

RESEARCH

Open Access



# The evolution of CHROMOMETHYLASES and gene body DNA methylation in plants

Adam J. Bewick<sup>1</sup>, Chad E. Niederhuth<sup>1</sup>, Lexiang Ji<sup>2</sup>, Nicholas A. Rohr<sup>1</sup>, Patrick T. Griffin<sup>1</sup>, Jim Leebens-Mack<sup>3</sup> and Robert J. Schmitz<sup>1\*</sup>

## Abstract

**Background:** The evolution of gene body methylation (gbM), its origins, and its functional consequences are poorly understood. By pairing the largest collection of transcriptomes (>1000) and methylomes (77) across Viridiplantae, we provide novel insights into the evolution of gbM and its relationship to CHROMOMETHYLASE (CMT) proteins.

**Results:** CMTs are evolutionary conserved DNA methyltransferases in Viridiplantae. Duplication events gave rise to what are now referred to as CMT1, 2 and 3. Independent losses of CMT1, 2, and 3 in eudicots, CMT2 and ZMET in monocots and monocots/commelinids, variation in copy number, and non-neutral evolution suggests overlapping or fluid functional evolution of this gene family. DNA methylation within genes is widespread and is found in all major taxonomic groups of Viridiplantae investigated. Genes enriched with methylated CGs (mCG) were also identified in species sister to angiosperms. The proportion of genes and DNA methylation patterns associated with gbM are restricted to angiosperms with a functional CMT3 or ortholog. However, mCG-enriched genes in the gymnosperm *Pinus taeda* shared some similarities with gbM genes in *Amborella trichopoda*. Additionally, gymnosperms and ferns share a CMT homolog closely related to CMT2 and 3. Hence, the dependency of gbM on a CMT most likely extends to all angiosperms and possibly gymnosperms and ferns.

**Conclusions:** The resulting gene family phylogeny of CMT transcripts from the most diverse sampling of plants to date redefines our understanding of CMT evolution and its evolutionary consequences on DNA methylation. Future, functional tests of homologous and paralogous CMTs will uncover novel roles and consequences to the epigenome.

**Keywords:** CHROMOMETHYLASE, Phylogenetics, DNA methylation, Whole-genome bisulfite sequencing, WGBS

## Background

DNA methylation is an important chromatin modification that protects the genome from selfish genetic elements, is important for proper gene expression, and is involved in genome stability. In plants, DNA methylation is found at cytosines (C) in three sequence contexts: CG, CHG, and CHH (H is any nucleotide, but G). A suite of distinct de novo and maintenance DNA methyltransferases establish and maintain DNA methylation at these three sequence contexts, respectively. CHROMOMETHYLASES (CMTs) are an important class of plant-specific DNA methylation enzymes, which are characterized by the presence of a

CHRromatin Organisation MODifier (CHROMO) domain between the cytosine methyltransferase catalytic motifs I and IV [1]. Identification, expression, and functional characterization of CMTs have been extensively performed in the model plant *Arabidopsis thaliana* [2–4] and in the model grass species *Zea mays* [5–7].

There are three CMT genes encoded in the *A. thaliana* genome: CMT1, CMT2, and CMT3 [2, 8–10]. CMT1 is the least studied of the three as a handful of *A. thaliana* accessions contain an Evelknievel retroelement insertion or a frameshift mutation truncating the protein, which has been suggested that CMT1 is non-essential [8]. The majority of DNA methylation at CHH sites (mCHH) within long transposable elements in pericentromeric regions of the genome is targeted by a CMT2-dependent pathway [3, 4]. Allelic variation at CMT2

\* Correspondence: schmitz@uga.edu

<sup>1</sup>Department of Genetics, University of Georgia, Athens, GA 30602, USA  
Full list of author information is available at the end of the article



has been shown to alter genome-wide levels of CHH DNA methylation (mCHH) and alleles of CMT2 may play a role in adaptation to temperature [11–13]. In contrast, DNA methylation at CHG (mCHG) sites is often maintained by CMT3 through a reinforcing loop with histone H3 lysine 9 di-methylation (H3K9me2) catalyzed by the KRYPTONITE (KYP)/SU(VAR)3-9 HOMOLOG 4 (SUVH4), SUVH5, and SUVH6 lysine methyltransferases [2, 6, 14, 15]. In *Z. mays*, ZMET2 is a functional homolog of CMT3 and catalyzes the maintenance of mCHG [5–7]. A paralog of ZMET2, ZMET5, contributes to the maintenance of mCHG to a lesser degree in *Z. mays* [5, 7]. Homologous CMTs have been identified in other flowering plants (angiosperms) [16–19], early diverging land plants (Embryophyta) – the moss *Physcomitrella patens* and the lycophyte *Selaginella moellendorffii* – and the green algae (Chlorophyta) *Chlorella sp.* NC64A and *Volvox carteri* [20]. The function of CMTs in non-angiosperms is poorly understood. However, in at least *P. patens* a CMT protein contributes to mCHG [21].

Within angiosperms (flowering plants), a large number of genes in angiosperms exclusively contain CG DNA methylation (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” or “gbM”) [22–25]. GbM genes are generally constitutively expressed, evolutionarily conserved, and typically longer than unmethylated genes [25–27]. How gbM is established and subsequently maintained is unclear. Recently it was discovered that CMT3 has been independently lost in two angiosperm species belonging to the Brassicaceae family of plants – *Conringia planisiliqua* and *Eutrema salsugineum* – and this coincides with the loss of gbM [19, 25], indicating CMT3 is required for maintenance of DNA methylation. This has led to a hypothesis that the evolution of gbM is linked to incorporation/methylation of histone H3 lysine-9 di-methylation (H3K9me2) in gene bodies with subsequent failure of INCREASED IN BONSAI METHYLATION 1 (IBM1) to demethylate H3K9me2 [19, 28]. This provides a substrate for CMT3 to bind and methylate DNA and through an unknown mechanism leads to mCG. MCG is maintained over evolutionary timescales by the CMT3-dependent mechanism and during DNA replication by the maintenance DNA METHYTRANSFERASE 1 (MET1). Methylated DNA then provides a substrate for binding by KRYPTONITE (KYP) and related family members through their SRA domains, which increases the rate at which H3K9 is di-methylated [29]. Finally, mCG spreads throughout the gene over evolutionary time [19].

Previous phylogenetic studies have proposed that CMT1 and CMT3 are more closely related to each other than to CMT2, and that ZMET2 and ZMET5 proteins are more closely related to CMT3 than to CMT1 or CMT2 [5], and

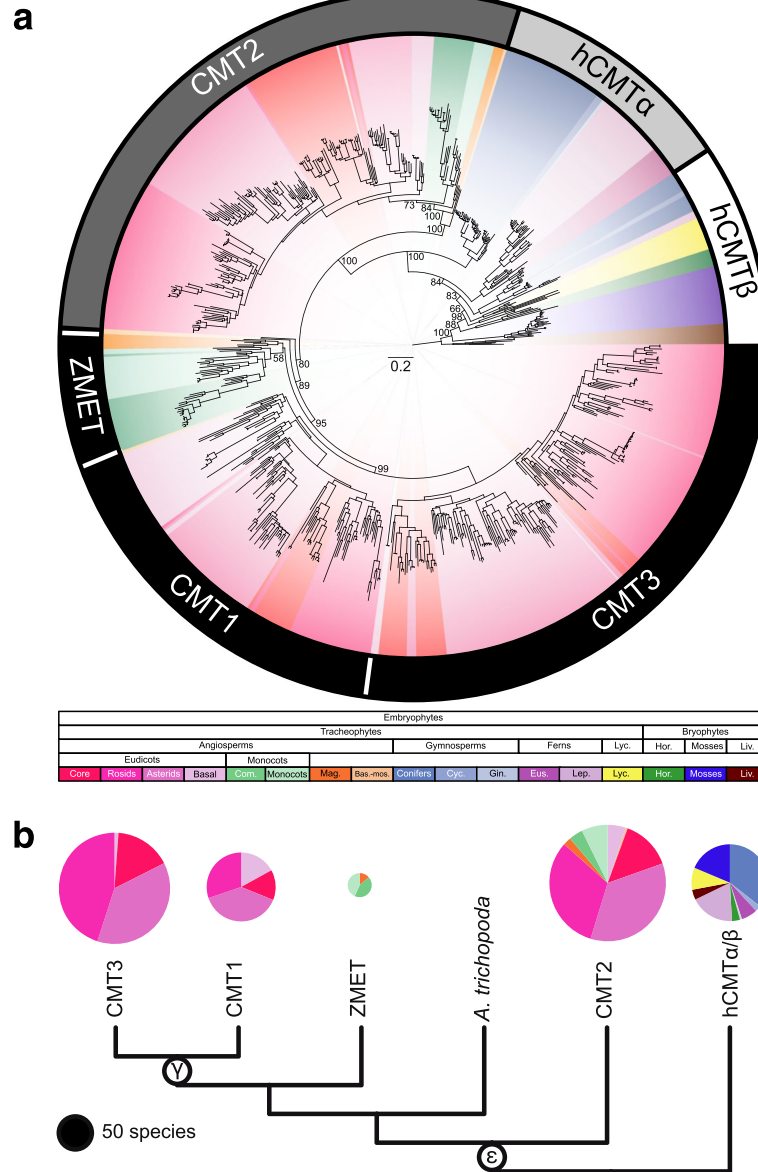
the placement of non-seed plant CMTs is more closely related to CMT3 [21]. However, these studies were not focused on resolving phylogenetic relationships within the CMT gene family, but rather relationships of CMTs between a handful of species. These studies have without question laid the groundwork to understand CMT-dependent DNA methylation pathways and patterns in plants. However, the massive increase in transcriptome data from a broad sampling of plant species together with advancements in sequence alignment and phylogenetic inference algorithms have made it possible to incorporate thousands of sequences into a single phylogeny, allowing for a more complete understanding of the CMT gene family. Understanding the evolutionary relationships of CMT proteins is foundational for inferring the evolutionary origins, maintenance, and consequences of genome-wide DNA methylation and gbM.

