

On the origin and evolutionary consequences of gene body DNA methylation

Adam J. Bewick^{a,1}, Lexiang Ji^{b,1}, Chad E. Niederhuth^{a,1}, Eva-Maria Willing^{c,1}, Brigitte T. Hofmeister^b, Xiuling Shi^a, Li Wang^d, Zefu Lu^a, Nicholas A. Rohr^a, Benjamin Hartwig^c, Christiane Kiefer^c, Roger B. Deal^e, Jeremy Schmutz^f, Jane Grimwood^f, Hume Stroud^g, Steven E. Jacobsen^{g,h}, Korbinian Schneeberger^c, Xiaoyu Zhang^d, and Robert J. Schmitz^{a,2}

^aDepartment of Genetics, University of Georgia, Athens, GA 30602; ^bInstitute of Bioinformatics, University of Georgia, Athens, GA 30602; ^cDepartment of Plant Developmental Biology, Max Planck Institute for Plant Breeding Research, 50829 Cologne, Germany; ^dDepartment of Plant Biology, University of Georgia, Athens, GA 30602; ^eDepartment of Biology, Emory University, Atlanta, GA 30322; ^fHudson Alpha Genome Sequencing Center, Huntsville, AL 35806; ^gDepartment of Molecular, Cell, and Developmental Biology, University of California, Los Angeles, CA 90095; and ^hHoward Hughes Medical Institute, University of California, Los Angeles, CA 90095

Edited by David C. Baulcombe, University of Cambridge, Cambridge, United Kingdom, and approved June 10, 2016 (received for review March 22, 2016)

In plants, CG DNA methylation is prevalent in the transcribed regions of many constitutively expressed genes (gene body methylation; gbM), but the origin and function of gbM remain unknown. Here we report the discovery that *Eutrema salsugineum* has lost gbM from its genome, to our knowledge the first instance for an angiosperm. Of all known DNA methyltransferases, only CHROMOMETHYLASE 3 (CMT3) is missing from *E. salsugineum*. Identification of an additional angiosperm, *Conringia planisiliqua*, which independently lost CMT3 and gbM, supports that CMT3 is required for the establishment of gbM. Detailed analyses of gene expression, the histone variant H2A.Z, and various histone modifications in *E. salsugineum* and in *Arabidopsis thaliana* epigenetic recombinant inbred lines found no evidence in support of any role for gbM in regulating transcription or affecting the composition and modification of chromatin over evolutionary timescales.

DNA methylation | gene body methylation | epigenetics | histone modifications | CHROMOMETHYLASE 3

In angiosperms, cytosine DNA methylation occurs in three sequence contexts: Methylated CG (mCG) is catalyzed by METHYLTRANSFERASE 1 (MET1), mCHG (where H is A/C/T) by CHROMOMETHYLASE 3 (CMT3), and mCHH by DOMAINS REARRANGED METHYLTRANSFERASE 2 (DRM2) or CHROMOMETHYLASE 2 (CMT2) (1). MET1 performs a maintenance function and is targeted by VARIANT IN METHYLATION 1 (VIM1), which binds preexisting hemimethylated CG sites. In contrast, DRM2 is targeted by RNA-directed DNA methylation (RdDM) for the de novo establishment of mCHH. CMT3 forms a self-reinforcing loop with the H3K9me2 pathway to maintain mCHG; however, considering that transformation of CMT3 into the *cmt3* background can rescue DNA methylation defects, it is reasonable to also consider CMT3 a de novo methyltransferase (2). Two main lines of evidence suggest that DNA methylation plays an important role in the transcriptional silencing of transposable elements (TEs): that TEs are usually methylated, and that the loss of DNA methylation (e.g., in methyltransferase mutants) is often accompanied by TE reactivation.

A large number of plant genes (e.g., ~13.5% of all *Arabidopsis thaliana* genes) also contain exclusively mCG in the transcribed region and a depletion of mCG from both the transcriptional start and stop sites (referred to as “gene body DNA methylation”; gbM) (Fig. 1A) (3–5). A survey of plant methylome data showed that the emergence of gbM in the plant kingdom is specific to angiosperms (6), whereas nonflowering plants (such as mosses and green algae) have much more diverse genic methylation patterns (7, 8). Similar to mCG at TEs, the maintenance of gbM requires MET1. In contrast to DNA methylation at TEs, however, gbM does not appear to be associated with transcriptional repression. Rather, genes containing gbM are ubiquitously expressed at moderate to high levels compared with non-gbM genes (4, 5, 9), and within gbM genes there is a correlation between transcript abundance and methylation levels (10, 11).

It has been proposed that gbM may be established by the de novo methylation activity of the RdDM pathway and subsequently maintained by MET1 independent of RdDM. In this “de novo” scenario, occasional antisense transcripts could form double-stranded RNA by pairing with sense transcripts, which could trigger the production of small interfering RNAs (siRNAs) to target DRM2 for de novo methylation in gene bodies. Although mechanistically feasible, it is difficult to explain why gbM is absent from many nonangiosperm plants (such as the moss *Physcomitrella patens*) with functional RdDM and MET1 pathways (12).

Alternatively, we propose that the establishment of gbM might involve the self-reinforcing loop between CMT3 and the histone H3 lysine 9 (H3K9) methyltransferase KRYPTONITE/SUVH4 (KYP) (13, 14) in addition to transcription, similar to a model proposed by Inagaki and Kakutani (15). CMT3 is recruited to chromatin by H3K9me2 for DNA methylation, which in turn recruits KYP for H3K9me2. Although mCHG and H3K9me2 are normally limited to heterochromatin, they accumulate ectopically in thousands of actively transcribed genes upon the loss of the H3K9 demethylase INCREASED IN BONSAI METHYLATION 1 (IBM1) (16). It therefore appears likely that mCHG and H3K9me2 also occur constantly (albeit transiently) in actively transcribed genes, but their accumulation is normally prevented by IBM1 (17). The transient presence of H3K9me2 in transcribed regions

Significance

DNA methylation in plants is found at CG, CHG, and CHH sequence contexts. In plants, CG DNA methylation is enriched in the transcribed regions of many constitutively expressed genes (gene body methylation; gbM) and shows correlations with several chromatin modifications. Contrary to other types of DNA methylation, the evolution and function of gbM are largely unknown. Here we show two independent concomitant losses of the DNA methyltransferase CHROMOMETHYLASE 3 (CMT3) and gbM without the predicted disruption of transcription and of modifications to chromatin. This result suggests that CMT3 is required for the establishment of gbM in actively transcribed genes, and that gbM is dispensable for normal transcription as well as for the composition and modification of plant chromatin.

Author contributions: A.J.B., C.E.N., S.E.J., K.S., and R.J.S. designed research; A.J.B., L.J., C.E.N., E.-M.W., B.T.H., X.S., L.W., Z.L., N.A.R., B.H., C.K., J.S., J.G., H.S., and X.Z. performed research; R.B.D., J.S., and J.G. contributed new reagents/analytic tools; A.J.B., L.J., C.E.N., E.-M.W., B.T.H., and R.J.S. analyzed data; and R.J.S. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Data deposition: The sequence reported in this paper has been deposited in the Gene Expression Omnibus (GEO) database, www.ncbi.nlm.nih.gov/geo (accession no. GSE75071).

¹A.J.B., L.J., C.E.N., and E.-M.W. contributed equally to this work.

²To whom correspondence should be addressed. Email: schmitz@uga.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1604666113/-DCSupplemental.

