

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.

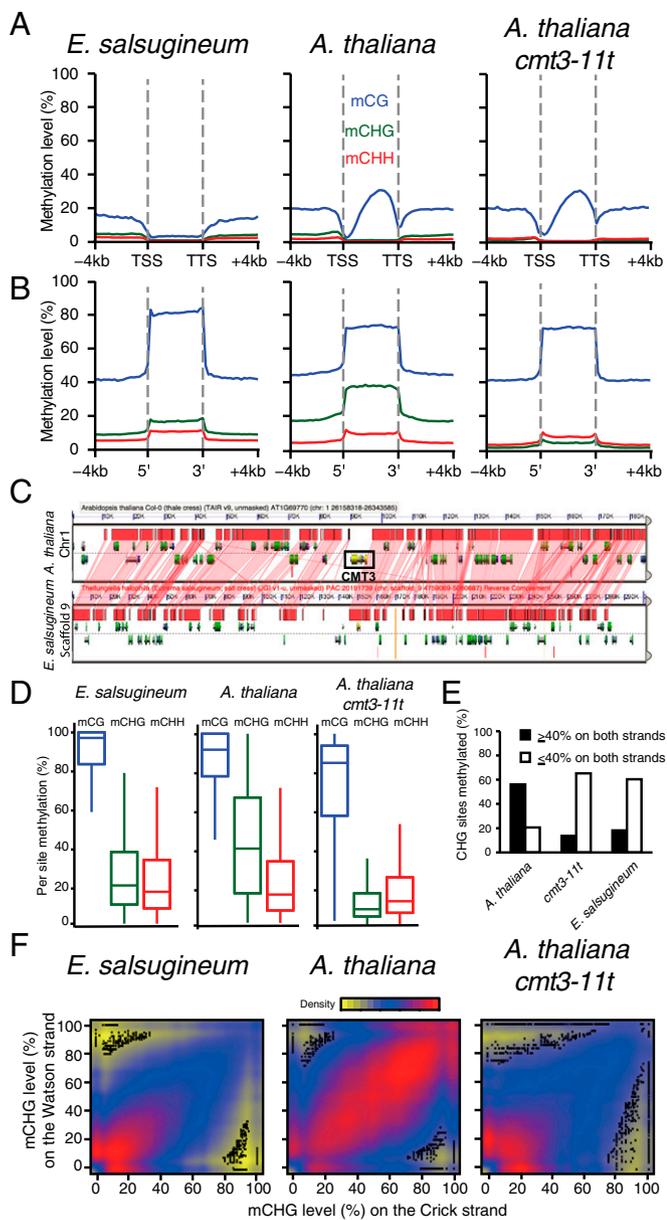


Fig. 1. CMT3 and gbM are absent in *E. salsugineum*. (A and B) Metagene plots of DNA methylation across (A) gene bodies and (B) repeats including 4 kb up- and downstream. TTS, transcriptional termination site. (C) A syntenic block of sequence is shared between *A. thaliana* and *E. salsugineum*. The 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