Here we investigate phylogenetic relationships of CMTs at a much larger evolutionary timescale using data generated from the 1KP Consortium ([www.onekp.com](http://www.onekp.com)). In the present study, we have identified and analyzed 771 messenger RNA transcripts from de novo assembled transcriptomes and primary coding sequences from annotated genomes belonging to the CMT gene family. A final set of sequences was obtained from an extensive taxonomic sampling of 443 different species including eudicots (basal, core, rosid, and asterid), basal angiosperms, monocots and monocots/commelinid, magnoliids, gymnosperms (conifers, Cycadales, Ginkgoales), monilophytes (ferns and fern allies), lycophytes, bryophytes (mosses, liverworts and hornworts), and green algae. CMT homologs identified across Viridiplantae (land plants and green algae) indicate that CMT genes originated prior to the origin of Embryophyta (land plants) ( $\geq 480$  million years ago [MYA]) [30–33]. In addition, phylogenetic relationships suggest at least two duplication events occurred within the angiosperm lineage giving rise to the CMT1, CMT2, and CMT3 clades. In the light of CMT evolution we explored patterns of genomic and genic DNA methylation levels in 77 species of Viridiplantae, revealing diversity of the epigenome within and between major taxonomic groups and the evolution of gbM in association with the origin of CMT3 and orthologous sequences.

## Results

### The origins of CHROMOMETHYLASES

CMTs are found in most major taxonomic groups of land plants and some algae: eudicots, basal angiosperms, monocots and commelinids, magnoliids, gymnosperms, ferns, lycophytes, mosses, liverworts, hornworts, and green algae (Fig. 1a, Additional file 1: Table S1). CMTs were not identified in transcriptome datasets for species sister to Viridiplantae, including those belonging to Glaucophyta, red algae, Dinophyceae, Chromista, and Euglenozoa. CMTs



**Fig. 1** Phylogenetic relationships of CMTs across Embryophyta. **a** CMTs are separated into four monophyletic clades based on bootstrap support and the relationship of *A. thaliana* CMTs: (1) the gbM-dependent CMT superclade with subclades CMT1, CMT3, ZMET, and *A. trichopoda*; (2) CMT2; and (3) homologous (hCMT) α and β. CMT1 and CMT3 clades only contain eudicot species of plants suggesting a eudicot-specific duplication event that occurred after the divergence of eudicots from monocots and monocots/commelinids. Sister to CMT1 and CMT3 is the monophyletic group ZMET, which contains monocots, monocots/commelinids, and magnoliids. CMT2 is sister to CMT1 and CMT3. Lastly, the polyphyletic hCMT clades are sister to all previously mentioned clades. hCMTα is sister to CMT2 and the CMT superclade and contains gymnosperm and ferns. hCMTβ contains gymnosperms, ferns, and other early diverging land plants. **b** A collapsed CMT gene family tree showing the seven clades described in (a). Pie charts represent species diversity within each clade and are scaled to the number of species. Two duplication events shared by all angiosperms (ε) and eudicots (γ) gave rise to what is now referred to as CMT1, CMT2, and CMT3. These duplication events correspond to what was reported by Jiao et al. (2011). Values at nodes in (a) and (b) represent bootstrap support from 1000 replicates and (a) was rooted to the clade containing all liverwort species

were identified in a few green algae species: *Picocystis salinarum*, *Cylindrocystis sp.*, and *Cylindrocystis brebissonii*. Additionally, functional CMTs – based on the presence/absence of characterized protein domains – were not identified from three species within the gymnosperm

order Gnetales. A transcript with a CHROMO and C-5 cytosine-specific DNA methylase domain was identified in *Welwitschia mirabilis* (Gnetales), but this transcript did not include a Bromo Adjacent Homology (BAH) domain. The BAH domain is an interaction surface that is required

to capture H3K9me2 and mutations that abolish this interaction causes a failure of a CMT protein (i.e. ZMET2) binding to nucleosomes and a complete loss of activity in vivo [6]. Therefore, although a partial sequence is present, it might represent a non-functional allele of a CMT. Alternatively, it might represent an incomplete transcript generated during sequencing and assembling of the transcriptome. Overall, the presence of CMT homologs across Viridiplantae and their absence from sister taxonomic groups suggest CMT evolved following the divergence of green algae [34, 35].

The relationships among CMTs suggest that CMT2 and the clade containing CMT1, CMT3, and ZMET arose from a duplication event at the base of all angiosperms (Fig. 1b). This duplication event might have coincided with the ancestral angiosperm whole-genome duplication (WGD) event,  $\epsilon$  [36]. In support of this hypothesis, collinearity in face of  $\sim 192$  million years of evolution is observed between CMT2 and CMT3/ZMET in a number of eudicot and monocot/commelinid species (Additional file 2: Figure S1) [36]. Genealogical relationships of CMTs largely recapitulate species relationships (Fig. 1a) [34, 37]. However, CMTs in gymnosperms and ferns are paraphyletic (Fig. 1a). Similarly, these homologous sequences might have been derived from a WGD (i.e.  $\zeta$ , the ancestral seed plant WGD), with one paralog being the ancestor to CMT1, CMT2, and CMT3 and ZMET [36]. Previously identified CMTs in *S. moellendorffii* [20] and *P. patens* [20, 38] were identified, which are sister to clades containing CMT1, CMT2, and CMT3 and ZMET (Fig. 1a). Additionally, CMTs previously identified in the green algae *Chlamydomonas reinhardtii*, *Chlorella sp.* NC64A, and *Volvox carteri* were excluded from phylogenetic analysis because they lacked the CHROMO and other domains typically associated with CMT proteins (Additional file 2: Figure S2). Furthermore, based on percent amino acid identity, *C. reinhardtii* and *V. carteri* CMT sequences are more homologous to MET1 than a CMT (Additional file 3: Table S2). Similar to *S. moellendorffii* and *P. patens* CMT sequences, green algae CMT sequences are sister to clades containing CMT1, CMT2, and CMT3 and ZMET (Additional file 2: Figure S3). The increase taxonomic sampling redefines relationships of CMTs in early diverged land plants and in Viridiplantae in general [18, 20, 21, 38, 39].

Further diversification of CMT proteins occurred in eudicots. CMT1 and CMT3 clades contain only sequences from eudicots (Fig. 1a and b). This relationship supports the hypothesis that CMT1 and CMT3 arose from a duplication event shared by all eudicots. Thus, CMT1 and CMT3 might be the result of the  $\gamma$  WGD event at the base of eudicots [36]. Collinearity between CMT1 and CMT3, despite  $\sim 125$  million years of divergence, further supports this hypothesis (Additional file 2:

Figure S4a). Not all eudicots possess CMT1, CMT2, and CMT3, but rather exhibit CMT gene content in the range of 0–3 (Additional file 2: Figure S5a). Also, many species possess multiple copies of CMT1, CMT2, or CMT3. The presence/absence of CMTs might represent differences in transcriptome sequencing coverage or spatial and temporal divergence of expression. However, CMT2, CMT3, and homologous proteins have functions in methylating a significant number of non-CG sites throughout the entire genome and thus are broadly expressed in *A. thaliana*, *Z. mays*, and other species [6, 8, 18, 25]. Additionally, eudicot species with sequenced and assembled genomes show variation in the presence/absence and copy number of CMTs. Hence the type of tissue(s) used in transcriptome sequencing (www.onekp.com) would have limited biases against CMTs, suggesting that the variation reflects presence/absence at a genetic level.