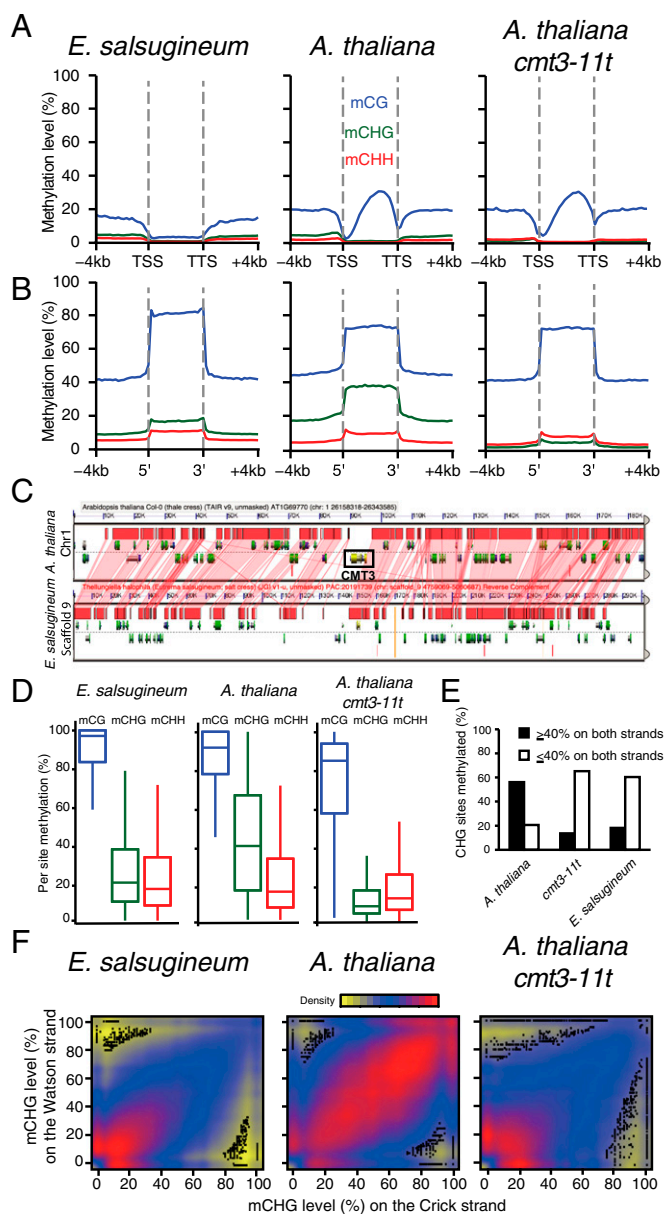


Fig. 1. CMT3 and gbM are absent in *E. salsugineum*. (A and B) Metagene plots of DNA methylation across (A) gene bodies and (B) repeats including 4 kb up- and downstream. TTS, transcriptional termination site. (C) A syntenic block of sequence is shared between *A. thaliana* and *E. salsugineum*. The black box in the *A. thaliana* block indicates the location of CMT3, which is absent in *E. salsugineum*. The red shaded areas indicate regions of shared synteny. (D) Boxplot representation of methylation levels of individual methylated cytosines within each sequence context. (E) A bar plot of methylation levels at symmetric CHG methylated sites in *A. thaliana*, *cmt3-11t*, and *E. salsugineum*. (F) Density plot representation of CMT3-dependent versus CMT3-independent CHG methylation.

could trigger CMT3-dependent methylation of CHG and other contexts. The MET1 pathway would then maintain rare methylation of CG sites in an H3K9me2-independent manner. Consistent with this possibility, gbM-containing genes are preferred targets for hypermethylation in the *ibm1* mutant (16). Last, in support of a role for CMT3 in this model, CMT3 is only present in angiosperms, which coincides with the emergence of gbM in the plant kingdom (6, 17).

The difficulty in addressing the origin of gbM is twofold. First, gbM is unaffected by the loss of RdDM or CMT3 in the short term, indicating that the maintenance activity of MET1 is sufficient for the persistence

of gbM over extended periods of time. Second, once gbM is lost in the *met1* mutant, it does not immediately return when MET1 is reintroduced by crossing, indicating that the establishment of gbM is a stochastic process that requires many generations (18).

Here we describe the results from a comparative epigenomics approach, where we sought to identify natural variation in plant methylomes (SI Appendix, Table S1) that were associated with genetic changes in key genes in DNA methylation pathways. The methylomes of the vast majority of the plant species are similar to *A. thaliana*, with high levels of mCG/mCHG/mCHH colocalized to repetitive sequences and gbM in moderately expressed genes (19). The only exception was *Eutrema salsugineum* (accession Shandong), a member of the Brassicaceae family that shares a common ancestor with *A. thaliana* and *Brassica* spp. ~47 and 40 million y ago, respectively (20). Comparisons between the *E. salsugineum* methylome and those of other plants revealed two major differences. First, *E. salsugineum* has lost gbM (Fig. 1A). In contrast to other species where thousands of active genes contained gbM (e.g., 4,934 in *A. thaliana*), only 103 *E. salsugineum* genes contained gbM based on our identification criteria (Methods and Dataset S1). A closer inspection of these 103 loci revealed that the distribution of mCG in these loci was not representative of gbM genes in other angiosperms (Fig. 2). Importantly, mCG was present at high levels in repetitive sequences in *E. salsugineum*, indicating that the absence of gbM in its genome was not due to the loss of MET1 activity (Fig. 1B). Second, the mCHG level in *E. salsugineum* repetitive sequences was much lower compared with other plant species (Fig. 1B). To further validate these results, we performed MethylC sequencing (MethylC-seq) using an additional *E. salsugineum* accession (Yukon), and the results showed that it too has lost gbM (Fig. 2 and SI Appendix, Fig. S1). A detailed analysis of *E. salsugineum* genes identified homologs of all known DNA methyltransferases in *A. thaliana* (e.g., MET1, CMT2, and DRM2), with the exception of CMT3. In addition, a comparison of the genomic regions between *E. salsugineum* syntenic to the *A. thaliana* CMT3 locus found no evidence for CMT3-related sequences at the syntenic location (Fig. 1C).

The methylome of *E. salsugineum*, with lower mCHG levels in repeats and a complete loss of gbM, was unique compared with 86 *A. thaliana* mutants for which methylome data are available (18). The absence of CMT3 from *E. salsugineum* is consistent with two characteristics of mCHG in its genome. First, the methylation level at individual CHG sites was significantly lower than any other species and was similar to the *A. thaliana cmt3* mutant (Fig. 1D). Second, because CHG is symmetrical, with a mirrored cytosine on the opposing strand, CMT3 activity results in high methylation of cytosines on both strands. In *E. salsugineum* the percentage of paired CHG sites that were highly methylated is significantly lower than wild-type *A. thaliana* and again similar to the *cmt3* mutant, suggesting that the mCHG in *E. salsugineum* is likely a result of RdDM activity (Fig. 1E and F). Taken together, these results indicated that *E. salsugineum* does not have CMT3 activity.

