient to rescue visible phenotypic defects caused by *dme* mutation. Although this truncated

270    form of DME protein demethylates the majority of the canonical DME target sites, it does so in a

271    less active and less efficient manner compared to the endogenous protein. We see two

272    possibilities that might explain this lower activity and efficiency: 1) Critical cis-elements

273    residing within introns or in 3'-end flanking sequences that are missing in the transgene might be

274    required for robust transgene expression. 2) The N-terminal region might be required for full

275    DME activity *in vivo*. Unfortunately, our attempt to assess the difference between DME$^{FL}$ and

276    nDME$^{CTD}$ was confounded by the possible interference from truncated NTD proteins due to T-

277    DNA insertion in *dme-2* background. We suspect this might contribute to the reduced DMRs

278    observed in *DME$^{FL}$*complemented endosperm. Therefore, we believe it is premature to draw any

279    conclusion based on direct comparisons between *DME$^{FL}$*- and *nDME$^{CTD}$*-complemented

280    endosperm methylomes (Fig. 2c).

281    Since the C-terminal AGB region is sufficient for DME's seed viability function in *planta*,

282    and can be recruited to most of the canonical DME target sites, the CTD polypeptide most likely

283    contains sufficient targeting information. *in vitro* studies of ROS1 suggest that the B region

284    containing the CXXC and RRM domains is essential for the glycosylase and lyase activities, and

285    might recognize modified DNA[38]. It is possible that the permuted CXXC domain is required to

286    direct the protein to the target sites, or is involved in discriminating methylated vs un-methylated

287    cytosines[39]. This is supported by mutation studies that implicate a potential role for this domain

288    in DME *in vivo* function, but not in vitro enzymatic activity (Huh and Hsieh, unpublished

13

289    results). Similarly, the role of the enigmatic divergent RNA-recognition motif (RRM) domain is

290    also not fully understood. Mutagenesis screens for residues required for demethylation activity in

291    bacteria identified multiple amino acid residues within the RRM domain[40]. Although the

292    involvement of RNA species in the active DNA demethylation process has not been firmly

293    established, an RRM protein ROS3 required for ROS1 demethylation suggests a potential role of

294    non-coding RNAs in the active DNA demethylation pathway in Arabidopsis[41]. While it is

295    tempting to speculate a role for RNA-binding, the DME RRM might also bind single-stranded

296    DNA with methylated bases.

297        Based on the reduced demethylation activity of nDME$^{CTD}$ on the canonical DME target sites,

298    we suspect the NTD region might be required for full and robust demethylation activity probably

299    to ensure that the imprinting network is properly activated and maintained (e.g., by subsequent

300    PRC2 activity). To achieve this, the DME NTD might function to assist the glycosylase enzyme

301    by tightly binding to DNA template for more complete and thorough demethylation. Supporting

302    such model, *in vitro* study of ROS1 activity on 5mC excision revealed that the basic repeats

303    (3DR, AT-hooks) region binds strongly to DNA template non-specifically, and removal of NTD

304    region impairs the sliding capacity of the protein on DNA template[42], and significantly reduced

305    ROS1 5mC excision activity[43]. We observed reduced degree of demethylation by nDME$^{CTD}$

306    regardless of target length (Supplemental Fig. 8), suggesting that NTD is needed for complete

307    demethylation in all the target sites.

308        Although DME preferentially targets smaller euchromatic transposons that flank coding

309    genes, it also targets gene-poor heterochromatin regions for demethylation [13]. The biological

310    significance of heterochromatin demethylation by DME is not known, but was speculated to

311    involve reinforcing DNA methylation in egg cell and subsequently in the embryo [13]. These

312     heterochromatin target sites are densely methylated, and demethylation by DME results in longer

313     DMRs between *dme-2* and wt endosperm. Interestingly, the number of longer DMRs is

314     significantly reduced between *dme-2* and $nDME^{CTD}$-complemented endosperm, suggesting that

315     removal of NTD region also reduces the processivity of demethylation in long target sites

316     (Supplemental Fig. 9a). Since heterochromatin regions are compacted, demethylation in such

317     loci will require substantial chromatin remodeling such as eviction of nucleosomes for DME to

318     gain access to the templates. It is tempting to speculate that the conserved motif in the DemeN

319     domain might recruit other factor(s) via protein interaction to remodel local chromatins to permit

320     DME demethylation. However, based on current data we cannot unequivocally ascribe NTD's

321     function due to lack of proper full length DME transgenic comparison. Nevertheless, our results

322     caution that peculiarity in certain genetic backgrounds (e.g., *dme-2*) might confound data

323     interpretation. Future work on DME functional study could benefit from the generation of a

324     clean loss-of-function background such as deleting the entire DME locus using CRISPR-assisted

325     genome editing techniques.

326     We envision a possible model where the AGB region is sufficient for directing DME to target

327     loci while NTD region is required for interacting with local chromatin environment, stabilizing

328     binding to chromosomal templates, and assisting demethylating flanking sequences. In the

329     absence of NTD, $nDME^{CTD}$ can still demethylate majority of target sites, but in a less-efficient

330     manner, likely due to the lack of non-specific DNA-binding by the basic AT-hook motifs. We

331     surveyed wt DMRs that are longer than 1.5 kb, and found that these regions are also $nDME^{CTD}$'s

332     DMRs, but are shorter in length (Supplemental Fig. 9b), possibly due to missing the DemeN

333     domain. If NTD is needed for longer and more robust demethylation, why ectopic expression of

334     NTD causes dominant negative (DN) effects on endogenous protein? Classical examples of

335    dominant negative mutation often involve protein-protein interactions that are disrupted by

336    mutated or truncated form of one particular partner or subunit. Although we do not have any

337    evidence to suggest DME might homodimerize to become active, any weak physical interaction

338    caused by ectopic NTD expression might induce conformational change that renders DME non-

339    functional. Unfortunately our attempt to assess whether the NTD of DME can interact with each

340    other was confounded by the self-activating activity of DME.2 NTD in yeast two-hybrid assay

341    when fused to the GAL4 DNA binding domain (data not shown). Their possible interaction will

342    need to be assessed by alternative strategies. Another possibility is that NTD binds and titrates

343    out an important interacting partner required to activate DME through conformational change

344    (allosteric interaction). By removing NTD, the AGB region is liberated from such

345    conformational constrain and can demethylate its target sites. It is also possible that the non-

346    specific DNA binding activity of NTD competes with DME for target sites, thereby reducing the

347    overall efficiency of DME.  The molecular underpinning of how NTD induces DN effect remains

348    to be elucidated. From an evolutionary viewpoint, the use of an active DME-based

349    demethylation appears to have been acquired early in the plant lineage. The presence of several

350    accessory domains in addition to the conserved core suggests adjustments to the chromatin and

351    methylation environment of the different species. The presence of additional domains such as the

352    DemeN and basic repeats in angiosperms and the permuted CXXC domain in streptophyta

353    lineage might reflect the adjustment to the unique methylation and chromatin environment of the

354    larger Streptophyta and land plant genomes.