The *Z. mays* in-paralogs ZMET2 and ZMET5, and closely related CMTs in other monocots, commelinids, and magnoliids form a well-supported monophyletic clade (Fig. 1a, Additional file 2: Figure S4b and c). In addition to *Z. mays*, in-paralogs were identified in *Sorghum bicolor* and *Brachypodium distachyon* (Fig. 1a, Additional file 2: Figure S4b). Relationships of *S. bicolor* and *Z. mays* ZMETs differed between gene and amino acid derived phylogenies (Fig. 1a, Additional file 2: Figure S4b). However, collinearity between paralogs of both species supports two independent duplications (Additional file 2: Figure S4c). Also, paralogous ZMETs are shared across species (Additional file 2: Figure S4b). These shared paralogs might have originated from a Poaceae-specific duplication event, which was followed by losses in some species. The contribution of each paralog to DNA methylation and other chromatin modifications remains unknown at this time.

Akin to eudicots, monocots and monocots/commelinids possess ZMET and CMT2 homologs (Additional file 2: Figure S5b). The model grass species *Z. mays* has lost CMT2, whereas the closely related species *S. bicolor* possess both ZMET and CMT2 (Additional file 1: Table S1). ZMET is not strictly homologous to CMT3 and represents a unique monophyletic group that is sister to both CMT1 and CMT3. However, ZMET2 is functionally similar to CMT3 in maintaining DNA methylation at CHG sites [6, 7]. Unlike CMT3, ZMET2 is associated with DNA methylation at CHH sites within some loci [7]. Given the inclusion of monocot and magnoliid species in the monophyletic ZMET clade, this dual-function is expected to be present in other monocot species and magnoliid species.

Overall, these redefined CMT clades, characterized across the land plant phylogeny, are well supported (Fig. 1a). Thus, the identification of novel CMTs in magnoliids, gymnosperms, lycophytes, hornworts, liverworts,

bryophytes, and green algae pushes the timing of evolution of CMT and potentially certain mechanisms maintaining mCHG and mCHH, prior to the origin of Embryophyta ( $\geq 480$  MYA) [30–33].

### Reduced selective constraint of CMT3 in the Brassicaceae affects gbM

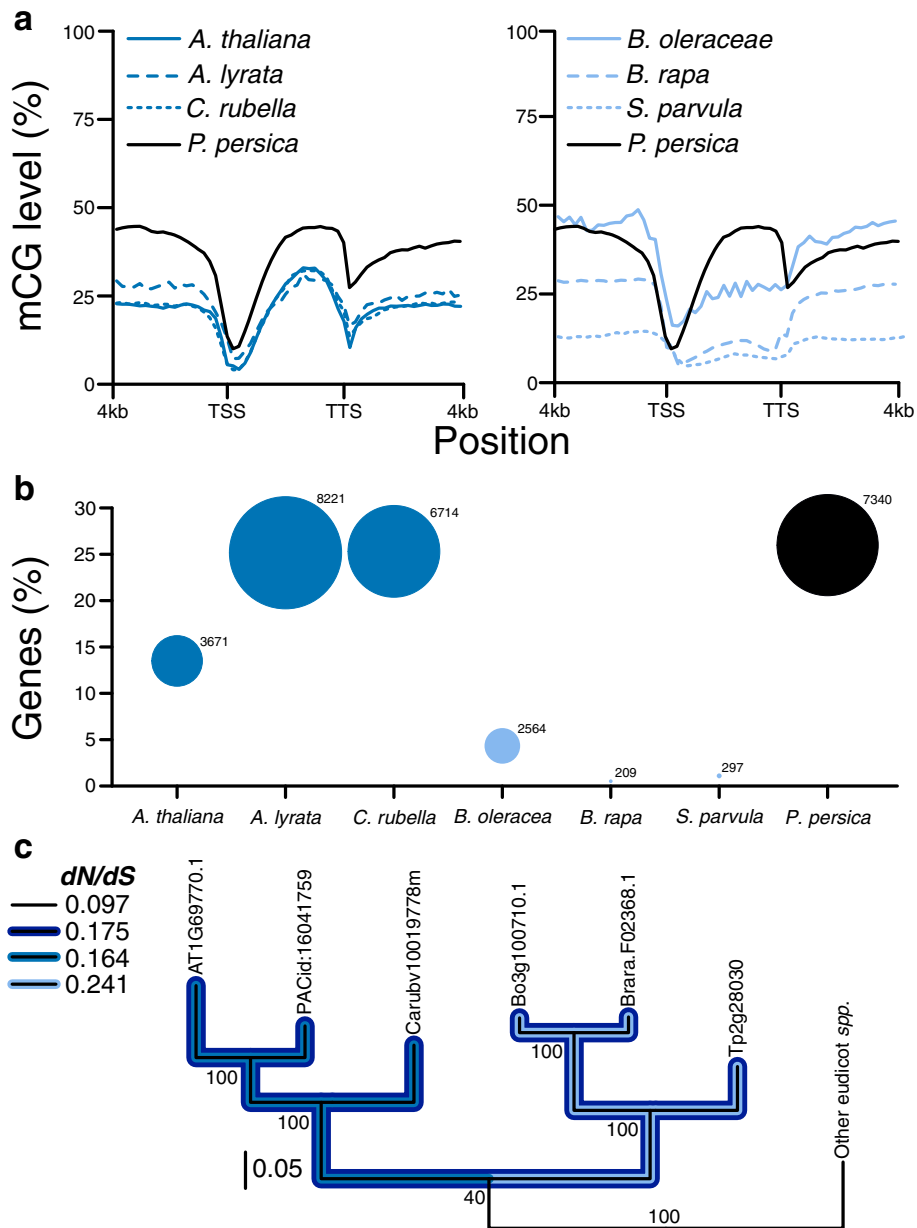
Recent work has described the DNA methylomes of 34 angiosperms, revealing extensive variation across this group of plants [25]. This variation was characterized in terms of overall genomic and genic levels of DNA methylation, genome-wide per-site levels of DNA methylation, and the number of DNA methylated genes. Overall and per-site levels of DNA methylation homogenizes the variation within a heterogeneous population of cells and cell types. Although ascribing genes as DNA methylated relies on levels of DNA methylation, this metric provides insights into the predominant DNA methylation pathway and expected relationship to genic characteristics [25–27]. The genetic underpinnings of this variation are not well understood, but some light has been shed through investigating DNA methylation within the Brassicaceae [19]. The Brassicaceae have reduced levels of genomic and genic levels of mCG, genome-wide per-site levels of mCHG, and numbers of gbM genes [19, 25]. Genome-wide per-site levels of mCG are not reduced compared with other angiosperms [19, 25]. In at least *C. planisiliqua* and *E. salsugineum* this reduction in levels of DNA methylation and numbers of gbM genes has been attributed to the loss of the CMT3 [19]. However, closely related species with CMT3 – *Brassica oleracea*, *Brassica rapa*, and *Schrenkiella parvula* – have reduced levels of mCG within gbM and numbers of gbM genes compared with the sister clade of *A. thaliana*, *Arabidopsis lyrata*, and *Capsella rubella* (Fig. 2a and b) and overall with other eudicots [19, 25]. Although CMT3 is present, changes at the sequence level, including the evolution of deleterious or functionally null alleles, could disrupt function to varying degrees. At the sequence level, CMT3 has a higher rate of molecular evolution – measured as  $dN/dS$  ( $\omega$ ) – in the Brassicaceae ( $\omega = 0.175$ ) compared with 162 eudicots ( $\omega = 0.097$ ), with further increases in the clade containing *B. oleracea*, *B. rapa*, and *S. parvula* ( $\omega = 0.241$ ) compared with the clade containing *A. thaliana*, *A. lyrata*, and *C. rubella* ( $\omega = 0.164$ ) (Fig. 2). A low background rate of molecular evolution suggests purifying selection acting to maintain low allelic variation across eudicots. Conversely, increased rates of molecular evolution can be a consequence of positive selection. However, a hypothesis of positive selection was not preferred to contribute to the increased rates of  $\omega$  in either Brassicaceae clade (Additional file 4: Table S3). Alternatively, relaxed selective constraint possibly resulted in an increased  $\omega$ , which might have introduced null alleles ultimately affecting function of CMT3, and in turn,

affecting levels of DNA methylation and numbers of gbM genes. Higher rates of molecular evolution in the clade containing *Brassica spp.* and *S. parvula* relative to all eudicots and other Brassicaceae are correlated with an exacerbated reduction in the numbers of gbM loci and their methylation levels, which suggests unique substitutions between clades or a quantitative effect to an increase in the number of substitutions. However, at least some substitutions affecting function are shared between Brassicaceae clades because both have reductions in per-site levels of mCHG [19].