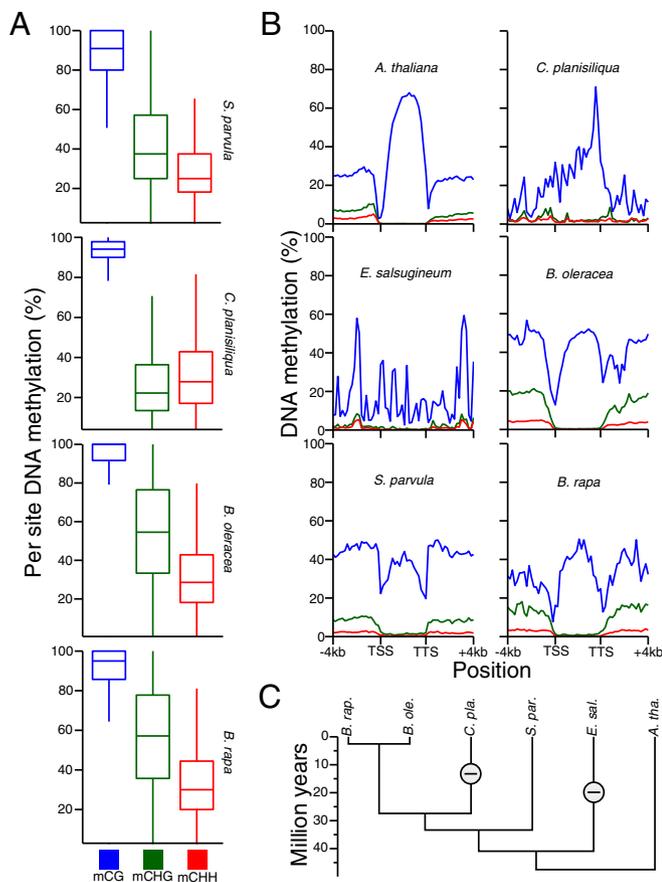


Fig. 2. CMT3 is required for CHG DNA methylation and establishment of gbM in angiosperms. (A) Similar to *E. salsugineum*, loss of CMT3 in *C. planisiliqua* leads to genome-wide reductions of CHG DNA methylation at a per-site level. Closely related species that possess CMT3 maintain higher per-site CHG DNA methylation compared with *C. planisiliqua*. (B) The loss of CMT3 also causes the loss of gene body methylation in *C. planisiliqua*. (C) The loss of CMT3 (-) in *E. salsugineum* and *C. planisiliqua* represents two independent events separated by at least 13.5 My; loss of CMT3 in *C. planisiliqua* and *E. salsugineum* occurred within the last 27.41 and 40.92 My, respectively (20). However, loss in *C. planisiliqua* could be more recent, ≤ 12.27 My (20).

evolutionary relationship between these Brassicaceae species (Fig. 2C), the most parsimonious explanation is two independent losses of CMT3. Additionally, synteny between *A. thaliana*, *C. planisiliqua*, and *E. salsugineum* supports the hypothesis of unique events that led to the deletion of CMT3 in *C. planisiliqua* and *E. salsugineum* (SI Appendix, Fig. S3). Taken together, these results strongly suggest that CMT3 is involved in the establishment of gbM in angiosperms.

In addition to the identification of CMT3 as the enzyme correlated with the establishment of gbM, the methylome data generated here also provided insights into why only a subset of the genes in each plant genome contained gbM. Previous studies have shown that gbM genes tend to be expressed at moderate to high levels in *A. thaliana* (4, 5, 9). Similar results were found in all other plant species that contained gbM (19). These results indicated that the CMT3/KYP pathway might preferentially target moderately to highly transcribed genes. The mechanistic basis for this preference is not yet clear. It is possible that H3K9me2-containing H3 may be occasionally misincorporated into active genes during transcription-coupled nucleosome turnover. Alternatively, nucleosome movement during transcription may lead to the spreading of mCHG/H3K9me2 in flanking regions into active genes, as was described at the *A. thaliana* BONSAI gene (19).

A significant fraction of genes expressed at moderate to high levels did not contain gbM, indicating that active transcription in itself may not be sufficient to trigger gbM. A comparison of the DNA sequence

content in gbM and unmethylated (UM) genes expressed at comparable expression levels (SI Appendix, Fig. S4) revealed a difference in the frequency of CAG/CTG sites per kb of gene length; gbM loci had higher densities of these sites at a frequency of 60.2/kb compared with 45.4/kb for UM genes. These results indicate that gbM loci are potentially predisposed to accumulation of gbM because of their base composition in combination with transcriptional activity.