355

356    **Materials and Methods**

357    **Molecular Cloning of Constructs Used in this Study.**

16

358     All general molecular manipulations followed standard procedures (Sambrook et al. 1989). Q5

359     High Fidelity DNA polymerase (NEB, Ipswich MA, USA) was used for PCR amplifications.

360     PCR products were purified using AMPure XP beads (Beckman Coulter, Indianapolis IN, USA).

361     The sequences of all plasmid constructs were confirmed by sequencing (Eton, Research Triangle

362     Park NC, USA). All PCR primers and double-stranded DNA fragments were synthesized by

363     Integrated DNA Technologies (Coralville IA, USA), and sequences are listed in Supplementary

364     File 1.

365     A binary plasmid vector, pFGAMh, was modified to facilitate the generation of plasmid

366     constructs using the Gibson assembly method. In brief, the replication origins and T-DNA

367     borders originated from pFGC5941 (GenBank Accession: AY310901). A hygromycin resistance

368     gene (HPTII) under the control of the mannopine synthase promoter was installed for selection

369     of transgenic seedlings. A Gateway attR cassette (rfa, Invitrogen, Carlsbad CA, USA), flanked

370     with unique restriction sites XhoI and XbaI-SpeI was placed upstream octopine synthase

371     polyadenylation signal (OCS3'). Plasmid pFGAMh, digested with restriction enzymes XhoI and

372     XbaI, was used to generate plasmids pDME:DME$^{CTD}$, pDME:nDME$^{CTD}$ and

373     pDME:GFP::DME$^{CTD}$ using the Gibson assembly method. The DME.2 upstream regulatory

374     sequence (DMEpro; 2895 bp upstream of DME.2 translation start codon ATG) was PCR-

375     amplified using primer pair VeDME/P3R and Col-0 gDNA as template. The coding sequence of

376     linker-AGB (with a 6-Ala linker to its N-terminus; 3174 bp), was PCR-amplified using primer

377     pair lnAGBF/CTDVeR and Col-0 cDNA as template. To bridge these two fragments (DMEpro

378     and linker-AGB), one of the following three DNA fragments was used in the assembly reactions.

379     For pDME:DME$^{CTD}$, a 50-bp fragment was generated by annealing DNA oligos ATGF and

380     ATGR. For pDME:SV40NLS::AGB, a 71-bp fragment was generated by annealing DNA oligos

17

381    S40F and S40R followed by two rounds of PCR reactions. For pDME:GFP::DME$^{CTD}$, a 761-bp

382    fragment was PCR-amplified using primer pair p3GFPF/dmGFPR and plasmid DNA pGFP-JS

383    (Jen Sheen,Massachusetts GeneralHospital, Boston MA, USA) as template.

384    An intermediate plasmid vector, DME-P3-attR-AGB, was generated by digesting plasmid

385    pDME:SV40NLS::AGB with restriction enzymes AflII and NcoI, and re-assembled with a 2800-

386    bp fragment, which was produced through overlap PCR with 3 primer pairs, upAflII/P3attR,

387    P3attF/attAGBR and attAGBF/dnNcoI, and Col-0 gDNA, attR cassette and Col-0 cDNA as

388    templates. The resulting plasmid DME-P3-attR-AGB bears (1) the same 2895-bp regulatory

389    sequence as the above constructs, (2) an attR cassette flanked by unique restrict sites XbaI and

390    BglII, and (3) AGB coding sequence (3156 bp). To generate pDME:DME$^{FL}$, plasmid DME-P3-

391    attR-AGB was digested with XbaI and BglII, and assembled with a 2985-bp sequence, which

392    was generated through overlap PCR using primer pairs S1-5e/IN3R and IN3F/S1-5R, and Col-0

393    gDNA as template. The resulting plasmid pDME:DME$^{FL}$ carries the complete DME.2 coding

394    sequence and intron 2 sequence (6075 bp) immediately downstream of the 2895-bp regulatory

395    sequence with no additional sequences.

396    The intermediate plasmid vector DME-P3-attR-AGB was digested with restriction enzymes

397    BglII and SpeI (to completely remove the AGB coding sequence), and re-assembled with a 786-

398    bp sequence, which included the coding sequence of GFP (with its start codon ATG changed to

399    TTG) and was PCR-amplified using primers ttGFPF and SpeGFPR and plasmid DNA pGFP-JS

400    as template. The resulting plasmid DME-P3-attR-GFP was used as an intermediate plasmid

401    vector to generate constructs pDME:DME$^{NTD}$::GFP and pDME:mDME$^{NTD}$::GFP. Plasmid DME-

402    P3-attR-GFP was digested with XbaI and BglII, and assembled with two DNA fragments: a

403    3289-bp sequence was PCR-amplified using primers S1-5F and dme2tR2 and Col-0 gDNA as

404   template and a 158-bp synthetic DNA fragment (FragQ20) (Integrated DNA Technologies,

405   Coralville IA, USA). The resulting construct pDME:DME$^{NTD}$::GFP included the 2895-bp

406   upstream regulatory sequence, the 3332-bp sequence downstream of translation start codon ATG,

407   the coding sequence of 6-Ala linker, and the coding sequence of GFP. Note the NTD coding

408   sequence included the first 86 bp of intron 4 of gene DME.2, and it was designed to mimic dme-

409   2 T-DNA insertion. To generate pDME:mDME$^{NTD}$::GFP, the sequence of the first 1012 amnio

410   acid residues of DME.2 protein was converted to DNA sequence using program EMBOSS

411   Backtranseq (http://www.ebi.ac.uk/Tools/st/emboss_backtranseq/) and the Homo sapiens codon

412   usage table. The sequence was then analyzed using online programs SoftBerry FSPLICE

413   (http://linux1.softberry.com/berry.phtml?topic=fsplice&group=programs&subgroup=gfind) and

414   NetPlantGene2 (http://www.cbs.dtu.dk/services/NetPGene/), and manually edited to disrupt

415   potential splicing donor sites or acceptor sites. The mDME$^{NTD}$ sequence (3036 bp) and upstream-

416   and downstream-overlapping sequence are broken into 4 fragments, and synthesized by

417   Integrated DNA Technologies (Coralville IA, USA). The 4 DNA fragments were assembled with

418   plasmid DME-P3-attR-GFP digested with XbaI and BglII, resulting construct

419   pDME:mNTDh::GFP.

420

421   **Whole-Genome Bisulfite Sequencing and DNA Methylome Analysis**

422   Genomic DNA were isolated from hand dissected, 7-9 DAP *dme-2* endosperm that has been

423   complemented by *DME$^{FL}$* or *nDME$^{CTD}$* (*dme-2/dme-2;DME$^{FL}$/DME$^{FL}$* or *dme-2/dme-2;*

424   *nDME$^{CTD}$ nDME$^{CTD}$*). Whole genome bisulfite sequencing library was constructed as described

425   before [13, 44]. Approximately 20-50 ng of purified genomic DNA was spiked with 0.5ng of

426   unmethylated cl857 *Sam7* Lambda DNA (Promega, Madison, WI) and sheared to about 300bp

19

427     using Covaris M220 (Covaris Inc., Woburn, Massachusetts) under the following settings: target

428     BP, 300; peak incident power, 75 W; duty factor, 10%; cycles per burst, 200; treatment time, 90

429     second; sample volume 50µl. The sheared DNA was cleaned up and recovered by 1.2x AMPure

430     XP beads then followed by end repaired and A-tailing (NEBNext Ultra II DNA Library Prep Kit

431     for Illumina, NEB) before ligation to NEBNext methylated multiplex adapters (NEBNext

432     Multiplex Oligos for Illumina, NEB) according to the manufacturer's instructions. Adaptor-

433     ligated DNA was cleaned up with 1x AMPure XP beads. The purified adaptor-ligated DNA was

434     spiked with 50ng of unmethylated cl857 *Sam7* Lambda DNA and subjected to one round of

435     sodium bisulfite conversion using the EZ DNA Methylation-Lightning Kit (Zymo Research

436     Corporation, Irvine, CA) as outlined in the manufacturer's instructions with 80 min of incubation

437     time. Half of the bisulfite-converted DNA molecules was PCR amplified with the following

438     condition: 2.5 U of ExTaq DNA polymerase (Takara), 5 ul of 10 x Extaq reaction buffer, 25 µM

439     dNTPs, 1 ul of index Primers (10 uM) in 50 uL reaction. The thermocyling condition was as

440     follows: 95 °C for 2 min and then 10 cycles each of 95 °C for 30 s, 65 °C for 30 s, and 72 °C for

441     60 s. The enriched libraries were purified twice with 0.8x (v/v) AMPure XP beads to remove any

442     adapter dimers. High throughput sequencing was performed by Novogene Corporation (USA).

