

1           **Widespread conservation and lineage-specific diversification of genome-wide DNA**  
2   **methylation patterns across arthropods**

3       Lewis, S.<sup>1,2,3</sup>, Ross L.<sup>5</sup>, Bain, S.A.<sup>5</sup>, Pahita, E.<sup>2,3</sup>, Smith, S.A.<sup>8</sup>, Cordaux, R.<sup>7</sup>, Miska, E.M.<sup>1,4</sup>, Lenhard,  
4   B.<sup>2,3</sup>, Jiggins, F.M.<sup>1\*†</sup> & Sarkies, P.<sup>2,3\*†</sup>

- 5           1) Department of Genetics, University of Cambridge  
6           2) MRC London Institute of Medical Sciences, Du Cane Road, London, W120NN  
7           3) Institute of Clinical Sciences, Imperial College London, Du Cane Road, London, W12 0NN  
8           4) Wellcome Trust/Cancer Research UK Gurdon Institute, Tennis Court Road, Cambridge  
9           5) Institute of Evolutionary Biology, Edinburgh, UK  
10          6) Department of Biomedical Sciences and Pathobiology, Virginia Maryland College of  
11             Veterinary Medicine, Virginia Tech, USA  
12          7) Laboratoire Ecologie et Biologie des Interactions Universite de Poitiers, France  
13          8) Department of Biomedical Sciences and Pathology, Virginia Maryland College of Veterinary  
14             Medicine, 205 Duck Pond Drive, Virginia Tech, Blacksburg, VA24061, USA

15          † Contributed equally

17          \*Correspondence to

18          Francis Jiggins, [fmj1001@cam.ac.uk](mailto:fmj1001@cam.ac.uk)

19          Peter Sarkies, [psarkies@imperial.ac.uk](mailto:psarkies@imperial.ac.uk)

21       **Abstract**

22       Cytosine methylation is an ancient epigenetic modification yet its function and extent within genomes  
23       is highly variable across eukaryotes. In mammals, methylation controls transposable elements and  
24       regulates the promoters of genes. In insects, DNA methylation is generally restricted to a small subset  
25       of transcribed genes, with both intergenic regions and transposable elements (TEs) depleted of  
26       methylation. The evolutionary origin and the function of these methylation patterns are poorly  
27       understood. Here we characterise the evolution of DNA methylation across the arthropod phylum.  
28       While the common ancestor of the arthropods had low levels of TE methylation and did not methylate

1 promoters, both of these functions have evolved independently in centipedes and mealybugs. In  
2 contrast, methylation of the exons of a subset of transcribed genes is ancestral and widely conserved  
3 across the phylum, but has been lost in specific lineages. Remarkably the same set of genes are  
4 likely to be methylated in all species that retained exon-enriched methylation. We show that these  
5 genes have characteristic patterns of expression correlating to broad transcription initiation sites and  
6 well-positioned nucleosomes, providing new insights into potential mechanisms driving methylation  
7 patterns over hundreds of millions of years.

## 8 **Author Summary**

9 Animals develop from a single cell to form a complex organism with many specialised cells.  
10 Almost all of the fantastic variety of cells must have the same sequence of DNA, and yet  
11 they have distinct identities that are preserved even when they divide. This remarkable  
12 process is achieved by turning different sets of genes on or off in different types of cell using  
13 molecular mechanisms known as “epigenetic gene regulation”.

14 Surprisingly, though all animals need epigenetic gene regulation, there is a huge diversity in  
15 the mechanisms that they use. Characterising and explaining this diversity is crucial in  
16 understanding the functions of epigenetic pathways, many of which have key roles in human  
17 disease. We studied how one particular type of epigenetic regulation, known as DNA  
18 methylation, has evolved within arthropods. Arthropods are an extraordinarily diverse group  
19 of animals ranging from horseshoe crabs to fruit flies. We discovered that the levels of DNA  
20 methylation and where it is found within the genome changes rapidly throughout arthropod  
21 evolution. Nevertheless, there are some features of DNA methylation that seem to be the  
22 same across most arthropods- in particular we found that there is a tendency for a similar  
23 set of genes to acquire methylation of DNA in most arthropods, and that this is conserved  
24 over 350 million years. We discovered that these genes have distinct features that might  
25 explain how methylation gets targeted. Our work provides important new insights into the  
26 evolution of DNA methylation and gives some new hints to its essential functions.