### Divergence of DNA methylation patterns within gene bodies during Viridiplantae evolution

Levels and distributions of DNA methylation within gene bodies are variable across Viridiplantae. Levels of mCG range from  $\sim 2\%$  in *S. moellendorffi* to  $\sim 86\%$  in *Chlorella sp.* NC64A (Fig. 3a). Other plant species fall between these two extremes (Fig. 4a) [25]. *Beta vulgaris* remains distinct among angiosperms and Viridiplantae with respect to levels of DNA methylation at all sequence contexts (Fig. 3a). Similarly, *Z. mays* is distinct among monocots (Fig. 3a). Gymnosperms and ferns possess similar levels of mCG to mCHG within gene bodies and levels of mCHG qualitatively parallel those of mCG similar to observations in recently published study (Fig. 3a, Additional file 2: Figure S6) [40]. A similar pattern is observed in *Z. mays*. However, this pattern is not shared by other monocots/commelinids included in our study [25]. High levels of mCHG are common across the gymnosperms and ferns investigated in this study and tend to be higher than levels observed in angiosperms (Fig. 3a, Additional file 2: Figure S6) [25]. DNA methylation at CG, CHG, and CHH sites within gene bodies was detected in the liverwort *Marchantia polymorpha* (Fig. 3a). Furthermore, DNA methylation at CG sites was not detected in the *P. patens* when all genes are considered (Fig. 3a). Overall, increased taxonomic sampling has revealed natural variation in DNA methylation patterns between and within groups of Viridiplantae.

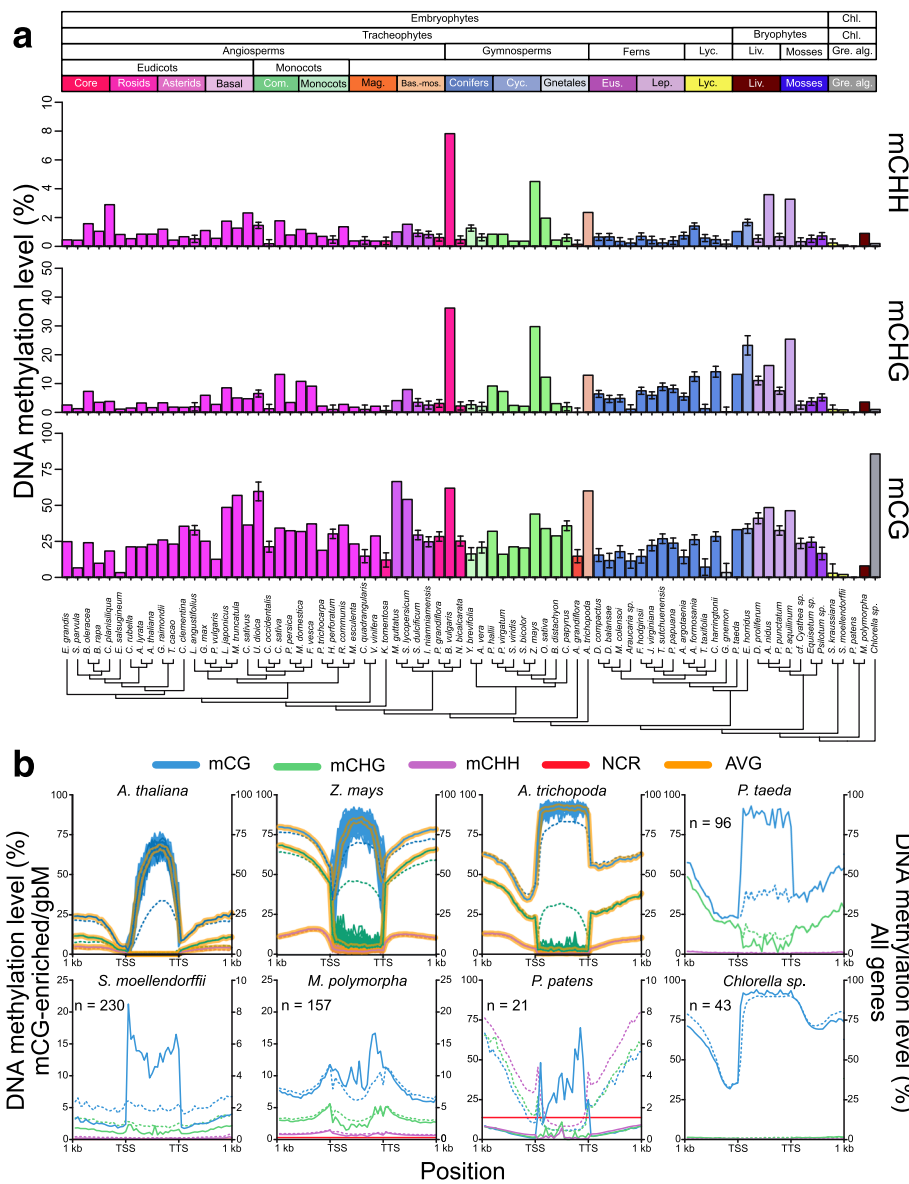
Despite the presence of mCG within the gene bodies of angiosperms, gymnosperms, ferns, lycophytes, liverworts, and green algae the distributions across gene bodies differ (Fig. 3b). In angiosperms (eudicots, commelinids, monocots, and early diverging angiosperm lineages) CG DNA methylation is depleted at the transcriptional start and termination sites (TSS and TTS, respectively) and gradually increases towards the center of the gene body (Fig. 3b). In *Amborella trichopoda*, sister lineage to all other extant angiosperms, levels of mCG decline sharply prior to the TTS (Fig. 3b). Similar to angiosperms, mCG is reduced at the TSS relative to the gene body in the gymnosperm *Pinus taeda* (Fig. 4b). However, mCG is approximately the same as the upstream region at the TTS



**Fig. 2** Non-neutral evolution of CMT3 in the Brassicaceae is correlated with reduced levels of genic mCG and numbers of gbM loci. **a** Distribution of mCG upstream, downstream, and within gene bodies of Brassicaceae species and outgroup species *Prunus persica*. MCG levels within gene bodies of Brassicaceae species are within the bottom 38% of 34 angiosperms. Data used represent a subset of that previously published [19, 25]. *TSS* transcriptional start site, *TTS* transcriptional termination site. **b** Similarly, the number of gbM genes within the genome of Brassicaceae species are within the bottom 15% of 34 angiosperms. The size of the circle corresponds to the number of gbM genes within each genome. Data used represent a subset of that previously published [19, 25]. **c** Changes at the amino acid level of CMT3 is correlated to reduced genic levels of DNA methylation and number of gbM genes in the Brassicaceae. An overall higher rate of molecular evolution, measured as the number of non-synonymous substitutions per non-synonymous site (*dN*) divided by the number of synonymous substitutions per synonymous site (*dS*) ( $\omega$ ), was detected in the Brassicaceae. Also, a higher rate ratio of  $\omega$  was detected in the Brassicaceae clade containing *B. rapa* and closely related species compared to the clade containing *A. thaliana* and closely related species. The higher rate ratio in the Brassicaceae, compared with the background branches, was not attributed to positive selection

(Fig. 3b). Additionally, DNA methylation at non-CG (mCHG and mCHH) sites is not reduced at the TSS and TTS. Little difference in mCG, mCHG, and mCHH within

gene bodies and upstream and downstream regions are observed in *S. moellendorffii* (Fig. 3b). Additionally, mCG, mCHG, and mCHH are not excluded from the TSS and