443     For each genotype, sequencing reads from three individual transgenic lines were combined.

444     Sequenced reads were mapped to the TAIR10 reference genomes and DNA methylation analyses

445     were performed as previously described (Supplementary Table 7) [13]. Fractional CG methylation

446     in 50-bp windows across the genome was compared between *dme*, wild-type (GSE38935 [13]),

447     $DME^{FL}$- $nDME^{CTD}$- complemented *dme-2* endosperm. Windows with a fractional CG

448     methylation difference of at least 0.3 in the endosperm comparison (Fisher's exact test p-value <

449     0.001) were merged to generate larger differentially methylated regions (DMRs) if they occurred

450   within 300 bp. DMRs were retained for further analysis if the fractional CG methylation across

451   the whole DMR was 0.3 greater in *dme* endosperm than in wild-type endosperm (Fisher's exact

452   test p-value < $10^{-10}$), and if the DMR is at least 100-bp long. The merged DMR lists are in the

453   Supplemental File 2. The *dme* and wild-type endosperm data used in this study were derived

454   from crossed between *Col* (female parent) and *Ler* (male parent) (GSE38935, [13]). To avoid

455   potential ecotype-specific methylation difference, *Ler* hyper-DMRs relative to Col-0 endosperms

456   (GSE52814, [45]) were identified using the same criteria as described above and excluded from

457   further analyses. For making the Venn diagram, merged DMR regions were converted into 50-bp

458   windows. Only windows with methylation scores in all samples were retained for comparison in

459   Venn diagram and boxplot analysis.

460

461   **Plant Materials and Complementation Assays**

462   We found we can easily obtained *dme-2/dme-2* Col-*gl* plants from *DME/dme-2* heterozygotes if

463   we rescued seeds prior to desiccation on MS sucrose plates. This is consistent with the report that

464   *fis* endosperm cellularization defect and embryo arrest can be rescued by culturing the

465   developing seeds in sucrose media because *fis* seeds have reduced hexose level [46]. Using this

466   method we generated multiple homozygous lines, and we did not detect any difference between

467   individuals in terms of normal seed rate or visible phenotype. The adult *dme-2/dme-2* plants are

468   morphologically indistinguishable from wild-type Col-*gl* plants but produce ~0.1% viable

469   mature seeds. These *dme-2/dme-2* plants are not due to genetic mutation or heritable aberrant

470   epigenetic effects that escape requirement of DME activity during gametogenesis because their

471   subsequent progeny are phenotypically normal and produces same level (~0.1%) of normal seeds.

21

472       The *DME/dme-2* heterozygous or *dme-2/dme-2* homozygous lines in Col-*gl* background were

473     subjected to Agrobacterium-mediated floral dipping transformation procedures [28]. Seeds were

474     sterilized by 30% bleach solution and screened for T1 transgenic plants on a 0.5x MS nutrient

475     medium with 1% sucrose, 0.8% agar and 40 µg/ml hygromycin. Germinated seedlings were

476     transferred to soil and grown in the growth room under 16 hours of light and 8 hours of dark

477     cycles at 23°C. Siliques from T1 transgenic plants were dissected 14-16 days after self-

478     pollination using a stereoscopic microscope (SteREO Discovery.V12, Carl Zeiss, Wetzlar,

479     Germany). The numbers of viable and aborted seeds in transgenic lines were statistically

480     analyzed with the $\chi$2 test. The probability that deviates from a 1:1 or 3:1 segregation ratio for

481     viable and aborted seeds was also calculated.

482

483     **RNA extraction, cDNA synthesis and quantitative PCR analysis**

484     Total RNA was extracted using TRIzol® Reagent (Invitrogen, Carlsbad, USA) and treated with

485     TURBO DNase (Ambion, Austin TX, USA) according to the manufacturers' instructions. For

486     cDNA synthesis, 5mg of total RNA was reverse-transcribed using Superscript III Reverse

487     Transcriptase and oligo(dT) primer (Invitrogen). The cDNA was treated with RNase H

488     (Invitrogen) at 37oC for 20min and diluted tenfold with H2O. For each 15-µl qPCR reaction, 1µl

489     of diluted cDNA was used. The quantitative PCR was run on ABI 7500 Fast Real-Time PCR

490     System (http://www.appliedbiosystems.com) using FastStart Universal SYBR Green Master Mix

491     (Roche, http://www.roche.com). The quantitative PCR primers are listed in Supplementary File 1.

492     The Ct values were normalized against *ACT2* (*At3g18780*) mRNA or *UBC* (*At5g25760*) mRNA.

493     The abundance of mRNAs was expressed as relative to controls, with control values set to 1. The

494     error bars represent the standard deviation of 4 biological replicates.

495

**Protein domain analysis and phylogenetic inference**

497 We utilized a domain-centric computational strategy to study DME and its related proteins.

498 Specifically, we identify DME homologs by using the iterative profile searches with PSI-BLAST

499 [47] from the protein non-redundant (NR) database at National Center for Biotechnology

500 Information (NCBI). Multiple sequence alignments were built by the Promals [48] program,

501 followed by careful manual adjustments. Consensus secondary structures were predicted using

502 the PSIPRED [49] JPred program [50]. Conserved domains were further characterized based on the

503 comparison to available domain models from pfam [51] and sequence/structural features. The

504 PhyML program [52] was used to determine the maximum-likelihood tree using the Jones–Taylor–

505 Thornton (JTT) model for amino acids substitution with a discrete gamma model (four categories

506 with gamma shape parameter: 1.096). The tree was rendered using MEGA Tree Explorer [53].

507

**Acknowledgments**

**Author contribution.**

23

518    C.Z., Y.-H.H., X.-Q.Z., J.H.H. and T.-F.H. designed the research. C.Z., Y.-H.H., X.-Q.Z.

519    performed the experiments. D.Z., L.M.I, and L.A. performed the evolutionary analysis. C.Z., Y.-

520    H.H., and T.-F.H. wrote the article. T.-F.H., C.Z., Y.-H.H., W.X., J.H.H. interpreted and

521    commented the article.