To obtain further evidence regarding the role of active transcription and sequence composition in the establishment of gbM, we took advantage of the availability of eighth-generation *met1* epigenetic recombinant inbred lines (epiRILs) (21) and determined the characteristics of genes where gbM was reestablished. The epiRILs contain mosaic methylomes, including genomic regions with normal methylation from the wild-type parent and gbM-free regions from the *met1* parent. Because there was very limited genetic variation between the wild-type and *met1* parents, we determined the parental origin of each genomic region according to gbM patterns, similar to what was previously performed for the *ddm1* epiRIL population (Fig. 3A and SI Appendix, Fig. S5) (22, 23). Although mCG was reestablished in repetitive sequences in *met1*-derived regions, the vast majority of genes remained mCG-free, indicating that in most cases gbM has not returned (Fig. 3B and C). A closer inspection identified some rare loci where gbM was partially restored. A total of 50, 10, and 29 genes were identified that accumulated $>5\%$ mCG methylation in *met1* epiRIL-1, -12, and -28 lines, respectively (Fig. 3D). The loci where mCG returned rarely overlapped in the three epiRILs, indicating that the establishment of gbM is a slow and stochastic process. Consistent with the results described earlier, the genes with partially restored gbM tend to be moderately expressed and possess higher frequencies of CAG/CTG sites, indicating that genes with active transcription and higher densities of CAG/CTG sites might be more susceptible to the establishment of gbM.

gbM has been proposed to function in several steps in gene expression, such as suppressing antisense transcription, impeding transcriptional elongation and thus negatively regulating gene expression, or affecting posttranscriptional RNA processing such as splicing (5, 9, 24, 25). However, comparisons between RNA-seq data from *met1* epiRILs and wild-type Col-0 revealed no evidence in support of these possibilities (SI Appendix, Tables S2–S4).

It is possible that the function of gbM has been masked by redundant contributions from other chromatin modification pathways. For example, H3K36me3 has been shown to function in suppressing cryptic transcriptional initiation and regulating splicing (26). We therefore created the triple mutant *met1 sdg7 sdg8*, in which both gbM and H3K36me3 were eliminated (27). RNA-seq experiments from *met1 sdg7 sdg8* and wild-type Col-0 again failed to identify significant differences in mRNA expression, antisense transcription, or splicing variants in a comparison between gbM loci and UM loci (SI Appendix, Tables S5–S7).

The effect of gbM on gene expression might only become detectable over an evolutionary timescale. The identification of plant species with no gbM provided a unique opportunity to test this possibility. We determined the gene set in the *E. salsugineum* genome that likely contained gbM before the loss of CMT3 by projecting orthologous *A. thaliana* gbM genes onto *E. salsugineum*, as they may have at one point possessed gbM in a common ancestor. RNA-seq analysis showed that *E. salsugineum* genes predicted to have contained gbM exhibited similar transcription levels to their *A. thaliana* orthologs (Fig. 4A). These results suggest that gbM may have a limited role in transcriptional regulation.

In addition to gene expression, gbM has also been proposed to function to prevent the histone variant H2A.Z from encroaching into gene bodies (9, 24). To test this hypothesis, chromatin immunoprecipitation sequencing (ChIP-seq) was performed using an H2A.Z antibody on chromatin isolated from *E. salsugineum* and *A. thaliana* leaves. The distribution pattern of H2A.Z in *A. thaliana* was found to be highly consistent with previously published results (9, 24): In gbM genes, H2A.Z was enriched at the 5' ends and absent from regions within genes that contained gbM; in non-gbM genes, H2A.Z could be found distributed throughout the gene bodies (Fig. 4B). Although gbM is presumably absent from *E. salsugineum* for a considerable amount of time, the distribution pattern of H2A.Z remained comparable to that in *A. thaliana*