The loss of CMT3 and gbM from *E. salsugineum* is consistent with the hypothesis that CMT3 is required for the establishment of gbM. To solidify this connection, we searched for additional angiosperms that do not possess CMT3. Curiously, we identified another Brassicaceae, *Conringia planisiliqua*, which is also missing CMT3. Methylome analysis of *C. planisiliqua* and other closely related Brassicaceae (*Brassica rapa*, *Brassica oleracea*, and *Schrenkiella parvula*), which all possess a CMT3, confirmed the presence of CHG methylation typical of CMT3 activity (Fig. 2A). However, the CHG methylation present in *C. planisiliqua* was similar to that observed in *E. salsugineum* and *cmt3* mutants (Fig. 1D), indicating that the CHG methylation detected is likely a result of RdDM and not maintenance by CMT3. Loci containing gbM were identified using the same methods defined previously (Methods), and CG, CHG, and CHH methylation metagene plots of these defined loci were generated for each of these Brassicaceae (Fig. 2B). All of the species that possess a functional CMT3 also possess gbM, whereas the two species that do not possess CMT3 (*E. salsugineum* and *C. planisiliqua*) do not possess patterns consistent with gbM loci (Fig. 2B). In addition, metagene plots of CG methylation across all genes reveal a complete absence of gbM in *C. planisiliqua* (SI Appendix, Fig. S2), which is similar to observations in *E. salsugineum* (Fig. 1A). Therefore, given the

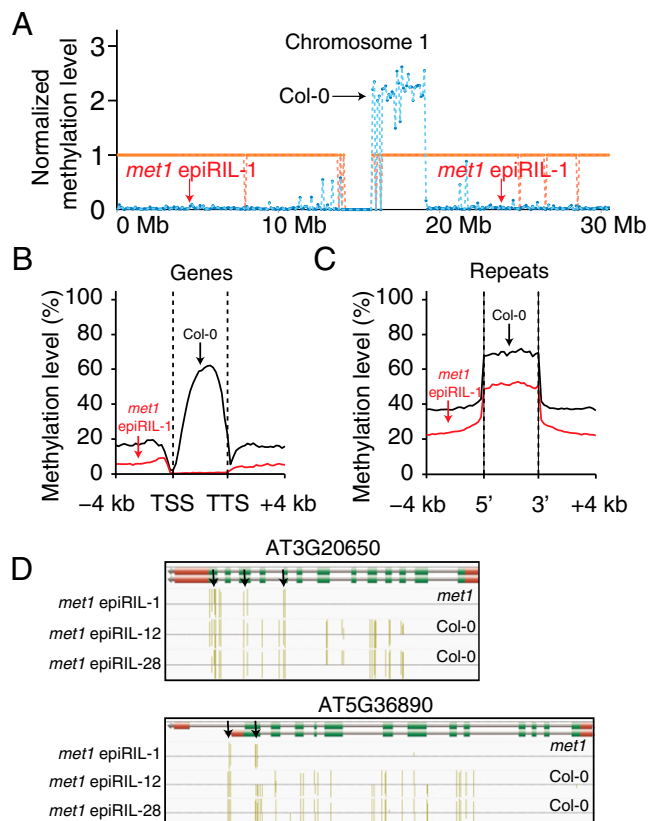


Fig. 3. De novo gbM accumulates incrementally over generational time. (A) A genetic map of *met1* epiRIL line 1 using only the methylation status of gbM loci from wild-type Col-0 as markers. The orange line indicates the expected heterozygous methylation levels from Col-0 and *met1* epiRIL-1. Thus, the blue line indicates inheritance of methylation from Col-0 (>1) or *met1* epiRIL-1 (<1). (B and C) Metaplots of CG methylation in (B) genes and (C) transposons including 4 kb up- and downstream of the TSS and TTS from the Col-0- or the *met1*-derived regions of the epiRIL. (D) Examples of loci in *met1* epiRIL-1 where gbM has partially returned in loci that are located in *met1*-derived regions of the genome. Black arrows indicate where mCG returned.

(Fig. 4B). Measuring enrichment of H2A.Z in two independent *met1* epIRILs revealed a similar result. The distribution and amplitude of H2A.Z were similar at a previously defined set of gbM loci regardless of whether the gene was inherited from the Col-0 or the *met1* parental genome (Fig. 4C). The mechanistic basis for the correlation between H2A.Z and transcription levels is not clear. Regardless, these results showed that the loss of gbM over an evolutionary timescale has no effect on the distribution of H2A.Z.

The loss of gbM in *E. salsguineum* might be compensated by a redistribution of other histone modifications across gene bodies. Therefore, additional ChIP-seq experiments using antibodies against H3K4me3, H3K9me2, H3K27me3, H3K36me3, and H3K56ac were performed to test whether the absence of gbM in *E. salsguineum* affected the distribution of these histone modifications. However, no differences in distributions were observed when compared against *A. thaliana* (Fig. 4 D and E).

We propose that gbM might represent a by-product of errant properties associated with enzymes that can establish DNA methylation, like CMT3, and enzymes that can maintain it, such as MET1. Loss of IBM1, a histone demethylase, results in immediate accumulation of H3K9me2 and CHG methylation in gene bodies (16, 17). In fact, DNA methylome profiling of an *ibm1-6* allele not only confirmed these results but also uncovered an increase of both CG and CHH methylation in gbM loci (*SI Appendix, Fig. S6*). This indicates that in the absence of active removal of H3K9me2 from gene bodies, methylation in all cytosine sequence contexts accumulates. Failure to properly remove H3K9me2 accumulation

from gene bodies of wild-type plants leads to recruitment of CMT3, which in turn methylates cytosines primarily in the CHG context but also enables methylation of CG and CHH sites (*SI Appendix, Fig. S6*). Once methylation is present in the gene bodies it spreads throughout the gene body, as methylated DNA serves as a substrate for the SRA domain-containing proteins KYP/SUVH4/SUVH5, which bind methylated cytosines and lead to continual methylation of H3K9 (28). The lack of gbM at the transcriptional start site (TSS) might be due to an inability of H2A.Z and H3K9me2 to co-occur in nucleosomes, which suggests that the primary role of H2A.Z is to prevent spreading of H3K9me2 into the TSS. This spreading mechanism is also consistent with the loci in *met1* epiRLs, where gbM partially returned, seemingly in a directional manner (Fig. 3D). Therefore, over evolutionary timescales, gbM accumulates and is maintained and tolerated by the genome, as thus far it appears to have no apparent functional role and no deleterious consequences.

It cannot be proven that gbM has no functional role in angiosperm genomes, as it is possible that it has a yet-undiscovered function or that it serves to redundantly perform functions with other transcriptional processes. However, the absence of gbM in *E. salsugineum* and *C. planisiliqua* and their perseverance as species are clear evidence that this feature of the DNA methylome is not required for viability. DNA methylation of gene bodies is also found in mammalian genomes,

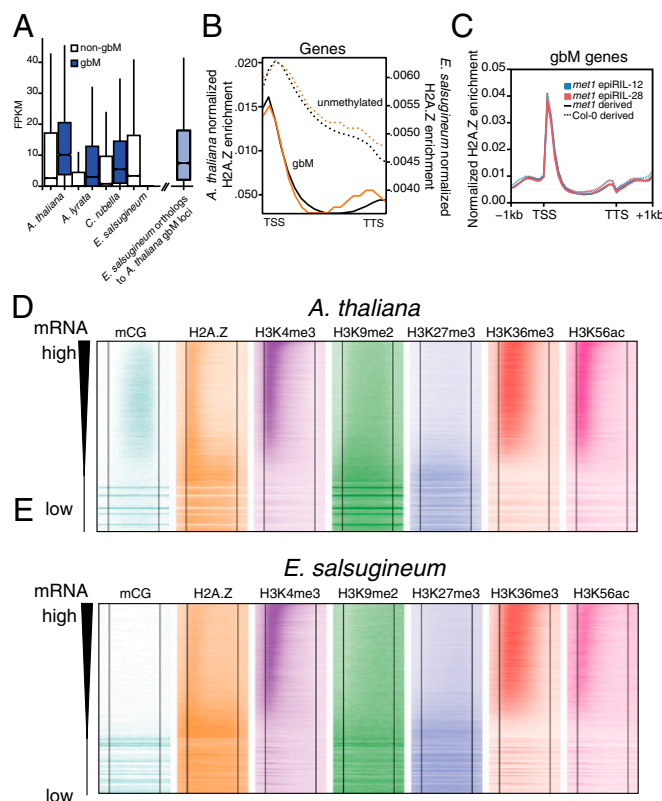


Fig. 4. Gene expression and histone modifications are not affected by the loss of gbM in *E. salsugineum*. (A) Comparison of gene expression levels between gbM and non-gbM within each species listed. Orthologs of gbM loci from *A. thaliana* were used to identify loci from *E. salsugineum*, and FPKM (fragments per kb of transcript per million mapped reads) values were plotted. (B) Metagene plots of H2A.Z enrichment in gbM versus unmethylated genes in *A. thaliana* (black line) and *E. salsugineum* (orange line). The y axis on the left is associated with *A. thaliana*, and the one on the right is associated with *E. salsugineum*. (C) Metagene plots of H2A.Z enrichment in two *met1* epiRILs of gbM loci derived from either the Col-0 (dotted lines) or the *met1* (solid lines) parent. (D and E) Heatmap representation of histone modification distributions and patterns in gene bodies of (D) *A. thaliana* and (E) *E. salsugineum*. The genes in each heatmap are ranked from highest to lowest expression levels. The vertical lines indicate the position of the transcriptional start and stop sites; 1 kb upstream of the TSS and downstream of the TTS are included in the heatmaps.