522    **References**

523    1.    Yang WC, Shi DQ, Chen YH. Female gametophyte development in flowering
524          plants. *Annu Rev Plant Biol* **61**, 89-108 (2010).
525

526    2.    Gehring M. Genomic imprinting: insights from plants. *Annu Rev Genet* **47**, 187-
527          208 (2013).
528

529    3.    Kohler C, Wolff P, Spillane C. Epigenetic mechanisms underlying genomic
530          imprinting in plants. *Annu Rev Plant Biol* **63**, 331-352 (2012).
531

532    4.    Jullien PE, Kinoshita T, Ohad N, Berger F. Maintenance of DNA Methylation
533          during the Arabidopsis Life Cycle Is Essential for Parental Imprinting. *Plant Cell*
534          **18**, 1360-1372 (2006).
535

536    5.    Gehring M, *et al.* DEMETER DNA glycosylase establishes MEDEA polycomb
537          gene self-imprinting by allele-specific demethylation. *Cell* **124**, 495-506 (2006).
538

539    6.    Xiao W, *et al.* Imprinting of the MEA Polycomb gene is controlled by antagonism
540          between MET1 methyltransferase and DME glycosylase. *Dev Cell* **5**, 891-901
541          (2003).
542

543    7.    Kinoshita T, *et al.* One-way control of FWA imprinting in Arabidopsis endosperm
544          by DNA methylation. *Science* **303**, 521-523 (2004).
545

546   8.    Tiwari S, *et al.* MATERNALLY EXPRESSED PAB C-TERMINAL, a novel
547         imprinted gene in Arabidopsis, encodes the conserved C-terminal domain of
548         polyadenylate binding proteins. *Plant Cell* **20**, 2387-2398 (2008).
549

550   9.    Choi Y, *et al.* DEMETER, a DNA Glycosylase Domain Protein, Is Required for
551         Endosperm Gene Imprinting and Seed Viability in *Arabidopsis. Cell* **110**, 33-42
552         (2002).
553

554   10.   Penterman J, Zilberman D, Huh JH, Ballinger T, Henikoff S, Fischer RL. DNA
555         demethylation in the Arabidopsis genome. *Proceedings of the National Academy*
556         *of Sciences of the United States of America* **104**, 6752-6757 (2007).
557

558   11.   Lister R, *et al.* Highly integrated single-base resolution maps of the Arabidopsis
559         genome. *Cell* **133**, 395-397 (2008).
560

561   12.   Gong Z, Morales-Ruiz T, Ariza RR, Roldan-Arjona T, David L, Zhu J-K. ROS1, a
562         Repressor of Transcriptional Gene Silencing in Arabidopsis, Encodes a DNA
563         Glycosylase/Lyase. *Cell* **111**, 803-814 (2002).
564

565   13.   Ibarra CA, *et al.* Active DNA demethylation in plant companion cells reinforces
566         transposon methylation in gametes. *Science* **337**, 1360-1364 (2012).
567

568   14.   Hsieh TF, *et al.* Regulation of imprinted gene expression in Arabidopsis
569         endosperm. *Proceedings of the National Academy of Sciences of the United*
570         *States of America* **108**, 1755-1762 (2011).
571

572   15.   Hsieh TF, *et al.* Genome-wide demethylation of Arabidopsis endosperm. *Science*
573         **324**, 1451-1454 (2009).
574

575   16.   Wang X, *et al.* RNA-binding protein regulates plant DNA methylation by

576    controlling mRNA processing at the intronic heterochromatin-containing gene
577    IBM1. *Proceedings of the National Academy of Sciences of the United States of*
578    *America* **110**, 15467-15472 (2013).

579

580  17.  Lei M, Zhang H, Julian R, Tang K, Xie S, Zhu JK. Regulatory link between DNA
581    methylation and active demethylation in Arabidopsis. *Proceedings of the National*
582    *Academy of Sciences of the United States of America* **112**, 3553-3557 (2015).

583

584  18.  Lang Z, *et al.* The methyl-CpG-binding protein MBD7 facilitates active DNA
585    demethylation to limit DNA hyper-methylation and transcriptional gene silencing.
586    *Mol Cell* **57**, 971-983 (2015).

587

588  19.  Qian W, *et al.* A histone acetyltransferase regulates active DNA demethylation in
589    Arabidopsis. *Science* **336**, 1445-1448 (2012).

590

591  20.  Ikeda Y, *et al.* HMG domain containing SSRP1 is required for DNA demethylation
592    and genomic imprinting in Arabidopsis. *Dev Cell* **21**, 589-596 (2011).

593

594  21.  Bustin M, Catez F, Lim JH. The dynamics of histone H1 function in chromatin.
595    *Mol Cell* **17**, 617-620 (2005).

596

597  22.  Fan Y, *et al.* Histone H1 depletion in mammals alters global chromatin structure
598    but causes specific changes in gene regulation. *Cell* **123**, 1199-1212 (2005).

599

600  23.  Graziano V, Gerchman SE, Schneider DK, Ramakrishnan V. Histone H1 is
601    located in the interior of the chromatin 30-nm filament. *Nature* **368**, 351-354
602    (1994).

603

604  24.  Hashimoto H, *et al.* Histone H1 null vertebrate cells exhibit altered nucleosome
605    architecture. *Nucleic Acids Res* **38**, 3533-3545 (2010).

606

607   25.   Rea M, *et al.* Histone H1 affects gene imprinting and DNA methylation in
608         Arabidopsis. *Plant J* **71**, 776-786 (2012).

609

610   26.   Brooks SC, Fischer RL, Huh JH, Eichman BF. 5-methylcytosine recognition by
611         Arabidopsis thaliana DNA glycosylases DEMETER and DML3. *Biochemistry* **53**,
612         2525-2532 (2014).

613

614   27.   Jang H, Shin H, Eichman BF, Huh JH. Excision of 5-hydroxymethylcytosine by
615         DEMETER family DNA glycosylases. *Biochem Biophys Res Commun* **446**, 1067-
616         1072 (2014).

617

618   28.   Clough SJ, Bent AF. Floral dip: a simplified method for Agrobacterium-mediated
619         transformation of Arabidopsis thaliana. *Plant J* **16**, 735-743 (1998).

620

621   29.   Park JS, *et al.* Control of DEMETER DNA demethylase gene transcription in
622         male and female gamete companion cells in Arabidopsis thaliana. *Proceedings*
623         *of the National Academy of Sciences of the United States of America* **114**, 2078-
624         2083 (2017).

625

626   30.   Kohler C, Hennig L, Bouveret R, Gheyselinck J, Grossniklaus U, Gruissem W.
627         Arabidopsis MSI1 is a component of the MEA/FIE Polycomb group complex and
628         required for seed development. *The EMBO journal* **22**, 4804-4814 (2003).

629

630   31.   Grossniklaus U, Vielle-Calzada J-P, Hoeppner MA, Gagliano WB. Maternal
631         control of embryogenesis by *MEDEA*, a polycomb-group gene in *Arabidopsis*.
632         *Science* **280**, 446-450 (1998).

633

634   32.   Luo M, Bilodeau P, Dennis ES, Peacock WJ, Chaudhury A. Expression and
635         parent-of-origin effects for FIS2, MEA, and FIE in the endosperm and embryo of

636      developing Arabidopsis seeds. *Proceedings of the National Academy of*
637      *Sciences of the United States of America* **97**, 10637-10642 (2000).

638

639  33.  Schoft VK, *et al.* Function of the DEMETER DNA glycosylase in the Arabidopsis
640      thaliana male gametophyte. *Proceedings of the National Academy of Sciences of*
641      *the United States of America* **108**, 8042-8047 (2011).