27

28

29

30

31

## 1 **Introduction**

2 In most organisms DNA bases are adorned with a variety of chemical modifications. Amongst the  
3 most common of these is methylation at the 5 position of cytosine (C5me), which is present from  
4 bacteria to humans (Ponger and Li, 2005; Casadesús and Low, 2006; Jurkowski and Jeltsch, 2011).  
5 In eukaryotes, a key property of cytosine DNA methylation is its ability to act epigenetically — that is,  
6 once introduced, methylation at specific cytosines can remain in place through cell division (Holliday  
7 *et al.*, 1987; Holliday, 2006). This relies on the activity of “maintenance” methyltransferases, DNMT1  
8 in animals (Law and Jacobsen, 2010), which recognise CG dinucleotides (CpG sites) where one  
9 strand is methylated and one strand unmethylated and catalyse the introduction of methylation on the  
10 unmethylated strand (Jeltsch, 2006). Meanwhile “de novo” methyltransferases act on unmethylated  
11 DNA. In animals this role is performed by DNMT3 enzymes, which introduce 5meC predominantly  
12 within CpG sites (Jeltsch, 2006). Mechanisms also exist to remove methylation from DNA, including  
13 the TET family of enzymes, which convert 5meC to a hydroxymethylated intermediate which can be  
14 removed by base excision repair or diluted out through cell division (Nashun, Hill and Hajkova, 2015).  
15 As the maintenance and de novo methylation of CG sequences occurs through the activity of  
16 homologous enzymes in plants and animals (Ponger and Li, 2005), CpG methylation was likely  
17 present among the earliest eukaryotic organisms.

18 In mammals, a key function of CpG methylation is to defend the genome against transposable  
19 elements (TEs) by preventing their transcription and transposition (Bird, 2002), and loss of DNA  
20 methylation leads to reactivation of TEs (Walsh, Chaillet and Bestor, 1998). CpG methylation targeted  
21 to the promoters of genes can also suppress transcription, typically resulting in stable silencing (Bird,  
22 2002). Another notable genome-wide pattern is the enrichment of CpG methylation within the exons  
23 of transcribed genes (Suzuki and Bird, 2008). In contrast to TE and promoter methylation, this is not  
24 associated with transcriptional silencing.

25 Whilst CpG methylation at both TEs and gene bodies is present in both plants and animals (Law and  
26 Jacobsen, 2010), across eukaryotic species both DNA methylation levels and the targets of  
27 methylation have evolved rapidly (Feng *et al.*, 2010; Zemach *et al.*, 2010). Most strikingly, CpG  
28 methylation has been independently lost altogether several times, coinciding with the loss of DNMT1  
29 and DNMT3 DNA methyltransferase enzymes (Ponger and Li, 2005; Feng *et al.*, 2010; Zemach *et al.*,

1 2010). Across eukaryotes, loss of CpG methylation tends to be accompanied by loss of the DNA  
2 alkylation repair enzyme ALKB2, which repairs damage caused by DNMTs introducing 3-  
3 methylcytosine into DNA. This suggests that some species correct DNA alkylation using ALKB2, and  
4 others avoid it altogether by losing the DNA methylation pathway (Rošić *et al.*, 2018). Even within  
5 species that retain DNA methyltransferases, the genomic distribution of CpG methylation differs  
6 widely (Feng *et al.*, 2010; Zemach *et al.*, 2010; Bewick *et al.*, 2017, 2019; Rošić *et al.*, 2018; de  
7 Mendoza *et al.*, 2019; de Mendoza, Pflueger and Lister, 2019). Such variability is surprising given the  
8 essential role of CpG methylation in genome regulation in mammals and plants, and there are few  
9 clues as to what factors drive the changes. Tracing the evolution of CpG methylation is currently  
10 challenging because detailed descriptions of DNA methylation patterns are patchy and focussed on  
11 model systems, leaving large parts of the phylogenetic tree underexplored.