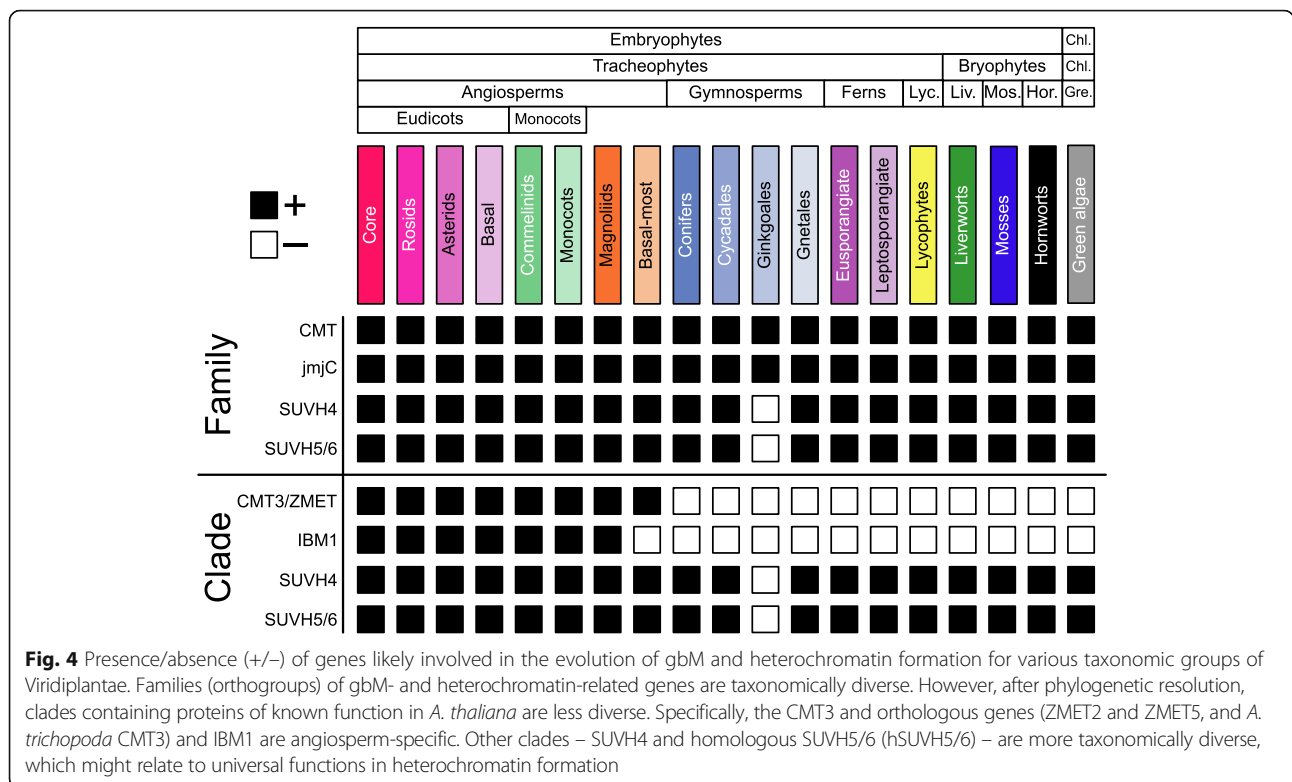


**Fig. 3** Variation in levels of DNA methylation within gene bodies across Viridiplantae. **a** DNA methylation at CG, CHG, and CHH sites within gene bodies can be found in the majority of species investigated. Variation of DNA methylation levels within gene bodies at all sequence contexts is high across all land plants and within major taxonomic groups. mCG levels are typically higher than mCHG, followed by mCHH. However, levels of mCG and mCHG within genes are similar in gymnosperms and ferns. Error bars represent 95% confidence intervals for species with low sequencing coverage. Cladogram was generated from Open Tree of Life [53]. **b** The distribution of DNA methylation within genes (all [dashed lines] and mCG-enriched/gbM [solid lines]) has diverged among taxonomic groups of Viridiplantae represented by specific species. Based on the distribution of DNA methylation, and number of mCG-enriched genes, gbM is specific to angiosperms. However, mCG-enriched genes in *P. taeda* share some DNA methylation characteristics to *A. trichopoda*. However, other characteristics associated with gbM genes remains unknown at this time for mCG-enriched genes in gymnosperms and other early diverging Viridiplantae. The yellow highlighted line represents the average from 100 random sampling of 100 gbM genes in angiosperms and was used to assess biases in numbers of mCG-enriched genes identified. NCR non-conversion rate, TSS transcriptional start site, TTS transcriptional termination site

TTS (Fig. 3b). As opposed to angiosperms and gymnosperms (*P. taeda*), mCG in *M. polymorpha* decreases towards the center of the gene body (Fig. 3b). This distribution also occurs for methylation at non-CG sites in *M. polymorpha* (Fig. 3b). Additionally, *M. polymorpha* has distinctive high levels of mCG, mCHG, and mCHH

surrounding the TSS and TTS (Fig. 3b). Finally, in *Chlorella sp.* NC64A, mCG is enriched to nearly 100% across the entire gene body (Fig. 3b).

The presence of mCG within gene bodies indicates that a gene could possess gbM. However, other types of DNA methylated genes contain high levels of mCG [25], thus



enrichment tests were performed to identify genes that are significantly enriched for mCG and depleted of non-CG methylation (i.e. gbM genes). Genes matching this DNA methylation enrichment profile were identified in species sister to angiosperms: gymnosperms, lycophytes, liverworts, mosses, and green algae (Fig. 3b). The proportion of genes within each genome or subset of the genome (*P. taeda*) was small compared to angiosperms with gbM (Fig. 3b). Furthermore, the number of gbM genes was comparable to angiosperms without gbM, which suggests these identified genes are the result of statistical noise (Additional file 2: Figure S7a). This is most likely the case for lycophytes, liverworts, mosses, and green algae, since the levels of mCG within genes bodies is highly skewed (Additional file 2: Figure S8). Additionally, the distribution of mCG and non-CG methylation across the gene body is unlike the distribution of gbM genes (Fig. 3b) [25]. However, the gymnosperm *P. taeda* shares some similarities to gbM genes of *A. trichopoda* (Fig. 3b). Hence, mCG-enriched genes identified in gymnosperms, lycophytes, liverworts, mosses, and green algae are most likely not gbM.

**Correlated evolution of CMT3 and the histone de-methylase IBM1 in angiosperms**

The exact mechanisms by which genes are targeted to become gbM and the establishment of DNA methylation

at CG sites is currently unknown. One proposed possibility is the failure of IBM1 to remove H3K9me2 within genes. This would provide the necessary substrate for CMT3 to associate with nucleosomes in genes. Due to the tight association between CMT3 and IBM1 (and SUVH4/5/6) these proteins might have evolved together. Resolution of phylogenetic relationships supports monophyly of IBM1 and orthologous sequences that is unique to angiosperms (Fig. 4, Additional file 2: Figure S9). Furthermore, high levels of mCHG and/or similar levels of mCHG to mCG in gymnosperms, ferns, *S. moellendorffii* (lycophyte), *M. polymorpha* (liverwort), and *P. patens* (moss) compared with angiosperms suggests a functionally homologous histone de-methylase is not present in these taxonomic groups and species. The absence of IBM1 in the basal angiosperm *A. trichopoda* and similarities of DNA methylation distribution between gbM genes and *P. taeda* mCG-enriched genes further supports a role of CMT3 and IBM1 in maintenance of mCG within gene bodies. Unlike CMT3 and IBM1, histone methylases SUVH4 and SUVH5/6 are common to all taxonomic groups investigated, which suggests common ancestry and shared functions of transposon silencing (Fig. 4, Additional file 2: Figures S8 and S9) [22, 41–43]. However, a Brassicaceae-specific duplication event gave rise to SUVH5 and SUVH6 and other Viridiplantae possess a homologous SUVH5/6 (hSUVH5/6)



(Additional file 2: Figure S10b and c). Additionally, a duplication event shared by all monocots generated paralogous hSUVH5/6 and an additional duplication event occurred in the Poaceae (Additional file 2: Figure S10d). The duplication event that gave rise to ZMET paralogs in the Poaceae might have also generated the paralogous hSUVH5/6. The diversity in levels and patterns of DNA methylation within gene bodies suggests corresponding changes in function of DNA methyltransferases and/or histone de-methylases during Viridiplantae divergence. Furthermore, monophyletic, angiosperm-specific clades of a gbM-dependent CMT and IBM1 suggest co-evolution of proteins involved in the gbM pathway.

## Discussion

CMTs are conserved DNA methyltransferases across Viridiplantae. Evolutionary phenomenon and forces have shaped the relationships of CMTs, which have most likely contributed to functional divergence among and within taxonomic groups of Viridiplantae. Duplication events have contributed to the unique relationships of CMTs and have given rise to clade-specific, family-specific, and species-specific CMTs. This includes the eudicot-specific CMT1 and CMT3, paralogous CMTs within monocots/commelinids (ZMETs), and the *Z. mays*-specific ZMET2 and ZMET5. The paralogous CMT1 and CMT3, and orthologous CMTs in monocots, monocots/commelinids, magnoliids, and basal angiosperms, form a superclade that is sister to CMT2. Homologous CMTs in gymnosperms and ferns are paraphyletic, and clades are sister to all CMTs – including CMT1, CMT2, CMT3, and ZMET – in angiosperms. CMTs have been shown to maintain methylation at CHG sites (CMT3 and ZMET5, and hCMT $\beta$  in *P. patens*) and methylate CHH sites within deep heterochromatin (CMT2) [2, 4, 6, 11, 14, 15, 20, 21], whereas CMT1 is non-functional in at least some *A. thaliana* accessions [8]. However, recent work has provided evidence for the requirement of CMT3 in the evolution of mCG within gene bodies, and specifically gbM, within angiosperms [19]. Additionally, non-neutral evolution of CMT3 can affect levels of genome-wide mCHG and within gene body mCG, and the number of gbM genes. Hence, functional divergence following duplication might be more widespread [44] and the exact fate of paralogous CMTs and interplay between paralogs in shaping the epigenome remain unknown at this time.