642

643  34.  Lang Z, *et al.* Critical roles of DNA demethylation in the activation of ripening-
644      induced genes and inhibition of ripening-repressed genes in tomato fruit.
645      *Proceedings of the National Academy of Sciences of the United States of*
646      *America* **114**, E4511-E4519 (2017).

647

648  35.  Choi Y, Harada JJ, Goldberg RB, Fischer RL. An invariant aspartic acid in the
649      DNA glycosylase domain of DEMETER is necessary for transcriptional activation
650      of the imprinted MEDEA gene. *Proceedings of the National Academy of Sciences*
651      *of the United States of America* **101**, 7481-7486 (2004).

652

653  36.  Iyer LM, Abhiman S, Aravind L. Natural history of eukaryotic DNA methylation
654      systems. *Progress in molecular biology and translational science* **101**, 25-104
655      (2011).

656

657  37.  Walsh P, Bursac D, Law YC, Cyr D, Lithgow T. The J-protein family: modulating
658      protein assembly, disassembly and translocation. *EMBO reports* **5**, 567-571
659      (2004).

660

661  38.  Hong S, Hashimoto H, Kow YW, Zhang X, Cheng X. The Carboxy-Terminal
662      Domain of ROS1 Is Essential for 5-Methylcytosine DNA Glycosylase Activity. *J*
663      *Mol Biol*, (2014).

664

665  39.  Long HK, Blackledge NP, Klose RJ. ZF-CxxC domain-containing proteins, CpG

666        islands and the chromatin connection. *Biochem Soc Trans* **41**, 727-740 (2013).
667

668    40.    Mok YG, *et al.* Domain structure of the DEMETER 5-methylcytosine DNA
669        glycosylase. *Proceedings of the National Academy of Sciences of the United*
670        *States of America* **107**, 19225-19230 (2010).
671

672    41.    Zheng X, *et al.* ROS3 is an RNA-binding protein required for DNA demethylation
673        in Arabidopsis. *Nature* **455**, 1259-1262 (2008).
674

675    42.    Ponferrada-Marin MI, Roldan-Arjona T, Ariza RR. Demethylation initiated by
676        ROS1 glycosylase involves random sliding along DNA. *Nucleic Acids Res* **40**,
677        11554-11562 (2012).
678

679    43.    Ponferrada-Marin MI, Martinez-Macias MI, Morales-Ruiz T, Roldan-Arjona T,
680        Ariza RR. Methylation-independent DNA binding modulates specificity of
681        Repressor of Silencing 1 (ROS1) and facilitates demethylation in long substrates.
682        *The Journal of biological chemistry* **285**, 23032-23039 (2010).
683

684    44.    Hsieh TF. Whole-genome DNA methylation profiling with nucleotide resolution.
685        *Methods in molecular biology* **1284**, 27-40 (2015).
686

687    45.    Pignatta D, Erdmann RM, Scheer E, Picard CL, Bell GW, Gehring M. Natural
688        epigenetic polymorphisms lead to intraspecific variation in Arabidopsis gene
689        imprinting. *Elife* **3**, e03198 (2014).
690

691    46.    Hehenberger E, Kradolfer D, Kohler C. Endosperm cellularization defines an
692        important developmental transition for embryo development. *Development* **139**,
693        2031-2039 (2012).
694

695    47.    Altschul SF, *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein

696         database search programs. *Nucleic Acids Res* **25**, 3389-3402 (1997).

697

698  48.    Pei J, Grishin NV. PROMALS: towards accurate multiple sequence alignments of

699         distantly related proteins. *Bioinformatics* **23**, 802-808 (2007).

700

701  49.    Buchan DW, Minneci F, Nugent TC, Bryson K, Jones DT. Scalable web services

702         for the PSIPRED Protein Analysis Workbench. *Nucleic Acids Res* **41**, W349-357

703         (2013).

704

705  50.    Cuff JA, Clamp ME, Siddiqui AS, Finlay M, Barton GJ. JPred: a consensus

706         secondary structure prediction server. *Bioinformatics* **14**, 892-893 (1998).

707

708  51.    Finn RD*, et al.* The Pfam protein families database: towards a more sustainable

709         future. *Nucleic Acids Res* **44**, D279-285 (2016).

710

711  52.    Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New

712         algorithms and methods to estimate maximum-likelihood phylogenies: assessing

713         the performance of PhyML 3.0. *Syst Biol* **59**, 307-321 (2010).

714

715  53.    Tamura K, Dudley J, Nei M, Kumar S. MEGA4: Molecular Evolutionary Genetics

716         Analysis (MEGA) software version 4.0. *Mol Biol Evol* **24**, 1596-1599 (2007).

717
718
719
720 **Figure Legends**
721

722 **Figure 1 Complementation of *dme* seed abortion phenotype by the truncated DME nAGB.**

723 (a) Siliques were dissected and photographed 14 days after self-pollination. In *dme-2/dme-2*

724 silique greater than 99% of seeds are aborted. A single copy of *nDME^CTD* transgene reduces seed

725 abortion rate to 50%; and in the *dme-2/dme-2; nDME^CTD/nDME^CTD* silique, all the *dme-2* seeds

726    are rescued and developed normally. Scale bar = 0.5 mm. (b) Complementation of *dme-2* seed

727    abortion phenotype by *nDME^{CTD} and DME^{FL}*. (c) The *nDME^{CTD}* transgene restores DME target

728    genes *FWA* and *FIS2* expression. WT: Col-0; *nDME^{CTD}*: *dme-2/dme-2; nDME^{CTD}/ nDME^{CTD}*;

729    *dme-2*: *dme-2/dme-2*. Total RNA was isolated from stage F1 to F12 floral buds.

730

731    **Figure 2 Endosperm methylome analysis.** (a) Genome browser snapshots of CG DNA

732    methylation at selected imprinted gene loci. Top two tracks are coding genes (magenta) and TEs

733    (orange) with Tair10 chromosome coordinates. For the bottom seven tracks, each track

734    represents fractional CG methylation levels for different genotype: black trace, *dme-2* endosperm;

735    dark green trace, WT endosperm; dark blue trace, *DME^{FL}*-complemented endosperm; dark purple

736    trace, *nDME^{CTD}*-complemented endosperm; light green trace, WT endosperm subtracted from

737    *dme-2* mutant endosperm; light blue trace, *DME^{FL}*-complemented endosperm subtracted from

738    *dme-2* endosperm; light purple trace, *nDME^{CTD}*-complemented endosperm subtracted form *dme-*

739    *2* endosperm. DNA CG hypomethylation at selected maternally expressed (*FIS2* and *SDC*) and

740    paternally expressed (*SUV7*, *YUC10, and PHE1*) imprinted genes is restored in *DME^{FL}*- and

741    *nDME^{CTD}*-complemented endosperm. (b) Boxplot of CG methylation levels among canonical

742    DME target sites in *dme-2* mutant (grey), WT (white), *DME^{FL}*- (blue), or *nDME^{CTD}* - (red)

743    complemented endosperm. (c) Venn Diagram (top panel) of CG hyper-DMRs in 50-bp windows

744    between *dme-2* endosperm relative to WT, *DME^{FL}*-complemented or *nDME^{CTD}*-complemented

745    endosperm. Boxplot (bottom panel) of CG methylation levels in *dme-2* mutant (grey), WT

746    (white), *DME^{FL}*- (blue) or *nDME^{CTD}*- (red) complemented endosperm in WT only (left panel),

747    *DME^{FL}* only, or *nDME^{CTD}* only (right panel) DMRs.