12 Here we study CpG methylation patterns across arthropods. Arthropods have been suggested to  
13 represent a very different system of CpG methylation from mammals (Keller, Han and Yi, 2016).  
14 Whilst the well-characterised model organism *Drosophila melanogaster* lacks DNA methylation  
15 altogether, DNA methyltransferases 1 and 3 were found in the honey bee *Apis mellifera* (Wang *et al.*,  
16 2006). Genome-wide CpG methylation mapping demonstrated that methylation was globally  
17 extremely low, and restricted to the bodies of a subset of transcribed genes (Lyko *et al.*, 2010;  
18 Zemach *et al.*, 2010). Subsequently, similarly restricted patterns of DNA methylation were found in  
19 other insects (Lyko *et al.*, 2010; Xiang *et al.*, 2010; Bonasio *et al.*, 2012; Wang *et al.*, 2013). Such  
20 patterns support the proposal that gene body methylation is conserved across eukaryotes while TE  
21 methylation has been lost altogether in arthropods (Zemach *et al.*, 2010; Keller, Han and Yi, 2016).  
22 However some insects show considerably higher levels of genome-wide methylation (Bewick *et al.*,  
23 2017), and variation in arthropod methylation levels also exists outside of insects (Kao *et al.*, 2016;  
24 Kvist *et al.*, 2018; de Mendoza, Pflueger and Lister, 2019; Liu *et al.*, 2019). There is also evidence of  
25 TE methylation in the desert locust *Schistocerca gregaria* (Lyko *et al.*, 2010). A thorough  
26 reconstruction of the evolution of methylation across the phylum is still lacking.

27 We set out to explore the evolution of arthropod methylation patterns by characterising genome-wide  
28 CpG methylation across the phylum. We show that TE methylation was ancestral to arthropods,  
29 although at a relatively low level. Methylation of protein-coding genes was also ancestral to  
30 arthropods, with similar subsets of genes being targeted for methylation across arthropods. Despite

1 these conserved features, we find many examples of diversification in methylation patterns across  
2 arthropods, in particular loss of gene methylation in crustaceans and gain of both promoter  
3 methylation and genome-wide TE methylation in the myriapod *Strigamia maritima* and the insect  
4 *Planococcus citri*. We find that methylation at genes, enriched within exons, is the most widely  
5 conserved feature of arthropod methylomes and we use comparative analysis to identify a link  
6 between exon methylation and nucleosome positioning. Overall, our findings demonstrate that while  
7 key features of global methylation patterns have been conserved across millions of years of arthropod  
8 evolution, the targets of DNA methylation can rapidly diverge within individual lineages.

9

## 10 **Results**

### 11 **Genome-wide levels of CpG methylation vary widely across the arthropods**

12 We carried out high-coverage whole-genome bisulphite sequencing (WGBS) on 14 species of  
13 arthropod and quantified the levels of DNA methylation with base-pair resolution. To examine  
14 genome-wide methylation levels we combined this data with published results from 15 additional  
15 species (Bewick *et al.*, 2017; Wu *et al.*, 2017; Kvist *et al.*, 2018) which were mostly sequenced at  
16 lower coverage. Estimates of genome-wide CpG methylation were then used to reconstruct ancestral  
17 methylation levels across the arthropod phylogeny. All 18 species of holometabolous insects had low  
18 levels of CpG methylation, and this was likely the ancestral state of this clade (Figure 1A and 1B).  
19 While CpG methylation rates in other arthropod clades tended to be higher, they varied considerably  
20 (Figure 1A and 1B). The ancestral arthropod likely had moderate methylation levels ( $8.59 \pm 4.8\%$ ;  
21 Figure 1A) but higher methylation levels evolved in *S. maritima*. Similarly, the ancestor of insects had  
22 methylation levels lower than some taxa such as *B. germanica* ( $3.9 \pm 3.3\%$  versus 12%) indicating that  
23 methylation level fluctuated throughout arthropod evolution.