2,000-bp segments, and 25,000 segments with at least 40 CpG sites were randomly chosen. Segments were ranked by weighted methylation in all contexts. The 2,000 segments with highest methylation and 2,000 segments with lowest methylation were defined as intergenic methylated and intergenic unmethylated, respectively. Orthologs between *A. thaliana* and *E. salsugineum* were split into two groups, gbM and UM, as defined by the methylation in *A. thaliana*. Coordinates for the genes were taken from the TAIR10 annotation of *A. thaliana* and Phytozome 10 annotation of *E. salsugineum* 173 version 1. All intergenic segments and orthologs were broken into 20 equally sized bins. Aligned ChIP reads starting each bin were summed and then normalized by bin length and nonclonal library size.

For heatmaps, the 95th percentile value of all bins was computed, and any bin with a value above this threshold was set equal to the threshold. Finally, the average bin value for bins of intergenic methylated regions was computed. This value was subtracted from all bins in the heatmap, and any bin value less than

zero was set equal to zero. Orthologs were ordered based on mRNA level in *A. thaliana*. For metagene plots, bin values were summed for each ortholog type, gbM and UM, then normalized for the number of genes in each group.

ACKNOWLEDGMENTS. We thank Zachary Lewis, Nathan Springer, and Dave Hall for comments and discussions, as well as Karen Schumaker (*E. salsugineum*), Marcus Koch (*C. planisiliqua*), and Jerzy Paszkowski (*met1* epiRILs) for seeds and the Georgia Genomics Facility and Georgia Advanced Computing Resource Center for technical support. This work was supported by the National Institutes of Health (R00GM100000), Pew Charitable Trusts, and the Office of the Vice President of Research at UGA (R.J.S.). C.E.N. was supported by National Science Foundation (NSF) Postdoctoral Fellowship IOS-1402183. Research in the X.Z. laboratory was supported by the NSF (MCB-0960425). B.T.H. was supported by a Scholars of Excellence graduate fellowship from the University of Georgia (UGA). H.S. was a Howard Hughes Medical Institute (HHMI) Fellow of the Damon Runyon Cancer Research Foundation (DRG-2194-14). S.E.J. is an Investigator of the Howard Hughes Medical Institute.

- Du J, Johnson LM, Jacobsen SE, Patel DJ (2015) DNA methylation pathways and their crosstalk with histone methylation. *Nat Rev Mol Cell Biol* 16(9):519–532.
- Chan SW-L, et al. (2006) RNAi, DRD1, and histone methylation actively target developmentally important non-CG DNA methylation in *Arabidopsis*. *PLoS Genet* 2(6):e83.
- Tran RK, et al. (2005) DNA methylation profiling identifies CG methylation clusters in *Arabidopsis* genes. *Curr Biol* 15(2):154–159.
- Zhang X, et al. (2006) Genome-wide high-resolution mapping and functional analysis of DNA methylation in *Arabidopsis*. *Cell* 126(6):1189–1201.
- Zilberman D, Gehring M, Tran RK, Ballinger T, Henikoff S (2007) Genome-wide analysis of *Arabidopsis thaliana* DNA methylation uncovers an interdependence between methylation and transcription. *Nat Genet* 39(1):61–69.
- Bewick AJ, et al. (2016) The evolution of CHROMOMETHYLASES and gene body DNA methylation in plants. *bioRxiv*, 10.1101/054924.
- Zemach A, McDaniel IE, Silva P, Zilberman D (2010) Genome-wide evolutionary analysis of eukaryotic DNA methylation. *Science* 328(5980):916–919.
- Feng S, et al. (2010) Conservation and divergence of methylation patterning in plants and animals. *Proc Natl Acad Sci USA* 107(19):8689–8694.
- Coleman-Derr D, Zilberman D (2012) Deposition of histone variant H2A.Z within gene bodies regulates responsive genes. *PLoS Genet* 8(10):e1002988.
- Dubin MJ, et al. (2015) DNA methylation in *Arabidopsis* has a genetic basis and shows evidence of local adaptation. *eLife* 4:e05255.
- Schmitz RJ, et al. (2013) Patterns of population epigenomic diversity. *Nature* 495(7440):193–198.
- Coruh C, et al. (2015) Comprehensive annotation of *Physcomitrella patens* small RNA loci reveals that the heterochromatic short interfering RNA pathway is largely conserved in land plants. *Plant Cell* 27(8):2148–2162.
- Du J, et al. (2012) Dual binding of chromomethylase domains to H3K9me2-containing nucleosomes directs DNA methylation in plants. *Cell* 151(1):167–180.
- Jackson JP, Lindroth AM, Cao X, Jacobsen SE (2002) Control of CpNpG DNA methylation by the KRYPTONITE histone H3 methyltransferase. *Nature* 416(6880):556–560.
- Inagaki S, Kakutani T (2012) What triggers differential DNA methylation of genes and TEs: Contribution of body methylation? *Cold Spring Harb Symp Quant Biol* 77:155–160.
- Miura A, et al. (2009) An *Arabidopsis* jmjC domain protein protects transcribed genes from DNA methylation at CHG sites. *EMBO J* 28(8):1078–1086.
- Saze H, Shiraishi A, Miura A, Kakutani T (2008) Control of genic DNA methylation by a jmjC domain-containing protein in *Arabidopsis thaliana*. *Science* 319(5862):462–465.
- Stroud H, Greenberg MVC, Feng S, Bernatavichute YV, Jacobsen SE (2013) Comprehensive analysis of silencing mutants reveals complex regulation of the *Arabidopsis* methylome. *Cell* 152(1–2):352–364.
- Niederhuth CE, et al. (2016) Widespread natural variation of DNA methylation within angiosperms. *bioRxiv*, 10.1101/045880.
- Arias T, Beilstein MA, Tang M, McKain MR, Pires JC (2014) Diversification times among *Brassica* (Brassicaceae) crops suggest hybrid formation after 20 million years of divergence. *Am J Bot* 101(1):86–91.
- Reinders J, et al. (2009) Compromised stability of DNA methylation and transposon immobilization in mosaic *Arabidopsis* epigenomes. *Genes Dev* 23(8):939–950.
- Colomé-Tatché M, et al. (2012) Features of the *Arabidopsis* recombination landscape resulting from the combined loss of sequence variation and DNA methylation. *Proc Natl Acad Sci USA* 109(40):16240–16245.
- Cortijo S, et al. (2014) Mapping the epigenetic basis of complex traits. *Science* 343(6175):1145–1148.
- Zilberman D, Coleman-Derr D, Ballinger T, Henikoff S (2008) Histone H2A.Z and DNA methylation are mutually antagonistic chromatin marks. *Nature* 456(7218):125–129.
- Regulski M, et al. (2013) The maize methylome influences mRNA splice sites and reveals widespread paramutation-like switches guided by small RNA. *Genome Res* 23(10):1651–1662.
- Lin C-H, Workman JL (2011) Suppression of cryptic intragenic transcripts is required for embryonic stem cell self-renewal. *EMBO J* 30(8):1420–1421.
- Xu Y, et al. (2014) *Arabidopsis* MRG domain proteins bridge two histone modifications to elevate expression of flowering genes. *Nucleic Acids Res* 42(17):10960–10974.
- Johnson LM, et al. (2007) The SRA methyl-cytosine-binding domain links DNA and histone methylation. *Curr Biol* 17(4):379–384.
- Baubec T, et al. (2015) Genomic profiling of DNA methyltransferases reveals a role for DNMT3B in genic methylation. *Nature* 520(7546):243–247.
- Kobayashi H, et al. (2012) Contribution of intragenic DNA methylation in mouse gametic DNA methylomes to establish oocyte-specific heritable marks. *PLoS Genet* 8(1):e1002440.
- Goodstein DM, et al. (2012) Phytozome: A comparative platform for green plant genomics. *Nucleic Acids Res* 40(Database issue):D1178–D1186.
- Schultz MD, et al. (2015) Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature* 523(7559):212–216.
- Martin M (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J* 17(1):10–12.
- Langmead B, Trapnell C, Pop M, Salzberg SL (2009) Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10(3):R25.
- Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J R Statist Soc B* 57(1):289–300.
- Schultz MD, Schmitz RJ, Ecker JR (2012) ‘Leveling’ the playing field for analyses of single-base resolution DNA methylomes. *Trends Genet* 28(12):583–585.
- Takuno S, Gaut BS (2012) Body-methylated genes in *Arabidopsis thaliana* are functionally important and evolve slowly. *Mol Biol Evol* 29(1):219–227.
- Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* 30(15):2114–2120.
- Trapnell C, et al. (2010) Transcript assembly and quantification by RNA-seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* 28(5):511–515.
- Trapnell C, et al. (2012) Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* 7(3):562–578.
- Quast C, et al. (2013) The SILVA ribosomal RNA gene database project: Improved data processing and web-based tools. *Nucleic Acids Res* 41(Database issue):D590–D596.
- Kent WJ (2002) BLAT—The BLAST-like alignment tool. *Genome Res* 12(4):656–664.
- Lyons E, Freeling M (2008) How to usefully compare homologous plant genes and chromosomes as DNA sequences. *Plant J* 53(4):661–673.
- Edgar RC (2004) MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32(5):1792–1797.
- Yang Z (2007) PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24(8):1586–1591.
- Li H, et al. (2009) The sequence alignment/map format and SAMtools. *Bioinformatics* 25(16):2078–2079.
- Quinlan AR, Hall IM (2010) BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* 26(6):841–842.