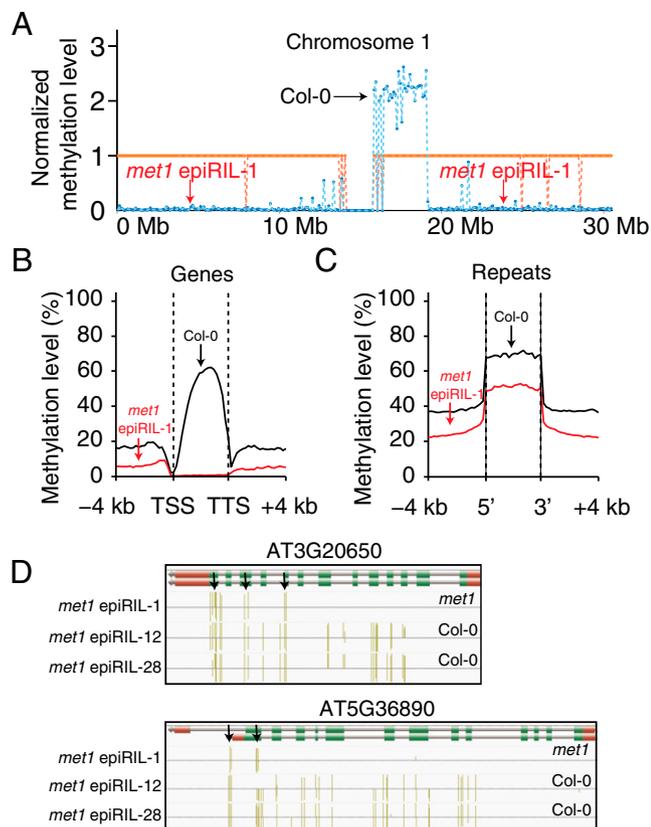


Fig. 3. De novo gbM accumulates incrementally over generational time. (A) A genetic map of *met1* epiRIL line 1 using only the methylation status of gbM loci from wild-type Col-0 as markers. The orange line indicates the expected heterozygous methylation levels from Col-0 and *met1* epiRIL-1. Thus, the 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. salsugineum* 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. salsugineum* affected the distribution of these histone modifications. However, no differences in distributions were observed when compared against *A. thaliana* (Fig. 4D 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* epiRILs, 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,

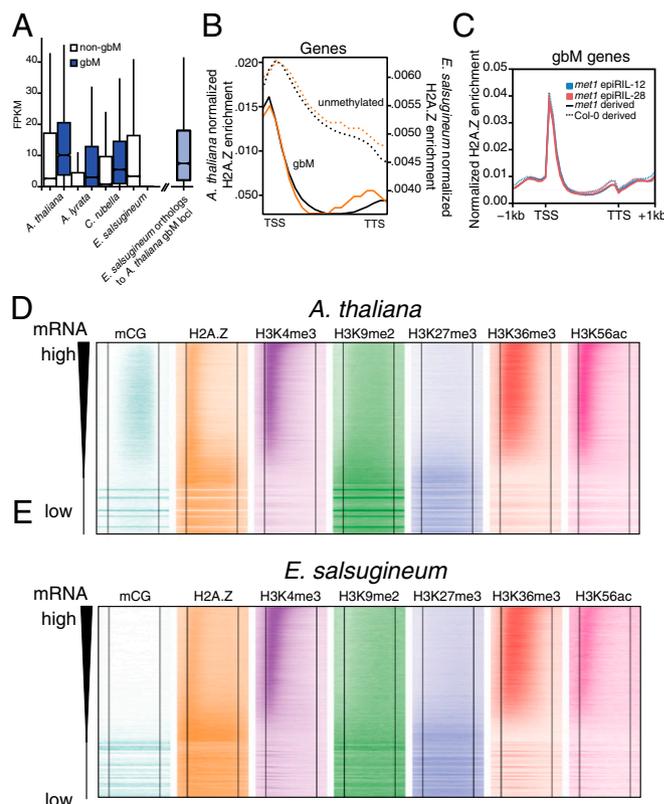


Fig. 4. Gene expression and histone modifications are not affected by the loss of gbM in *E. salsugineum*. (A) Comparison of gene expression levels between gbM and non-gbM within each species listed. Orthologs of gbM loci from *A. thaliana* were used to identify loci from *E. salsugineum*, and FPKM (fragments 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 up-stream of the TSS and downstream of the TTS are included in the heatmaps.

although its distribution throughout transcribed regions and its mechanism for establishment are distinct from angiosperms. In mammals, the methyltransferase Dnmt3B binds H3K36me3 through its PWWP domain and catalyzes de novo methylation in gene bodies, which is then maintained by the maintenance methyltransferase Dnmt1 (29). Intriguingly, loss of methylation in gene bodies in a mouse methylation mutant did not result in global reprogramming of the transcriptome, as the correlation of mRNA levels between mutant and wild type was 0.9611 (30). Only a handful of differentially expressed loci were found. Therefore, although there is abundant DNA methylation of gene bodies in mammals and certain plant species, the evidence available thus far does not provide strong support for a functional role in gene regulation. Future studies of this enigmatic feature of genomes will be required to understand whether gene body methylation has a function or is simply a by-product of active DNA methylation targeting systems.