748

749 **Figure 3 Expression of DME NTD region in wild-type central cell induces *dme*-like seed**

750 **abortion phenotype. (a)** Confocal microscopy image of ovule in F12 floral bud shows the

751 expression of mDME$^{NTD}$-GFP in the central cell. Scale bar, 20 µm. **(b-c)** Ectopic expression of

752 *DME$^{NTD}$* in WT central cell induces *dme-2* like seed abortion phenotype in silique (**b**) and in

753 developing seeds (**c**). Total RNA was isolated from stage F1 to F12 floral buds from independent

754 lines with different seed abortion ratios (**d**) to assess transgene and endogenous DME expression.

755 (**e**) Endogenous DME transcript levels in independent transgenic lines are comparable to the

756 control line, but the transgene expression level varies among these independent lines with

757 different seed abortion rates. Error bars indicate SD. NS, $p > 0.2$ (Ctrl vs 23), $p > 0.5$ (Ctrl vs 15),

758 $p > 0.3$ (Ctrl vs 25), $p > 0.4$ (Ctrl vs 8), not significant (two-tailed t test). (**f**) Correlation analysis

759 shows that the transcript abundance of the transgene, but not that of the endogenous DME

760 transcripts, correlates with seed abortion rates (by linear regression).

761

762 **Figure 4. Evolution of plant DME-like proteins.** A phylogenetic tree was reconstructed using

763 the PhyML program. Only node supporting values >0.80 from ML bootstrap analyses are shown.

764 The representative domain architectures of DME homologs in major plant clades are shown

765 along the tree, demonstrating domain fusions during evolution. Domain abbreviations: DemeN,

766 N-terminal domain of DEME-like proteins in angiosperms; DnaJ, DnaJ molecular chaperone

767 homology domain (Pfam: PF00226); FCL, [Fe4S4] cluster loop motif (also called Iron-sulfur

768 binding domain of endonuclease III; Pfam: PF10576); HhH-GL, HhH-GPD superfamily base

769 excision DNA repair protein (Pfam: PF00730); PHD, PHD finger (Pfam: PF00628); RRM, RNA

770 recognition motif (Pfam: PF00076); Tudor, Tudor domain (Pfam: PF00567).

771

772    **Supplemental Information**

773    **Figure Legends**

774    **Fig. S1. Diagrams of DME protein structure and transgene constructs.**

775    (a) DME protein domain architecture. The positions of conserved domains along DME protein.

776    Numbers represent amino acid position relative to the translation start sites. DME.1 is shorter

777    than DME.2 by 258 amino acids due to alternative splicing, missing the very N-terminal DemeN

778    domain. DemeN is a domain of unknown function conserved among angiosperm DME-like

779    protein. 3DR is the stretch of basic rich amino acid direct repeats, resembling AT-hook motifs,

780    and serves as a nuclear localization signal; per-CXXC is the permuted CXXC zinc finger motif;

781    RRM is the RNA recognition motif; FCL is a [Fe4S4] cluster loop following the HhH module.

782    The *dme-2* allele harbors a T-DNA insertion in region A at amino acid position 1012. ID1 and

783    ID2 are variable, low complexity sequences between the glycosylase domain and the conserved

784    B region. (b) Transgene constructs used in this study. DMEpro refers to the upstream regulatory

785    sequence (2895 bp upstream of the translation start codon ATG) of DME.2. SV40NLS:

786    PKKKRKV. A polypeptide linker comprising 6 alanine residues is placed between any protein

787    fragment fusions.

788

789    **Fig. S2. DNA methylomes of three independent $nDME^{CTD}$-complemented *dme-2* endosperm.**

790    (a) Venn diagram showing partial overlap of *dme* CG hyper-DMRs relatives to each nAGB-

791    complemented endosperm ($nDME^{CTD}$-*1* to $nDME^{CTD}$-3). (b) Boxplot of CG methylation levels

792    among canonical DME target sites in *dme-2* mutant (black), $nDME^{CTD}$-1 (pink), $nDME^{CTD}$-*2*

793    (magenta), or $nDME^{CTD}$-*3* (red) complemented endosperm, in $nDME^{CTD}$-*1* specific (left panel),

794    $nDME^{CTD}$-*2* specific (middle panel), and $nDME^{CTD}$-*3* specific DMRs. These results show that the

33

795     combined DMRs are more or less hypomethylated in each independent line compared to *dme-2*

796     endosperm.

797

798     **Fig. S3. DNA methylomes of three independent DME$^{FL-}$complemented *dme-2* endosperm.** (a)

799     Venn diagram showing partial-overlap of *dme* CG hyper-DMRs relatives to each *DME$^{FL}$*-

800     complemented endosperm (*DME$^{FL}$-1* to *DME$^{FL}$-3*). (b) Boxplot of CG methylation levels among

801     canonical DME target sites in *dme-2* mutant (black), *DME$^{FL}$*-1 (light blue), *DME$^{FL}$*-2 (medium

802     blue), or *DME$^{FL}$-3* (dark blue) complemented endosperm, in *DME$^{FL}$-1* specific (left panel),

803     *DME$^{FL}$-2* specific (middle panel), and *DME$^{FL}$-3* specific DMRs. These results show that the

804     combined DMRs are more or less hypomethylated in each independent line compared to *dme-2*

805     endosperm.

806

807     **Fig. S4. The DMRs of *dme* relative to WT endosperm or nDME$^{CTD}$- complemented**

808     **endosperm.** Venn Diagram (top) and Boxplot analysis (bottom) of CG hyper-DMRs in 50-bp

809     windows between *dme-2* endosperm relative to *nDME$^{CTD}$*-complemented or WT endosperm. CG

810     methylation levels of DMRs unique to *nDME$^{CTD}$*-complemented endosperm are also

811     demethylated in the WT endosperm (left panel). Similarly, DMRs unique to WT endosperm are

812     demethylated in *nDME$^{CTD}$*-complemented endosperm (right).

813

814     **Fig. S5. DME$^{FL}$ and nDME$^{CTD}$ transgenes are expressed at comparable levels among**

815     **independent complementation lines.** *DME$^{FL}$* and *nDME$^{CTD}$* expression levels are comparable

816     between the four of the six complementation lines used in the methylome study. Total RNA was

817     isolated from stage F1 to F12 floral buds. The results show that there is no significant difference

34

818    in expression level between these two transgenes(t-test, *p*>0.4).

819

820    **Fig. S6. The effects of T-DNA insertion on endogenous DME transcript abundance in *dme-***

821    ***2/dme-2* plants.** Total RNA was isolated from stage F1 to F12 floral buds. Equal amount of total

822    RNA from WT and *dme-2/dme-2* were used for reverse transcription and quantitative PCR. Six

823    paired of primers (PN1-PN6) correspond to the N-terminal region before the T-DNA insertion

824    site, and three pairs of C-terminal region primers (PC1-PC3) were used to assess endogenous

825    DME transcript level in *dme-2/dme-2* mutant plants. The position of each primer pair is indicated

826    in the DME diagram where T-DNA insertion site is shown.