Supporting Information

Plant material. Leaf tissue was flash frozen in liquid nitrogen for all experiments including ChIP-seq, RNA-seq and MethylC-seq. DNA was isolated using a Qiagen Plant DNeasy kit (Qiagen, Valencia, CA) following the manufacturer's recommendations. RNA was isolated using TRIzol (Thermo Scientific, Waltham, MA) following the manufacturer's instructions. Specimens of *Conringia planisiliqua* (B-2011-0093, HEID921022) were cultivated as described below. We collected the 11th leaf from two plants resulting in two biological replicates. The DNA was isolated from the young leaf tissue flash frozen in liquid nitrogen using a Plant DNeasy Mini Kit from Qiagen following the manufacturer's protocol.

Sequencing, assembly and annotation of *C. planisiliqua* genome. Specimens of *Conringia planisiliqua* (B-2011-0093, HEID921022, kindly provided by Marcus Koch, University of Heidelberg/Germany) were grown in the greenhouse (temperature day 20°C/night 18°C; cycles of 16 h light and 8 h darkness, if required day length was extended to 16 h by additional light sources (Osram HQI-BT 400W/D (white), NAV400W (orange/red)) or alternatively plants were shaded from light to keep day length constant. Plants were cultivated on soil (Balster, Type Mini Tray with 1 kg/m³ added fertilizer and 1 kg/m³ Osmocote (Scouts)), watered daily and fertilized weekly with Wuxal Super 8-8-6.

A whole genome shotgun library with a 300 bp insert size was generated following the manufacturer's protocol and paired-end sequenced on an Illumina HiSeq2500 with 101 bp read lengths. We obtained 28,942,646 read-pairs. A genome size estimation based on kmers indicated a genome size of 260 Mb, which is slightly higher than what has been estimated based on flow cytometry (180 Mb). We generated a whole genome assembly using Platanus (v1.2, (1)). The final assembly consisted of 12,970 scaffolds larger than 500 bp with L50 of 37 kb and a N50 of 729 kb summing up to a total length of 121 Mb. For homology-based annotation we aligned the protein sequences of *A. thaliana*, *Arabidopsis lyrata*, *S. parvula*, *E. salsugineum*, *B. rapa*, *Arabis alpina* and *Aethionema arabicum* against the scaffolds using scipio (1.4, (2)). The BLAT alignment files were filtered using a Perl script ("filterPSL.pl") provided with Augustus v3.0 (3) using minCover = 80 and minId = 80 as parameters. The filtered alignment file was then converted into GFF format using another Perl script ("blat2hints.pl") provided with Augustus and was used as input for Augustus *de novo* gene prediction (v3.0.1) using the adapted training parameters for *Arabidopsis* (--species = arabidopsis). Augustus predicted 29,373 genes. In order to check for completeness of our gene set, we blasted the protein sequences against the set of core eukaryotic genes of *A. thaliana* (4) (<http://korflab.ucdavis.edu/datasets/cegma/#SCT2>) and found for 455 out of 458 a blast hit with an e-value <1e-50 and an identity >70%. Reciprocal best BLAST was performed between *C. planisiliqua* protein coding genes and *A. thaliana* transposable elements (TEs) [The Arabidopsis Information Resource (TAIR)10]. Using an e-value cutoff of ≤1e-06 we identified 795 TEs annotated as protein coding genes; these TEs were removed prior to any analyses relating to gbM.

MethylC-seq library construction. All libraries were prepared as described by Urich et al. (5) with the exception of *Conringia planisiliqua*. Genomic DNA (gDNA) was quantified using the Qubit BR assay (ThermoFisher Scientific, U.S.A.) and the quality assessed by standard agarose gel electrophoresis. Bisulfite libraries were constructed using the NEXTflex Bisulfite-Seq Kit (Bioo Scientific) with 500 ng input gDNA being fragmented by COVARIS S2 and then increased in concentration with AMPure Beads (0.8 volume, Beckman Coulter, U.S.A.) and eluted in water. Sheared Lambda DNA was spiked as to record the bisulfite conversion rate. For ligation, the adapter concentration was diluted 1:1. Library fragments were then two times purified with 1 volume AMPure Beads and finally eluted in 10 mM Tris (pH 8.0). Bisulfite conversion of library fragments was performed as outlined in the EZ DNA Methylation-Gold Kit (Zymo Research Corporation, U.S.A.). Converted DNA strands were amplified with 15 PCR cycles, purified with AMPure beads followed by quality assessment with capillary electrophoresis (D 1000 Bioanalyser Assay, Agilent, U.S.A.) and quantified by fluorometry (Qubit HS assay; ThermoFisher Scientific, U.S.A.).

ChIP-seq library construction. ChIP experiments were performed as described in (6). Immunoprecipitated DNA was end repaired using the End-It DNA Repair Kit (Epicentre, Madison, WI) according to the manufacturer's instructions. DNA was purified using Sera-Mag (Thermo Scientific, Waltham, MA) at a 1:1 DNA to beads ratio. The reaction was then incubated for 10 minutes at room temperature, placed on a magnet to immobilize the beads, and the supernatant was removed. The samples were washed two times with 500 µl of 80% ethanol, air-dried at 37°C and then resuspended in 50 µl of 10 mM Tris-Cl (pH 8.0). Finally, the samples were incubated at room temperature for 10 minutes, placed on the magnet, and the supernatant was transferred to a new tube, which contained reagents for "A-tailing." A-tailing reactions were performed at 37°C according to the manufacturer's instructions (New England Biolabs, Ipswich, MA). The samples were cleaned using Sera-Mag beads as previously described. Next, adapter ligation was performed using Illumina Truseq Universal adapters and T4 DNA ligase (New England Biolabs, Ipswich, MA) overnight at 16°C. A double clean-up using the Sera-Mag beads was performed to remove any adapter-adapter dimers and the elution was used for 15 rounds of PCR. Lastly, samples were cleaned up one final time using the procedures described above.