Methods

MethylC-Seq Analysis. Genomes and annotations were downloaded from Phytozome 10.3 (phytozome.jgi.doe.gov/pz/portal.html) (31). Transposon annotations for *A. thaliana* were downloaded from TAIR (<https://www.arabidopsis.org>). MethylC-seq reads were processed and aligned, and sites were called using published methods (32). The genome was converted into a forward-strand reference (all Cs to Ts) and a reverse-strand reference (all Gs to As). Cutadapt 1.1.0 (33) was used to trim adaptor sequences, and then Bowtie 1.1.0 (34) aligned reads to the two converted reference genomes. Only uniquely aligned reads were retained, and the nonconversion rate was calculated from unmethylated reads aligned to the chloroplast genome or spiked in unmethylated lambda DNA. Summary statistics for each methylome sequenced are presented in *SI Appendix, Table S1*. A binomial test was then applied to each cytosine and corrected for multiple testing using the Benjamini–Hochberg false discovery rate (35). A minimum of three reads is required at each site to call it methylated. For the *cmt3-11t* and *met1-3* methylation data, raw reads were downloaded from a previously published study (18).

Metagene Plots. The gene body was divided into 20 windows. Additionally, regions 4,000 bp upstream and downstream were each divided into 20 200-bp windows. Weighted methylation levels were calculated for each window (36). For gene bodies, only cytosine sites within coding sequences were included. For repeat bodies, all cytosine sites were included. The mean weighted methylation for each window was then calculated for all genes and plotted in R (<https://www.r-project.org/>).

Analysis of mCHG Characteristic of CMT3 Activity. To create the plots for Fig. 1F, both strands of each CHG sequence were required to have a minimum coverage of at least three reads, and at least one of the CHG sites was identified as methylated, as described previously in *MethylC-Seq Analysis*. Methylation levels were calculated for each strand of the symmetric CHG site, and values were used to create a density plot using R. To perform the analysis presented in Fig. 1E, both strands of each CHG dinucleotide were required to have a minimum coverage of at least three reads, and at least one of the CHG sites was identified as methylated.

Identification of CG gbM Genes. Genes were identified as gbM (*Dataset S1*) using a modified version of the binomial test used by Takuno and Gaut (37). This approach tests for enrichment of CG, CHG, and CHH against a background level calculated from the entire set of genes. To do this, the total number of cytosines in each context (CG, CHG, CHH) with mapped reads and the total number of methylated cytosines called were calculated for the coding regions of each gene. Because species differ in genome size, TE content, and other factors, this analysis was restricted to coding sequences, and a single universal background level of methylation was calculated by combining data from all species and determining the percentage of methylated cytosines in each context for the coding regions. A one-tailed binomial test was then applied to each gene for each context testing against the background methylation level. Tested across tens of thousands of genes, there will be some false positives at the extremes. To control for this type I error, *q* values were calculated from *P* values by adjusting for multiple testing using the Benjamini–Hochberg false discovery rate. gbM genes were defined as having reads aligning to at least 20 CG sites, and a CG methylation *q* value <0.05 and CHG and CHH methylation *q* values >0.05. Thus, this test identifies gene-coding sequences that are enriched for mCG and depleted in mCHG and mCHH.

RNA-Seq Data Analysis. Raw FASTQ reads were trimmed for adapters and preprocessed to remove low-quality reads using Trimmomatic version 0.32

(38). Reads were aligned using TopHat version 2.0.13 (39) supplied with a reference GFF file and the following arguments: -l 50000 -b2-very-sensitive -b2-D 50. Transcripts were then quantified using Cufflinks version 2.2.1 supplied with a reference GFF file (40). For the *A. thaliana* samples, rRNA contaminants were removed using an *A. thaliana* rRNA database (41) using BLAT version 35 (42). Differentially expressed genes were determined by Cufflinks, requiring statistically significant changes and also by requiring a twofold (log2) change in gene expression.