DNA methylation within genes is common in Viridiplantae. However, certain classes of DNA methylated genes might be specific to certain taxonomic groups within the Viridiplantae. gbM is a functionally enigmatic class of DNA methylated genes, which is characterized by an enrichment of mCG and depletion of non-mCG within transcribed regions and depletion of DNA

methylation from all sequence contexts at the TSS and TTS. These genes are typically constitutively expressed, evolutionary conserved, housekeeping genes, which compose a distinct proportion of protein-coding genes [19, 25–27, 45]. gbM genes have been mostly studied in angiosperms and evidence for the existence of this class of DNA methylated gene outside of angiosperms is limited [40]. However, in the present study, genes matching the DNA methylation profile of gbM genes – enrichment of mCG and depletion of non-mCG – were identified in taxonomic groups sister to angiosperms: gymnosperms, lycophytes, liverworts, mosses, and green algae. It is unclear if these genes are gbM genes in light of findings in angiosperms [19, 25–27]. For example, the low proportion of mCG-enriched genes supports the absence of gbM in gymnosperms, lycophytes, liverworts, mosses, and green algae. Additionally, the distribution of mCG among all genes and across the gene body of mCG-enriched genes supports the absence of gbM in lycophytes, liverworts, mosses, and green algae. However, similar distributions of mCG between gbM genes in the basal angiosperm *A. trichopoda* and mCG-enriched genes in the gymnosperm *P. taeda* are observed, which support the presence of gbM in this species and possibly other gymnosperms. Also, a small proportion of mCG-enriched genes in gymnosperms are homologous to gbM genes in *A. thaliana* (Additional file 2: Figure S7b). With that being said, sequence conservation is not the most robust indicator of gbM [25]. gbM genes compose a unique class of genes with predictable characteristics [19, 25–27]. Through comparison of mCG-enriched genes identified in early diverging Viridiplantae to angiosperms with and without gbM, there is stronger support that this epigenetic feature is unique to angiosperms. However, future work including deeper WGBS and RNA sequencing, and additional and improved genome assemblies – especially for gymnosperms and ferns – will undoubtedly contribute to our understanding of the evolution of gbM.

gbM is dependent on the CHG maintenance methyltransferase CMT3 or an orthologous CMT in angiosperms. Support for the dependency of gbM on CMT3 comes from the naturally occurring  $\Delta cmt3$  mutants *C. planisilqua* and *E. salsugineum* which is correlated with the lack of gbM genes [19, 25]. The independent loss of CMT3 has also affected mCHG with low overall and per-site levels recorded for these species [25]. Both species belong to the Brassicaceae family and other species within this family show reduced numbers of gbM genes compared to other eudicot and angiosperm species [25]. Although CMT3 is present in these species, relaxed selective constraint might have introduced alleles, which functionally compromise CMT3 resulting in decreased per-site levels of mCHG and the number of gbM genes [25]. The functional consequences of CMT3 non-neutral evolution are shared and have diverged between

clades of Brassicaceae, which might reflect shared ancestry between clades and the unique evolutionary history of each clade, respectively. Furthermore, more relaxed selective constraint – as in the *Brassica spp.* and *S. parvula* – is correlated with a more severe phenotype relative to the other Brassicaceae clade. The dependency of gbM on a CMT protein might extend into other taxonomic groups of plants. Phylogenetic relationships of CMTs found in Viridiplantae and the location of *A. thaliana* CMTs support a eudicot-specific, monophyletic CMT3 clade. The CMT3 clade is part of a superclade, which includes a monophyletic clade of monocot (ZMET) and magnoliid CMTs and a CMT from the basal angiosperm *A. trichopoda*. Thus, the CMT-dependent gbM pathway might be specific to angiosperms. However, a homologous, closely related CMT in gymnosperms and ferns (i.e. hCMT $\alpha$ ) might have a similar function. It is conceivable that other proteins and chromatin modifications that interact with CMTs and non-CG methylation are important for the evolution of gbM and thus have evolved together. Specifically, IBM1 that de-methylates H3K9me2 and SUVH4/5/6 that binds to mCHG are associated with CMT3 and the maintenance of mCHG. One proposed model for the evolution of gbM requires failure of IBM1 and rare misincorporation of H3K9me2, which initiates mCHG by CMT3, and further di-methylation of H3K9 by SUVH4/5/6 [19, 28]. IBM1 shares similar patterns and taxonomic diversity as CMT3 and orthologous CMTs involved in gbM. Also, unlike most angiosperms investigated to date – with *A. trichopoda* as the exception – gymnosperms and ferns do not possess an IBM1 ortholog, hence IBM1 might be important for the distribution of mCG within gene bodies. Furthermore, the lack of IBM1 in *A. trichopoda* and *P. taeda* might explain some similarities shared between gbM genes and mCG-enriched genes with respect to the deposition of mCG, respectively. However, the exact relationship between gbM and IBM1 is unknown and similarities in underlying nucleotide composition of genes might affect distribution of mCG. Overall, the patterns of DNA methylation within gene bodies and the phylogenetic relationships of CMTs support a CMT3 and orthologous CMT-dependent mechanism for the maintenance of gbM in angiosperms, which is stochastically initiated by IBM1.

## Conclusions

In summary, we present the most comprehensive CMT gene-family phylogeny to date. CMTs are ancient proteins that evolved prior to the diversification of Embryophyta. A shared function of CMTs is the maintenance of DNA methylation at non-CG sites, which has been essential for DNA methylation at long transposable elements in the pericentromeric regions of the genome [6, 14, 15]. However, CMTs in some species of eudicots have been

shown to be important for mCG within gbM genes [19]. Refined relationships between CMT1, CMT2, CMT3, ZMET, and other homologous CMT clades have shed light on hypothesized models for the evolution of gbM and provided a framework for further research on the role of CMTs in establishment and maintenance of DNA methylation and histone modifications. Patterns of DNA methylation within gene bodies have diverged between Viridiplantae. Other taxonomic groups do not share the pattern of mCG associated with gbM genes in the majority of angiosperms, which further supports specificity of gbM in angiosperms. However, genic DNA methylation commonalities between angiosperms and other taxonomic groups were identified. DNA methylation within gene bodies and its consequences of or relationship to expression and other genic features has been extensively studied in angiosperms [25] and shifting focus to other taxonomic groups of plants for deep methylome analyses will aid in understanding the shared consequences of genic DNA methylation. Understanding the evolution of additional chromatin modifiers will undoubtedly unravel the epigenome and reveal unique undiscovered mechanisms.

## Methods

### 1KP sequencing, transcriptome assembling, and orthogrouping

The One Thousand Plants (1KP) Consortium includes assembled transcriptomes and predicted protein-coding sequences from a total of 1329 species of plants (Additional file 1: Table S1). Additionally, gene annotations from 24 additional species – *Arabidopsis lyrata*, *Brachypodium distachyon*, *Brassica oleracea*, *Brassica rapa*, *Citrus clementina*, *Capsella rubella*, *Cannabis sativa*, *Cucumis sativus*, *Eutrema salsugineum*, *Fragaria vesca*, *Glycine max*, *Gossypium raimondii*, *Lotus japonicus*, *Malus domestica*, *Marchantia polymorpha*, *Medicago truncatula*, *Panicum hallii*, *Panicum virgatum*, *Pinus taeda*, *Physcomitrella patens*, *Ricinus communis*, *Setaria viridis*, *Selaginella moellendorffii*, and *Zea mays* – were included (<https://phytozome.jgi.doe.gov/pz/portal.html> and <http://pinenome.org/pinerefseq/>). The CMT gene family was extracted from the previously compiled 1KP orthogroupings using the *A. thaliana* gene identifier for CMT1 (AT1G80740), CMT2 (AT4G19020), and CMT3 (AT1G69770). A single orthogroup determined by the 1KP Consortium included all three *A. thaliana* CMT proteins and a total of 5383 sequences. Sequences from species downloaded from Phytozome, that were not included in sequences generated by 1KP, were included to the gene family through reciprocal best BLAST with *A. thaliana* CMT1, CMT2, and CMT3. In total, the CMT gene family included 5449 sequences from 1043 species. We used the protein structure of *A. thaliana* as a reference to filter the sequences

found within the CMT gene family. Sequences were retained if they included the same base PFAM domains as *A. thaliana* – CHROMO (pfam00385), BAH (pfam01426), and C-5 cytosine-specific DNA methylase (pfam00145) domains – as identified by Interproscan [46]. These filtered sequences represent a set of high-confidence, functional, ideal CMT proteins, which included 771 sequences from 432 species, and were used for phylogenetic analyses. All sequences used in this study can be found in Additional files 5–12.

### Phylogeny construction

To estimate the gene tree for the CMT sequences, a series of alignment and phylogenetic estimation steps were conducted. An initial protein alignment was carried out using Pasta with the default settings [47]. The resulting alignment was back-translated using the coding sequence (CDS) into an in-frame codon alignment. A phylogeny was estimated by RAxML [48] (-m GTRGAMMA) with 1000 rapid bootstrap replicates using the in-frame alignment, and with only the first and second codon positions. Long branches can affect parameter estimation for the substitution model, which can in turn degrade phylogenetic signal. Therefore, phylogenies were constructed with and without green algae species, and were rooted to the green algae clade or liverworts, respectively. The species *Balanophora fungosa* has been reported to have a high substitution rate, which can also produce long branches, and was removed prior to phylogenetic analyses. Identical workflows were used for jumonji (jnjC) domain-containing (i.e. IBM1), SUVH4, and SUVH5/6 gene families.