RNA-seq library construction. RNA-seq libraries were constructed using Illumina TruSeq Stranded RNA LT Kit (Illumina, San Diego, CA) following the manufacturer's instructions with limited modifications. The starting quantity of total RNA was adjusted to 1.3 µg, and all volumes were reduced to a third of the described quantity.

Sequencing. MethylC-Seq libraries for the species survey were sequenced using the Illumina HiSeq 2500 (Illumina, San Diego, CA). Sequencing of libraries was performed up to 101 cycles. Wild-type, *met1* epiRILs and *lbm1-6*

methylomes and transcriptomes were sequenced to 150 bp on an Illumina NextSeq500 (Illumina, San Diego, CA) with the exception of wild-type and *sdg7/sdg8/met1* transcriptomes, which were sequenced using the Illumina HiSeq2000 (Illumina, San Diego, CA).

References

1. Kajitani R, et al. (2014) Efficient *de novo* assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res* 24(8):1384–1395.
2. Keller O, Odronitz F, Stanke M, Kollmar M, Waack S (2008) Scipio: Using protein sequences to determine the precise exon/intron structures of genes and their orthologs in closely related species. *BMC Bioinformatics* 9(1):278.
3. Stanke M, et al. (2006) AUGUSTUS: *Ab initio* prediction of alternative transcripts. *Nucleic Acids Res* 34(Web Server issue):W435–W439.
4. Parra G, Bradnam K, Ning Z, Keane T, Korf I (2009) Assessing the gene space in draft genomes. *Nucleic Acids Res* 37(1):289–297.
5. Urich MA, Nery JR, Lister R, Schmitz RJ, Ecker JR (2015) MethylC-seq library preparation for base-resolution whole-genome bisulfite sequencing. *Nat Protoc* 10(3):475–483.
6. Schubert D, et al. (2006) Silencing by plant Polycomb-group genes requires dispersed trimethylation of histone H3 at lysine 27. *EMBO J* 25(19):4638–4649.

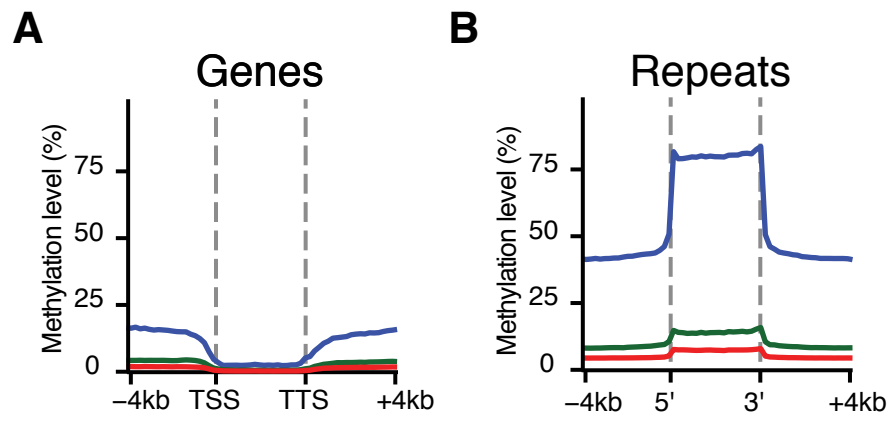


Figure S1. Metagene plots of DNA methylation across **(A)** gene bodies and **(B)** repeats including 4 kb up- and down-stream of the Yukon accession of *E. salisugineum*. (Blue lines = CG methylation, green lines = CHG methylation and red lines = CHH methylation).

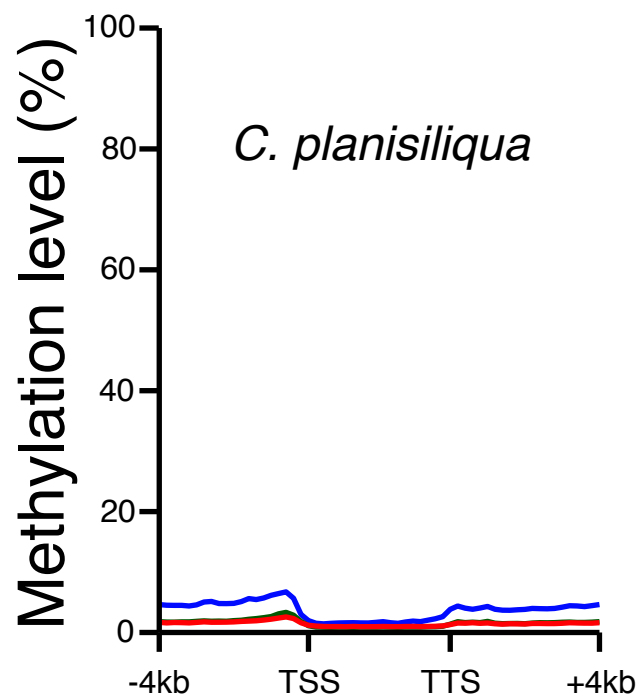


Figure S2. A metagene plot for methylation at CG (blue), CHG (green) and CHH (red) across all genes. Similar to *E. salisugineum* a complete depletion of gbM is observed.