Intron Retention Analysis. The *A. thaliana* TAIR10 GFF file was downloaded and intron entries were created. All gene annotation entries were removed except “gene,” “mRNA,” and “intron.” Then, introns were renamed to exons before using TopHat and Cufflinks. Gene expression values and detection of differentially expressed genes were determined using the same process and criteria as described above for the detection of differentially expressed mRNAs.

Antisense Transcription Analysis. Using .bam files generated from RNA-seq alignments as described previously, reads that mapped to convergently transcribed genes were removed from all subsequent analysis. The TAIR10 GFF file was modified to generate an “antisense transcription annotation GFF file” by reversing the strand orientation of all annotated features. Using this file coupled with the filtered .bam files, we determined the prevalence of antisense transcription using the same process and criteria as described above for the detection of differentially expressed mRNAs.

Syntenic Analysis Between *A. thaliana* and *E. salsugineum*. Whole-genome synteny between *A. thaliana* and *E. salsugineum* was determined using CoGe’s (Comparative Genomics) SynFind program (43). *A. thaliana* chromosome 1, which harbors CMT3, and *E. salsugineum* scaffold 9 were identified as being syntenic. Finer synteny analysis was performed on *A. thaliana* chromosome 1 and 150 kb upstream and downstream of CMT3, and *E. salsugineum* scaffold 9, using CoGe’s GEvo program (43).

Identifying Orthologs and Estimating Evolutionary Rates. Reciprocal best BLAST with an *e*-value cutoff of $\leq 1e-08$ was used to identify orthologs. Individual protein pairs were aligned using multiple sequence comparison by log-expectation (MUSCLE) (44) and back-translated into codon alignments using the coding sequence. Insertion–deletion (indel) sites were removed from both sequences, and the remaining sequence fragments were concatenated into a contiguous sequence. A ≥ 30 -bp and ≥ 300 -bp cutoff for retained fragment length after indel removal and concatenated sequence length was implemented, respectively. Subsequently, substitution rates were calculated for sequence pairs between *A. thaliana* and *E. salsugineum* using the yn00 method implemented in phylogenetic analysis by maximum likelihood (PAML) (45).

Analysis of a *met1* epiRIL Methylome and Generation of Epigenetic Maps. Each chromosome was split every 100 kb, and weighted methylation levels were computed from gbM loci in each bin for each *met1* epiRIL and Col-0. Only cytosines in coding regions were used to compute methylation levels. For each bin, the midpoint of the methylation level in Col-0 was defined as the heterozygous methylation level. Methylation levels from each *met1* epiRIL sample were divided by heterozygous methylation levels to normalize each bin. The horizontal line ($Y = 1$) was defined as the heterozygous line, and bins without any gbM loci were not assigned any values. Normalized *met1* epiRIL methylation levels above the heterozygous line were regarded as being derived from the Col-0 parent, whereas data points below the line were regarded as being derived from the *met1* parent. Bins located near cross-over events were excluded from subsequent analyses.

ChIP-Seq Data Analysis. Raw ChIP reads were trimmed for adapters and low-quality bases using Trimmomatic version 0.32. Reads were trimmed for TruSeq version 3 single-end adapters with a maximum of two seed mismatches, palindrome clip threshold of 30, and simple clip threshold of 10. Additionally, leading and trailing bases with quality less than 10 were removed; reads shorter than 50 bp were discarded. Trimmed reads were mapped to the TAIR10 genome using Bowtie2 version 2.2.3 with default options. Mapped reads were sorted using SAMtools version 1.2 (46) and then clonal duplicates were removed using SAMtools version 0.1.9. Remaining reads were converted to BED format with Bedtools version 2.21.1 (47).

Generation of H2A.Z Heatmaps and Metagene Plots. For heatmaps and metagene plots, intergenic regions, defined as the region between genes excluding 100 bp upstream and downstream of genes, were determined for both *A. thaliana* and *E. salsugineum* using the TAIR10 and Phytozome 10 *E. salsugineum* 173 version 1 annotations, respectively. Intergenic regions were broken into

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.