827

828    **Fig. S7. Alignment of angiosperm DME-like proteins showing the conserved DemeN**

829    **domain and the basic rich 3DR repeats.** Bioinformatics analysis using available DME-like

830    sequences identified a ~ 120-amino-acid-long conserved region at the very N-termini among

831    DME-like proteins in angiosperms. This sequence is characterized by a highly conserved

832    WxPxTPxK motif that might function in protein-protein interactions. Further toward the C-

833    terminus is a stretch of basic amino acids rich region that serves as a nuclear localization signal.

834    This sequence consists of three direct repeats (3DR) reminiscent of the AT-hook motifs that may

835    bind DNA.

836

837    **Fig. S8.** Boxplot of CG methylation levels among canonical DME target sites in different DMR

838    length category, in *dme-2* mutant (black), wild-type (white), or *nDME^{CTD}*-complemented

839    endosperm

840

841     **Fig. S9.** (a) Merged DMR length distribution in WT and $nDME^{CTD}$-complemented endosperm. (b)

842     Genome Browser examples of long WT DMRs. Tracks are as labeled. The DMR regions are

843     indicated as horizontal bars according to their length in each sample (bottom two tracks). Even

844     though $nDME^{CTD}$ complemented endosperm lack longer DMRs, these regions are also shorter

845     DMRs in $nDME^{CTD}$-complemented endosperm.

846

847

848

849

850 **Table 1.  Rescue of the reduced paternal *dme-2* allele transmission by the *nDME^{CTD}***
851 **transgene.**
852

| Female | Male parent | F1, DME/*dme-2* | F1, *DME/dme-2*; *nDME^{CTD}* | *nDME^{CTD}* transmission rate (%) | *p* for 1:1† |
|--------|-------------|-----------------|-------------------------------|--------------------------------|------------|
| Col-0 | *dme-2/dme-2; nDME^{CTD}* /~ Line 1 | 32 | 62 | 66 | 2.0E-3 |
| Col-0 | *dme-2/dme-2; nDME^{CTD}* /~ Line 2 | 3 | 50 | 94.3 | 1.1E-10 |
| Col-0 | *dme-2/dme-2; nDME^{CTD}* /~ Line 3 | 8 | 34 | 81 | 6.0E-5 |
| Col-0 | *dme-2/dme-2; nDME^{CTD}* /~ Line 4 | 9 | 44 | 83 | 1.5E-6 |
| † Probability  that that the deviation from the indicated segregation ration (1:1 inheritance of paternal genome with or without *nDME^{CTD}* transgene in the F1 generation) is due to chance. | | | | | |

853

# Figure 1

## a



*dme-2/dme-2*

*dme-2/dme-2; pDME:nDME^{CTD}/~*

*dme-2/dme-2; pDME:nDME^{CTD}/ pDME:nDME^{CTD}*

## b



Proportion of viable seeds (%)

## c



*FWA*

*FIS2*

# Figure 2

# Figure 3

## a



## b



*DME/DME*

*DME/dme-2*

*DME/DME ;*
*DME^{NTD} /~*

## c



*DME/DME*　　*dme-2/dme-2*

*DME/DME; DME^{NTD}/~*

## d

| Sample | Proportion of aborted seeds (%) |
|---|---|
| Control | 0 |
| Line 23 | 52 |
| Line 15 | 27 |
| Line 25 | 0 |
| Line 8 | 0 |

## e



NS

Relative expression level

Ctrl　　23　　15　　25　　8

☐ Endogenous DME　　■ Transgene

## f



Relative expression level

$y = -0.0008x + 0.9979$
$R^2 = 0.03$

$y = 0.0147x + 0.1945$
$R^2 = 0.98$

Seed abortion rate (%)

☐ Endogenous DME　　■ Transgene
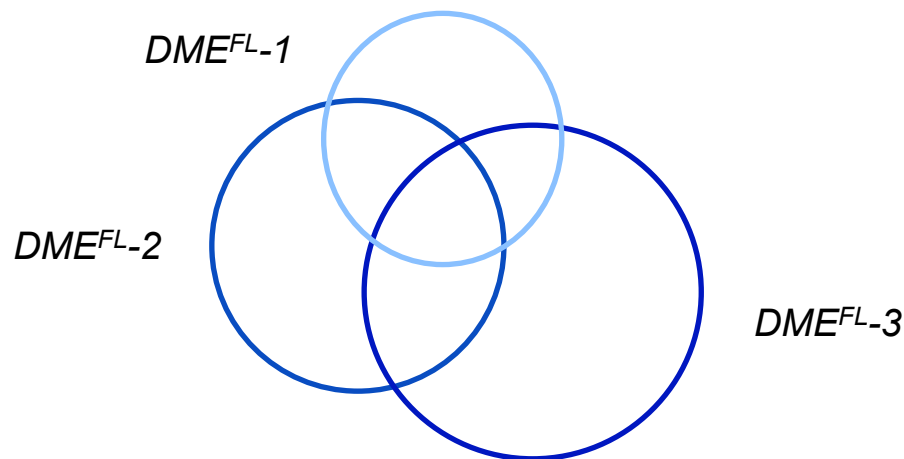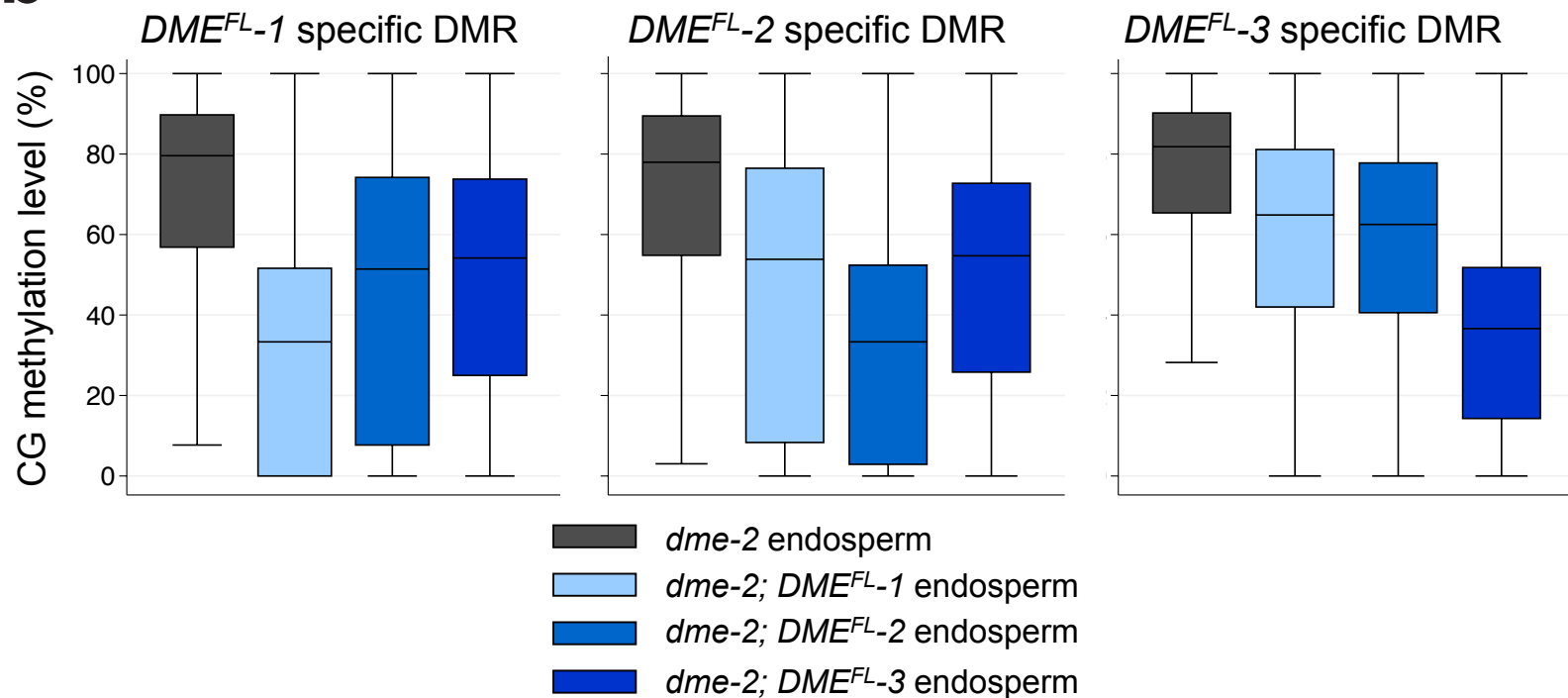
# Figure 4