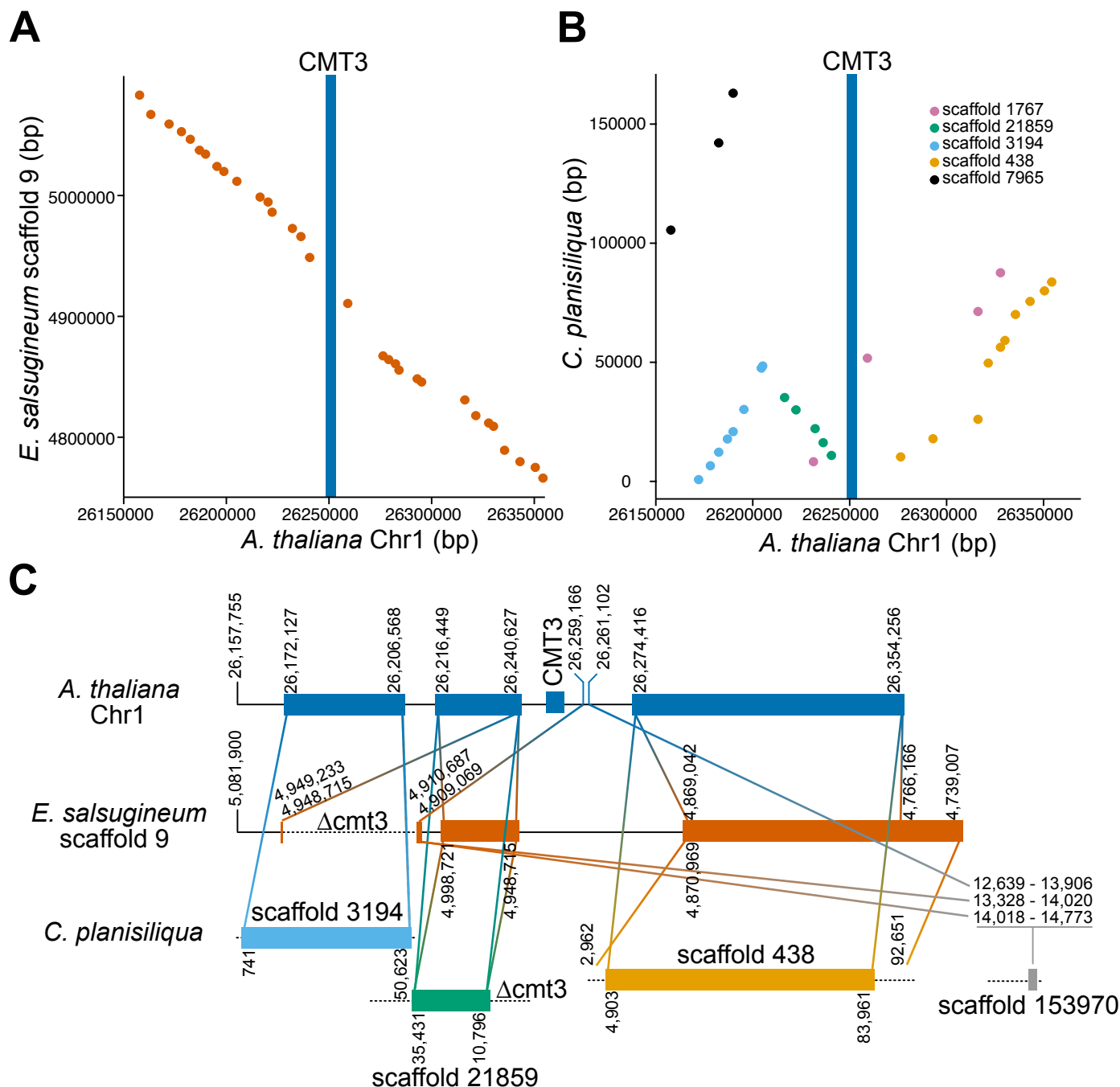


Figure S3. Synteny suggests unique deletions of CMT3 occurred in *E. salsugineum* and *C. planisiliqua*. Macro-synteny between *E. salsugineum* and *C. planisiliqua* to *A. thaliana* is well conserved **(A)** Scaffold 9 in *E. salsugineum* expands a large region of chromosome 1 in *A. thaliana*, which harbors CMT3. **(B)** However, several scaffolds in the *C. planisiliqua* assembly flank CMT3 in *A. thaliana*. **(C)** Micro-synteny immediately flanking CMT3 in *A. thaliana* suggests unique deletions occurred in *E. salsugineum* and *C. planisiliqua*. We hypothesize that an inversion followed by an expansion of ~80 kb occurred immediately upstream of CMT3 in *E. salsugineum*, which is not shared with *C. planisiliqua*. In *C. planisiliqua*, scaffold 21859 is expected to harbor CMT3 at base-pair positions 1 to 3074. However, this region is not syntenic with *A. thaliana* (dashed line) and homology-based searches revealed no CMT3 locus within this region. Scaffold 153970 was revealed to not be syntenic with *A. thaliana* chromosome 1 or *E. salsugineum* scaffold 9. However, BLAST suggests some regions may be homologous. Chromosome and scaffold orientation is based on the sequence assemblies.

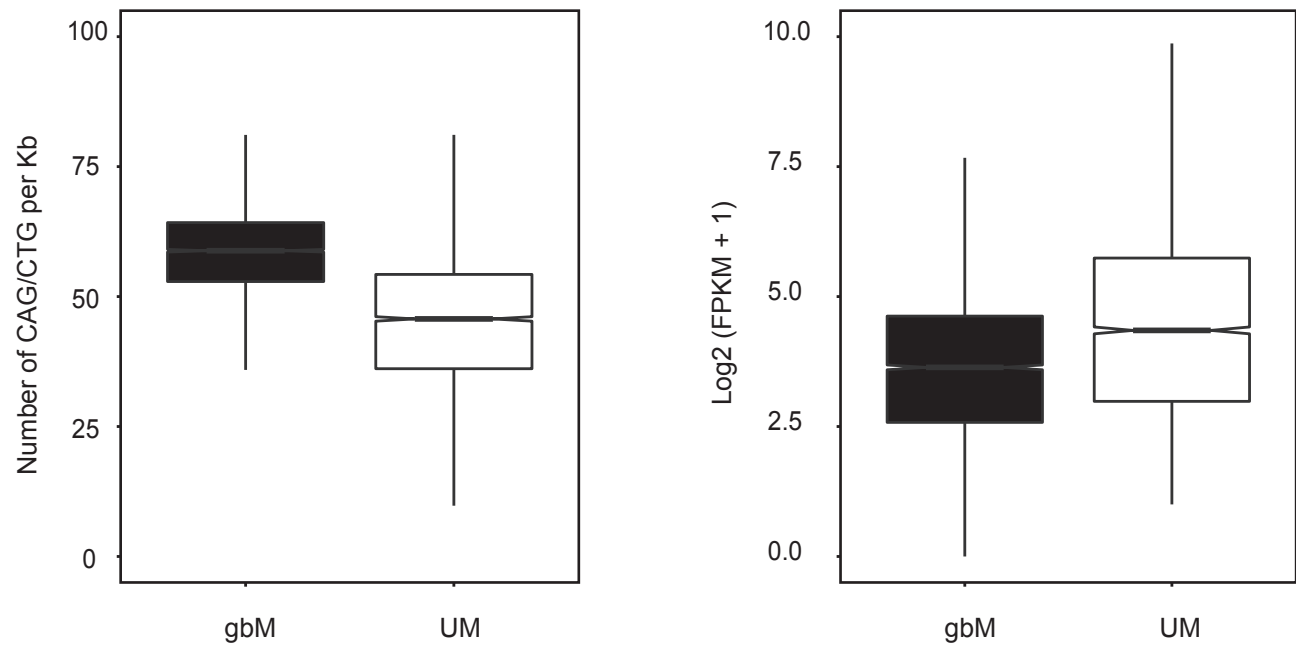


Figure S4. (A) Genes with gbM have higher densities of CAG/CTG sites per kilobase of gene length as compared to UM genes. **(B)** The expression levels of gbM and UM genes used in this the analysis in part **(A)**.