### Codon analysis

A similar methodology as described above was used to construct phylogenetic trees for testing hypotheses on the rates of evolution in a phylogenetic context. However, the program Gblocks [49] was used to identify conserved codons. The parameters for Gblocks were kept at the default settings, except allowing for 50% gapped positions. The program Phylogenetic Analysis by Maximum Likelihood (PAML) [50] was used to test branches (branch test) and sites along branches (branch-site test) for deviations from the background rate of molecular evolution ( $\omega$ ) and for deviations from the neutral expectation, respectively. Branches tested and a summary of each test can be found in Additional file 4: Table S3.

### MethylC-seq

MethylC-seq libraries were prepared according to the following protocol [51]. For *A. thaliana*, *A. trichopoda*, *Chlorella* sp. NC64A, *M. polymorpha*, *P. patens*, *P. taeda*, *S. moellendorffii*, and *Z. mays* reads were mapped to the

respective genome assemblies. *P. taeda* has a large genome assembly of ~23 Gbp divided among ~14 k scaffolds ([http://dendrome.ucdavis.edu/ftp/Genome\\_Data/genome/pinerefseq/Pita/v1.01/README.txt](http://dendrome.ucdavis.edu/ftp/Genome_Data/genome/pinerefseq/Pita/v1.01/README.txt)). Due to computational limitations imposed by the large genome size only 4 Gbp of the *P. taeda* genome assembly was used for mapping, which includes 2411 (27%) of the high quality gene models. Before mapping for species with only transcriptomes, each transcript was searched for the longest open reading frame from all six possible frames, and only transcripts beginning with a start codon and ending with one of the three stop codons were kept. All sequencing data for each species were aligned to their respective transcriptome or species within the same genus using the methylpy pipeline [52]. GEO or SRA accessions of MethylC-seq data used in this study can be found in Additional file 13: Tables S4 and Additional file 14: Tables S5. Weighted methylation was calculated for each sequence context (CG, CHG, and CHH) by dividing the total number of aligned methylated reads by the total number of methylated plus un-methylated reads. Since per site sequencing coverage was low – on average  $\sim 1\times$  – subsequent binomial tests could not be performed for the majority of species to bin genes as gbM [25]. To investigate the effect of low coverage on estimates of DNA methylation levels, we compared levels of DNA methylation of  $1\times$  randomly sampled MethylC-seq reads to actual levels for 32 angiosperm species, *S. moellendorffii* (lycophyte), *M. polymorpha* (liverwort), and *Chlorella* sp. NC64A (green algae) [19, 20, 25, 40]. Specifically, a linear model was constructed between deep ( $x$ ) and  $1\times$  ( $y$ ) sequencing coverage, which was then used to extrapolate levels of DNA methylation and 95% confidence intervals from low sequence coverage species (Additional file 2: Figure S11, Additional file 14: Table S5).

### Genic DNA methylation analyses and metaplots

DNA methylation was estimated as weighted DNA methylation, which is the total number of aligned DNA methylated reads divide by the total number of methylated plus un-methylated reads. This metric of DNA methylation was estimated for each sequence context within coding regions. For *P. taeda* only high quality gene models were used, since low quality models cannot distinguish between pseudogenes and true protein-coding genes. For genic metaplots, the gene body – start to stop codon – was divided into 20 windows. Additionally, for species with assembled and annotated genomes, regions 1000 or 4000 bp upstream and downstream were divided into 20 windows. Weighted DNA methylation was calculated for each window. The mean weighted methylation for each window was then calculated for all genes and plotted in R v3.2.4 (<https://www.r-project.org/>).

### mCG-enrichment test

Sequence context enrichment for each gene was determined through a binomial test followed by Benjamini–Hochberg false discovery rate [25, 26]. A context-specific background level of DNA methylation determined from the coding sequence was used as a threshold in determining significance. Genes were classified as mCG-enriched/gbM if they had reads mapping to at least 10 CG sites and a *q*-value < 0.05 for mCG and a *q*-value > 0.05 for mCHG and mCHH.

### Additional files

- Additional file 1: Table S1.** Taxonomic, sequence name, and phylogenetic summary of sequences used in this study. (XLSX 82 kb)
- Additional file 2: Figures S1 to S11** with legends. (PDF 23575 kb)
- Additional file 3: Table S2.** Homology-based assessment of CMTs identified in green algae. (XLSX 43 kb)
- Additional file 4: Table S3.** A summary of hypotheses tested using the program PAML. (XLSX 35 kb)
- Additional file 5:** Nucleotide CMT sequences in fasta format. (FAA 619 kb)
- Additional file 6:** Amino acid CMT sequences in fasta format. (FNA 1803 kb)
- Additional file 7:** Nucleotide SUVH4 sequences in fasta format. (FAA 314 kb)
- Additional file 8:** Amino acid SUVH4 sequences in fasta format. (FNA 900 kb)
- Additional file 9:** Nucleotide SUVH5/6 sequences in fasta format. (FAA 479 kb)
- Additional file 10:** Amino acid SUVH5/6 sequences in fasta format. (FNA 1381 kb)
- Additional file 11:** Nucleotide jmjC sequences in fasta format. (FAA 1770 kb)
- Additional file 12:** Amino acid jmjC sequences in fasta format. (FNA 5153 kb)
- Additional file 13: Table S4.** DNA methylation levels estimated from deep and sampled WGBS from 34 Viridiplantae species. (XLSX 61 kb)
- Additional file 14: Table S5.** DNA methylation levels of species sequenced in this study. (XLSX 57 kb)

### Acknowledgements

The authors thank Nathan Springer for comments and discussions. Kevin Turner (UGA greenhouse), and Michael Wenzel and Ron Determann (Atlanta Botanical Gardens) for plant tissue; Brigitte T. Hofmeister for establishment and maintenance of the genome browser; and the Georgia Genomics Facility (GGF) for sequencing. Computational resources were provided by the Georgia Advanced Computing Resource Center (GACRC). The authors thank Gane Ka-Shu Wong and the 1000 Plants initiative (1KP, onekp.com) for advanced access to transcript assemblies.

### Funding

The work was conducted by the National Science Foundation (NSF) (MCB-1402183 and IOS-1339194) and by The Pew Charitable Trusts to RJS.

### Availability of data and materials

Genome browsers for all methylation data used in this paper are located at Plant Methylation DB (<http://schmitzlab.genetics.uga.edu/plantmethyomes>). Sequence data for MethylC-seq are located at the Gene Expression Omnibus, accession GSE81702.

### Authors' contributions

Conceptualization: AJB and RJS; performed experiments: AJB, CEN, LJ, and NAR; data analysis: AJB, CEN, and LJ; writing – original draft: AJB; writing – review and editing: AJB, JL-M, and RJS; resources: JL-M and RJS. All authors read and approved the final manuscript.

### Competing interests

The authors declare that they have no competing interests.

### Ethics approval

Ethics approval was not needed for this study.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

### Author details

<sup>1</sup>Department of Genetics, University of Georgia, Athens, GA 30602, USA.

<sup>2</sup>Institute of Bioinformatics, University of Georgia, Athens, GA 30602, USA.

<sup>3</sup>Department of Plant Biology, University of Georgia, Athens, GA 30602, USA.