# Supplementary Figure S1

# Supplementary Figure S2



**a**

$nDME^{CTD}$-1

$nDME^{CTD}$-2

$nDME^{CTD}$-3

**b**

$nDME^{CTD}$-1 specific DMR

$nDME^{CTD}$-2 specific DMR

$nDME^{CTD}$-3 specific DMR

CG methylation level (%)

*dme-2* endosperm

*dme-2; nDME$^{CTD}$-1* endosperm

*dme-2; nDME$^{CTD}$-2* endosperm

*dme-2; nDME$^{CTD}$-3* endosperm

# Supplementary Figure S3

**a**



**b**

# Supplementary Figure S4



CG hyper-DMRs
*dme-2* vs *nDME^CTD*

CG hyper-DMRs
*dme-2* vs WT

21421   18731   33425

CG methylation level (%)

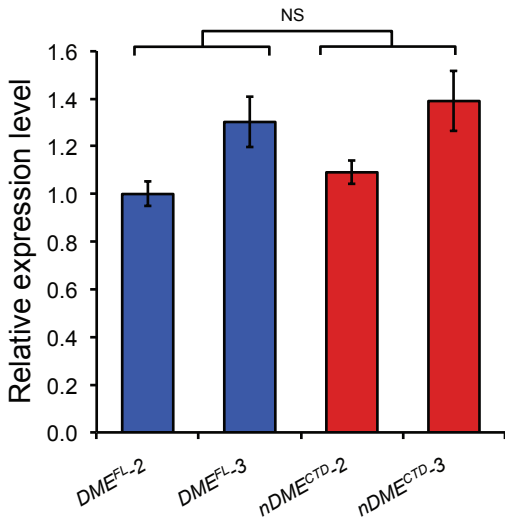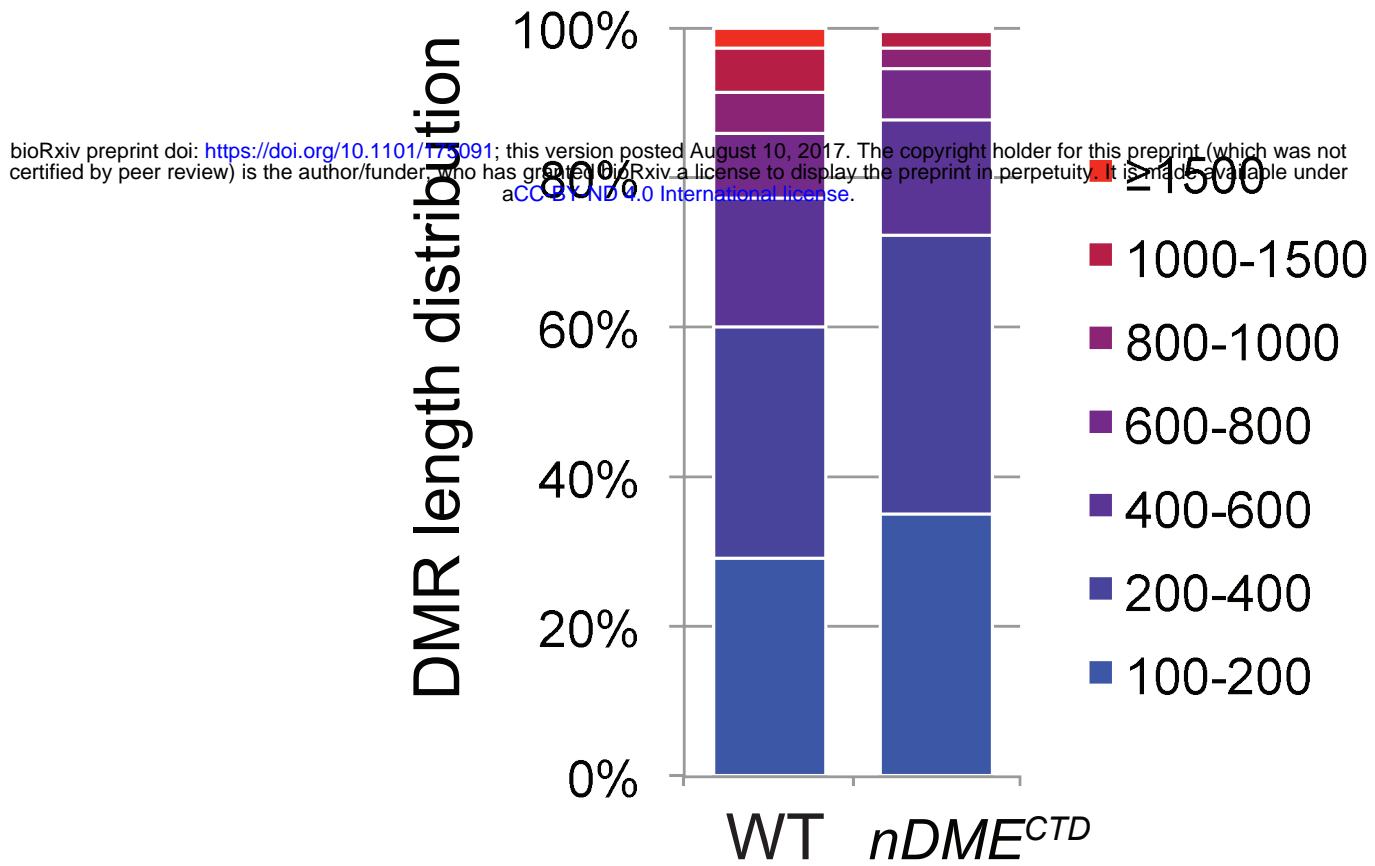■ *dme-2* endosperm
□ WT endosperm
■ *dme-2; nDME^CTD* endosperm

# Supplementary Figure S5

# Supplementary Figure S6

# Supplementary Figure S7



WxPxTPxK motif

Basic amino acid rich stretch

**Supplementary Figure S8**



*Legend:*
- *dme-2* endosperm
- WT endosperm
- *dme-2; nDME$^{CTD}$* endosperm

x-axis: WT DMR size (bp)
y-axis: CG methylation level (%)

# Supplementary Figure S9

## a

## b