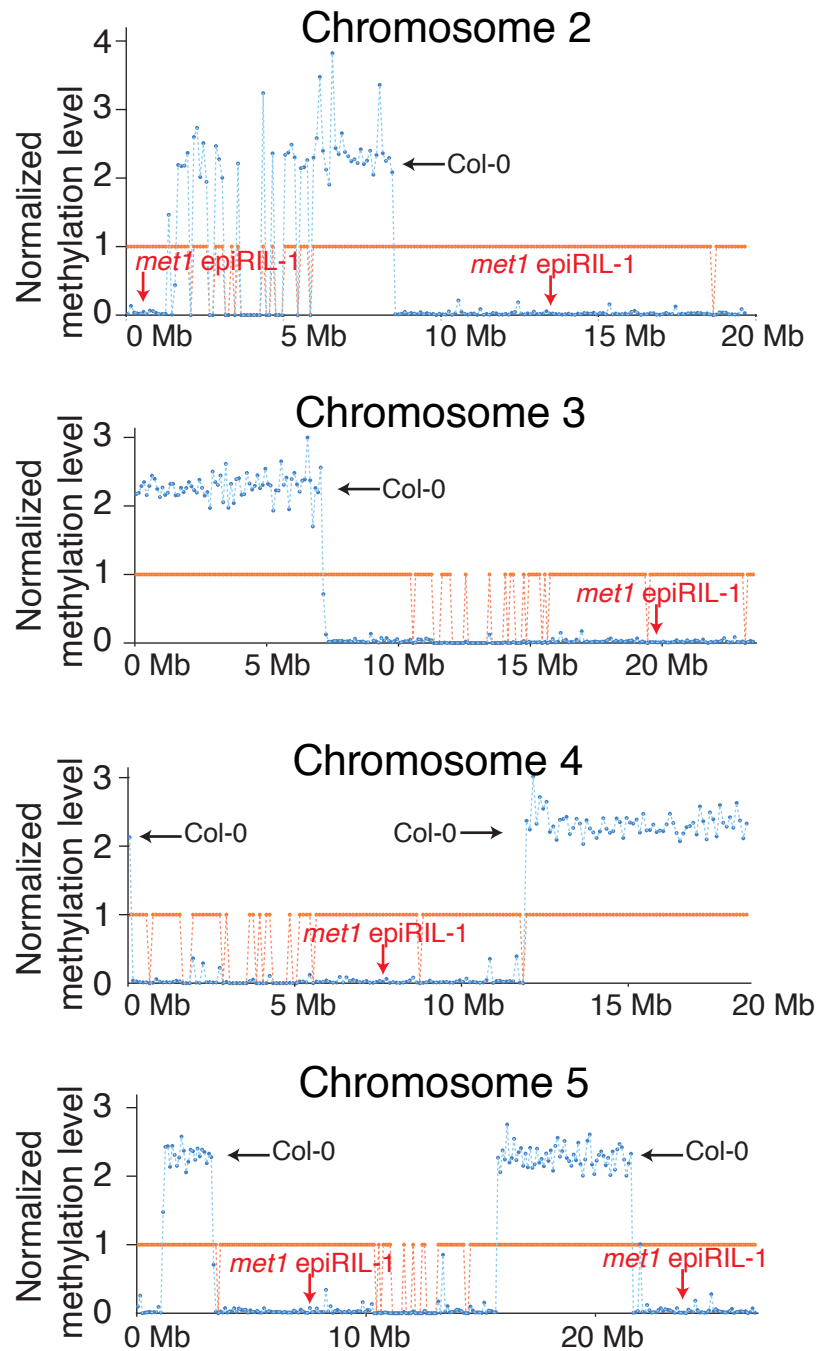


Figure S5. A genetic map of the *met1* epiRIL line 1 using only the methylation status of gbM loci from wild-type Col-0 as markers. The orange midpoint line indicates the expected heterozygous levels and the blue line indicates observed results from the *met1* epiRIL that was analyzed. Chromosomes 2-5 are shown here and chromosome 1 is shown in the main figure.

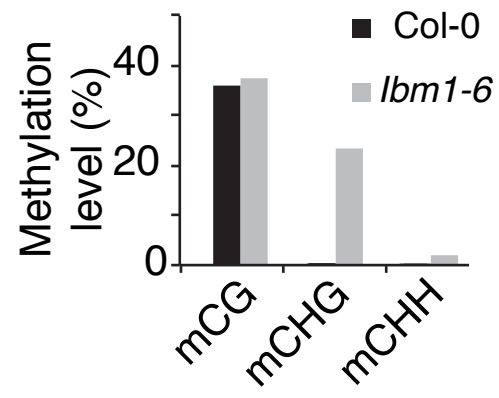


Figure S6. Methylation levels increased in all sequence contexts in gbM loci in *ibm1-6*.

Table S1. Sequencing summary statistics

Methylomes			
Name		Non-clonal unique reads	Non-conversion (%)
<i>A. lyrata</i>		44,841,137	0.32
<i>A. thaliana</i>		47,350,526	0.49
<i>B. oleracea</i>		70,770,309	0.34
<i>B. rapa</i>		43,120,536	0.33
<i>C. planisiliqua</i>		24,293,276	0.70
<i>C. rubella</i>		28,247,233	0.35
<i>E. salsugineum</i> (Shandong)		52,132,901	0.40
<i>E. salsugineum</i> (Yukon)		19,782,414	0.34
<i>met1</i> epiRIL-1		11,635,897	0.22
<i>met1</i> epiRIL-12		12,214,044	0.20
<i>met1</i> epiRIL-28		13,740,036	0.21
Col-0		16,257,718	0.29
<i>met1-3</i>		80,255,893	0.05
<i>ibm1-6</i>		21,399,108	0.17
<i>cmt3-11t</i>		60,299,089	0.61
Transcriptomes			
Comparison	Name	Number of mapped reads	Percent Mapped
Species comparisons	<i>A. thaliana</i> (Same as Col-0 rep 2)	24,156,251	96.6%
	<i>A. lyrata</i>	6,675,951	96.5%
	<i>C. rubella</i>	7,067,624	96.3%
	<i>E. salsugineum</i> (Shandong)	48,362,506	96.6%
	<i>E. salsugineum</i> (Yukon)	13,935,485	89.3%
<i>met1</i> epiRIL comparisons	Col-0 rep1	18,065,032	95.4%
	Col-0 rep2	24,156,251	96.6%
	Col-0 rep3	20,987,335	96.5%
	<i>met1</i> epiRIL-1 rep1	28,378,452	96.5%
	<i>met1</i> epiRIL-1 rep2	25,535,770	97.3%
	<i>met1</i> epiRIL-1 rep3	21,952,024	97.8%
<i>met1/sdg7/sdg8</i> comparisons	WT rep1	24,947,538	99.0%
	WT rep2	23,222,299	99.2%
	WT rep3	24,520,765	99.1%
	<i>met1/sdg7/sdg8</i> rep1	7,313,207	97.9%
	<i>met1/sdg7/sdg8</i> rep2	12,657,371	98.1%
	<i>met1/sdg7/sdg8</i> rep3	15,941,588	98.5%
	<i>met1/sdg7/sdg8</i> rep4	14,687,477	97.9%
	<i>met1/sdg7/sdg8</i> rep5	20,360,538	98.5%

Table S2. Comparison of gene expression between wild type and *met1* epiRIL-1 in *met1* derived regions

Category	Total number of genes	Differentially expressed	Constant
gbM genes	3,471	6	3,465
Unmethylated genes	3,124	46	3,078

P-value = 2.55E-09

Chi-square test indicates gbM loci have significantly less differentially expressed loci in *met1* epiRIL-1 compared to unmethylated loci

Table S3. Comparison of intron retention between wild-type and *met1* epiRIL-1 in *met1* derived regions

Category	Total number of genes	Differentially retained	Constant
gbM genes	3,471	5	3,466
Unmethylated genes	3,124	27	3,097

P-value = 2.64E-05

Chi-square test indicates gbM loci have significantly less numbers of retained intron reads in *met1* epiRIL-1 compared to unmethylated loci

Table S4. Comparison of antisense transcription between wild-type and *met1* epiRIL-1 in *met1* derived regions

Category	Total number of genes	Differentially expressed	Constant
gbM genes	3,471	0	3,471
Unmethylated genes	3,124	0	3,124

Table S5. Comparison of gene expression between wild type and *sdg7/sdg8/met1*

Category	Total number of genes	Differentially expressed	Constant
gbM genes	4,934	386	4,548
Unmethylated genes	4,408	1,113	3,295

P-value = 3.79E-116

Chi-square test indicates gbM loci have significantly less numbers of differentially expressed loci in *sdg7/sdg8/met1* compared to unmethylated loci

Table S6. Comparison of intron retention between wild type and *sdg7/sdg8/met1*

Category	Total number of genes	Differentially retained	No retention
gbM genes	4,934	121	4,813
Unmethylated genes	4,408	499	3,909

P-value = 3.19E-66

Chi-square test indicates gbM loci have significantly less numbers of retained intron reads in *sdg7/sdg8/met1* compared to unmethylated loci

Table S7. Comparison of antisense transcription between wild-type and *sdg7/sdg8/met1*

Category	Total number of genes	Differentially expressed	Constant
gbM genes	4,934	1	4,933
Unmethylated genes	4,408	6	4,402

P-value = 0.04

Chi-square test indicates gbM loci have no significant difference (cutoff=0.01) of antisense transcription in *sdg7/sdg8/met1* compared to unmethylated loci