Received: 31 January 2017 Accepted: 17 March 2017

Published online: 01 May 2017

### References

- Bartee L, Malagnac F, Bender J. Arabidopsis cmt3 chromomethylase mutations block non-CG methylation and silencing of an endogenous gene. *Genes Dev.* 2001;15:1753–58.
- Jackson JP, Lindroth AM, Cao X, Jacobsen SE. Control of CpNpG DNA methylation by the KRYPTONITE histone H3 methyltransferase. *Nature.* 2002;416:556–60.
- Zemach A, Kim MY, Hsieh PH, Coleman-Derr D, Eshed-Williams L, Thao K, et al. The Arabidopsis nucleosome remodeler DDM1 allows DNA methyltransferases to access H1-containing heterochromatin. *Cell.* 2013;153:193–205.
- Stroud H, Do T, Du J, Zhong X, Feng S, Johnson L, et al. Non-CG methylation patterns shape the epigenetic landscape in Arabidopsis. *Nat Struct Mol Biol.* 2014;21:64–72.
- Papa CM, Springer NM, Muszynski MG, Meeley R, Kaeppler SM. Maize chromomethylase Zea methyltransferase2 is required for CpNpG methylation. *Plant Cell.* 2001;13:1919–28.
- Du J, Zhong X, Bernatavichute YV, Stroud H, Feng S, Caro E, et al. Dual binding of chromomethylase domains to H3K9me2-containing nucleosomes directs DNA methylation in plants. *Cell.* 2012;151:167–80.
- Li Q, Eichten SR, Hermanson PJ, Zaunbrecher VM, Song J, Wendt J, et al. Genetic perturbation of the maize methylome. *Plant Cell.* 2014;26:4602–16.
- Henikoff S, Comai L. A DNA methyltransferase homolog with a chromodomain exists in multiple polymorphic forms in Arabidopsis. *Genetics.* 1998;149:307–18.
- Finnegan EJ, Kovac KA. Plant DNA methyltransferases. *Plant Mol Biol.* 2000;43:189–201.
- McCallum CM, Comai L, Greene EA, Henikoff S. Targeted screening for induced mutations. *Nat Biotechnol.* 2000;18:455–57.
- Shen X, De Jonge J, Forsberg SK, Pettersson ME, Sheng Z, Hennig L, et al. Natural CMT2 variation is associated with genome-wide methylation changes and temperature seasonality. *PLoS Genet.* 2014;10:e1004842.
- Bewick AJ, Schmitz RJ. Epigenetics in the wild. *elife.* 2015;4:e07808.
- Dubin MJ, Zhang P, Meng D, Remigereau MS, Osborne EJ, Paolo Casale F, et al. DNA methylation in Arabidopsis has a genetic basis and shows evidence of local adaptation. *elife.* 2015;4:e05255.
- Du J, Johnson LM, Groth M, Feng S, Hale CJ, Li S, et al. Mechanism of DNA methylation-directed histone methylation by KRYPTONITE. *Mol Cell.* 2014;55:495–504.
- Du J, Johnson LM, Jacobsen SE, Patel DJ. DNA methylation pathways and their crosstalk with histone methylation. *Nat Rev Mol Cell Biol.* 2015;16:519–32.
- Hou PQ, Lee YI, Hsu KT, Lin YT, Wu WZ, Lin JY, et al. Functional characterization of Nicotiana benthamiana chromomethylase 3 in developmental programs by virus-induced gene silencing. *Physiol Plant.* 2013;150:119–32.
- Garg R, Kumari R, Tiwari S, Goyal S. Genomic survey, gene expression analysis and structural modeling suggest diverse roles of DNA methyltransferases in legumes. *PLoS One.* 2014;9:e88947.
- Lin YT, Wei HM, Lu HY, Lee YI, Fu SF. Developmental- and tissue-specific expression of NbCMT3-2 encoding a chromomethylase in Nicotiana benthamiana. *Plant Cell Physiol.* 2015;56:1124–43.

19. Bewick AJ, Ji L, Niederhuth CE, Willing EM, Hofmeister BT, Shi X, et al. On the origin and evolutionary consequences of gene body DNA methylation. *Proc Natl Acad Sci U S A*. 2016;113:9111–16.
20. Zemach A, McDaniel IE, Silva P, Zilberman D. Genome-wide evolutionary analysis of eukaryotic DNA methylation. *Science*. 2010;328:916–9.
21. Noy-Malka C, Yaari R, Itzhaki R, Mosquna A, Auerbach Gershovitz N, Katz A, et al. A single CMT methyltransferase homolog is involved in CHG DNA methylation and development of *Physcomitrella patens*. *Plant Mol Biol*. 2014;84:719–35.
22. Tran RK, Henikoff JG, Zilberman D, Ditt RF, Jacobsen SE, Henikoff S. DNA methylation profiling identifies CG methylation clusters in *Arabidopsis* genes. *Curr Biol*. 2005;15:154–9.
23. Zhang X, Yazaki J, Sundaresan A, Cokus S, Chan SW, Chen H, et al. Genome-wide high-resolution mapping and functional analysis of DNA methylation in *Arabidopsis*. *Cell*. 2006;126:1189–201.
24. Zilberman D, Gehring M, Tran RK, Ballinger T, Henikoff S. Genomewide analysis of *Arabidopsis thaliana* DNA methylation uncovers an interdependence between methylation and transcription. *Nat Genet*. 2007;39:61–9.
25. Niederhuth CE, Bewick AJ, Ji L, Alabady MS, Kim KD, Li Q, et al. Widespread natural variation of DNA methylation within angiosperms. *Genome Biol*. 2016;17:194.
26. Takuno S, Gaut BS. Body-methylated genes in *Arabidopsis thaliana* are functionally important and evolve slowly. *Mol Biol Evol*. 2012;1:219–27.
27. Takuno S, Gaut BS. Gene body methylation is conserved between plant orthologs and is of evolutionary consequence. *Proc Natl Acad Sci U S A*. 2013;110:1797–802.
28. Inagaki S, Kakutani T. What triggers differential DNA methylation of genes and TEs: contribution of body methylation? *Cold Spring Harb Symp Quant Biol*. 2012;77:155–60.
29. Johnson LM, Du J, Hale CJ, Bischof S, Feng S, Chodavarapu RK, et al. SRA- and SET-domain-containing proteins link RNA polymerase V occupancy to DNA methylation. *Nature*. 2014;507:124–8.
30. Kenrick P, Crane PR. The origin and early evolution of plants on land. *Nature*. 1997;389:33–9.
31. Wellman CH, Osterloff PL, Mohiuddin U. Fragments of the earliest land plants. *Nature*. 2003;425:282–5.
32. Steemans P, Herisse AL, Melvin J, Miller MA, Paris F, Verniers J, et al. Origin and radiation of the earliest vascular land plants. *Science*. 2009;324:353.
33. Rubinstein CV, Gerrienne P, de la Puente GS, Astini RA, Steemans P. Early Middle Ordovician evidence for land plants in Argentina (eastern Gondwana). *New Phytol*. 2010;188:365–9.
34. Stiller JW, Hall BD. The origin of red algae: Implications for plastid evolution. *Proc Natl Acad Sci U S A*. 1997;94:4520–5.
35. Bhattacharya D, Medlin L. Algal phylogeny and the origin of land plants. *Plant Physiol*. 1998;116:9–15.
36. Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, et al. Ancestral polyploidy in seed plants and angiosperms. *Nature*. 2011;473:97–100.
37. Wickett NJ, Mirarab S, Nguyen N, Warnow T, Carpenter E, Matasci N, et al. Phylotranscriptomic analysis of the origin and early diversification of land plants. *Proc Natl Acad Sci U S A*. 2014;111:E4859–68.
38. Malik G, Dangwal M, Kapoor S, Kapoor M. Role of DNA methylation in growth and differentiation in *Physcomitrella patens* and characterization of cytosine DNA methyltransferases. *FEBS J*. 2012;279:4081–94.
39. Feng S, Cokus SJ, Zhang X, Chen PY, Bostick M, Goll MG, et al. Conservation and divergence of methylation patterning in plants and animals. *Proc Natl Acad Sci U S A*. 2010;107:8689–94.
40. Takuno S, Ran J-H, Gaut BS. Evolutionary patterns of genic DNA methylation vary across land plants. *Nature Plants*. 2016;2:15222.
41. Lippman Z, May B, Yordan C, Singer T, Martienssen R. Distinct mechanisms determine transposon inheritance and methylation via small interfering RNA and histone modification. *PLoS Biol*. 2003;1:E67.
42. Qin FJ, Sun QW, Huang LM, Chen XS, Zhou DX. Rice SUVH histone methyltransferase genes display specific functions in chromatin modification and retrotransposon repression. *Mol Plant*. 2010;3:773–82.
43. Stroud H, Greenber MVC, Feng S, Bernatavichute YV, Jacobsen SE. Comprehensive analysis of silencing mutants reveals complex regulation of the *Arabidopsis* methylome. *Cell*. 2013;152:352–64.
44. Lynch M, Conery JS. The evolutionary fate and consequences of duplicate genes. *Science*. 2000;290:1151–5.
45. Aceituno FF, Moseyko N, Rhee SY, Gutiérrez RA. The rules of gene expression in plants: organ identity and gene body methylation are key factors for regulation of gene expression in *Arabidopsis thaliana*. *BMC Genomics*. 2008;9:438. doi:10.1186/1471-2164-9-438.
46. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics*. 2014;30:1236–40.
47. Mirarab S, Nguyen N, Guo S, Wang LS, Kim J, Warnow T. PASTA: ultra-large multiple sequence alignment for nucleotide and amino-acid sequences. *J Comput Biol*. 2015;22:377–86.
48. Stamatakis A. RAxML Version 8: A tool for phylogenetic analysis and postanalysis of large phylogenies. *Bioinformatics*. 2014;30:1312–3.
49. Castresana J. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol*. 2000;17:540–52.
50. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol*. 1997;24:1586–91.
51. Urich MA, Nery JR, Lister R, Schmitz RJ, Ecker JR. MethylCseq library preparation for base-resolution whole-genome bisulfite sequencing. *Nat Protoc*. 2015;10:475–83.
52. Schultz MD, He Y, Whitaker JW, Hariharan M, Mukamel EA, Leung D, et al. Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature*. 2015;523:212–6.
53. Hinchliff CE, Smith SA, Allman JF, Burleigh JG, Chaudhary R, Coghill LM, et al. Synthesis of phylogeny and taxonomy into a comprehensive tree of life. *Proc Natl Acad Sci U S A*. 2015;112:12764–9.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at  
[www.biomedcentral.com/submit](http://www.biomedcentral.com/submit)