24 To investigate the evolution of the DNA methylation machinery across arthropods, we searched the  
25 genomes of these species for homologues of the genes encoding the methyltransferases DNMT1-3.  
26 We confirmed the genes all encoded a full cytosine methyltransferase domain, and where we did not  
27 find annotated homologues we directly search the genomic DNA for unannotated genes. In each  
28 species we found a single copy of DNMT2, which methylates tRNAs (Goll *et al.*, 2006) (Figure 1C).  
29 DNMT1 was present in all species apart from the five Diptera (Figure 1C). The loss of this gene was

1 associated with the loss of CpG methylation (Figure 1c), with methylation rates in *D. melanogaster* not  
2 significantly different from the unmethylated DNA spike-in included in each reaction. DNMT3 was  
3 absent from the genomes of 14 species, with inspection of the tree suggesting at least eight  
4 independent losses (Figure 1C). Several of these species possessed moderate levels of CpG  
5 methylation (Figure 1B), indicating that DNMT1 alone can be sufficient for introducing genome-wide  
6 DNA methylation, consistent with earlier studies in arthropods and nematodes (Xiang *et al.*, 2010;  
7 Bewick *et al.*, 2017; Rošić *et al.*, 2018).

8 Across the eukaryotes ALKB2, which repairs DNA alkylation damage introduced by DNMTs, tends to  
9 be lost from the same taxa as DNMT1 and 3 (Rošić *et al.*, 2018). Arthropods exhibited many  
10 exceptions to this general rule—there have been at least five losses of ALKB2 but only one of these is  
11 associated with the loss of DNMT1 and 3 (Figure 1C). However, we found that species with ALKB2  
12 possessed higher methylation levels (Figure 1 Supplement; phylogenetic mixed model:  $p=0.0182$ ),  
13 suggesting ALKB2 is dispensable in species with low levels of DNA methylation.

#### 14 **Rapid loss and gain of TE methylation across arthropods**

15 In mammals, plants and nematodes, transposable elements (TEs) are preeminent targets of DNA  
16 methylation, but previous studies have shown that TE methylation is rare in holometabolous insects  
17 (Feng *et al.*, 2010; Lyko *et al.*, 2010; Zemach *et al.*, 2010; Bonasio *et al.*, 2012; Wang *et al.*, 2013).  
18 However, DNA methylation has been found at TEs in some arthropods (Falckenhayn *et al.*, 2013; Kao  
19 *et al.*, 2016; de Mendoza, Pflueger and Lister, 2019; Liu *et al.*, 2019). To explore the distribution of TE  
20 methylation across arthropods we annotated transposable elements using RepeatMasker analysis of  
21 the entire genome, and removed annotations that did not contain Pfam domains derived from  
22 transposable elements. We focused on 14 species that represent the diversity of arthropods, and  
23 have assembled and annotated genomes (see Fig 2B).

24 Compared to unannotated regions of the genome, TEs were strongly enriched for methylation in *S.*  
25 *maritima* and *P. citri*, and weakly enriched in several other species (Figure 2B,C). This pattern is  
26 reflected in the distribution of methylation across TEs — this is skewed towards 0% for most species,  
27 but in *S. maritima* and *P. citri* the large majority of TEs are methylated (Figure 2A,B; Figure 2  
28 Supplement). In these two species there was a sharp drop in methylation rates at the boundary of the  
29 TE (Figure 2D). In agreement with earlier studies (Lyko *et al.*, 2010; Bonasio *et al.*, 2012), the

1 methylation rate of TEs was low in holometabolous insects. However, outside of this group there was  
 2 moderate methylation of TEs in chelicerates (*L. polyphemus* and *P. tepidariorum*), the crustacean *P.*  
 3 *hawaiensis* and hemimetabolous insects (*B. germanica* and *A. pisum*) (Figure 1A,C). To further  
 4 quantify the extent of TE methylation, we clustered TEs into highly- and lowly-methylated groups in  
 5 each species separately, and calculated the proportion of TEs that were assigned to the highly-  
 6 methylated group (Table 1). The large majority of TEs were targeted by methylation in *S. maritima*  
 7 and *P. citri*, while in all other species under 15% of TEs were methylated. Ancestral state  
 8 reconstruction suggested that a low level of TE methylation was present in the ancestral arthropod,  
 9 but was lost in the ancestor of holometabolous insects (Figure 2A).

10

11 **Table 1. Proportion of Genes and TEs that are highly methylated**

Species	TEs <sup>a</sup>		Genes	
	Number	Proportion methylated <sup>b</sup>	Number	Proportion methylated <sup>b</sup>
<i>Acyrtosiphon pisum</i>	293	0.017	13147	0.171
<i>Apis mellifera</i>	7	0.143	10066	0.272
<i>Armadillidium vulgare</i>	655	0.020	4703	0.019
<i>Blattella germanica</i>	276	0.145	9272	0.387
<i>Bombus terrestris</i>	78	0.128	8550	0.069
<i>Heliconius melpomene</i>	34	0.088	11583	0.077
<i>Ixodes scapularis</i>	212	0.033	5775	0.219
<i>Limulus polyphemus</i>	342	0.117	7227	0.265
<i>Nicrophorus vespilloides</i>	9	0.111	12305	0.032
<i>Parasteatoda tepidariorum</i>	622	0.032	9742	0.243
<i>Parhyale hawaiensis</i>	89	0.079	3302	0.028
<i>Planococcus citri</i>	361	0.751	34044	0.099
<i>Strigamia maritima</i>	719	0.758	12898	0.326

12 <sup>a</sup> TEs with annotated TE-associated domains (see Methods); <sup>b</sup> the proportion falling into the highly  
 13 methylated group after clustering each feature type within each species

14

### 15 **Methylation at exons is conserved across most arthropods**

16 We next investigated methylation at genes across arthropods. In all but one of the species we tested,  
 17 mean methylation levels across exons were significantly higher than unannotated regions of the  
 18 genome (Figure 3B). The exception was *P. hawaiensis*, where exons are significantly less methylated  
 19 than unannotated regions of the genome (Figure 3B). There is a significant difference between  
 20 methylation at exons and introns in *P. hawaiensis* ( $p=0.001$ , paired t test). In the species with exon

1 methylation, the distribution of methylation suggested that a subset of genes is targeted for  
2 methylation (Figure 2C). When clustered into highly and lowly methylated genes, the proportion of  
3 methylated genes varied similarly to mean methylation across genes (Table 1).

4 To investigate the distribution of methylation within genes, we compared the methylation levels at  
5 exons and introns in each species. Methylation was higher at exons in the majority of species,  
6 suggesting that the gene body methylation in arthropods is due to targeting of methylation to exons.  
7 However, there was little difference between exons and introns for the two crustaceans, *P.*  
8 *hawaiensis* and *A. vulgare* (Figure 3C; supplemental Figure S3). Given that *P. hawaiensis* exons are  
9 depleted for methylation relative to the genome-wide background while *A. vulgare* exons are only  
10 slightly greater than the background, this may reflect an ancient loss of gene body methylation in the  
11 ancestor of these species. Among species with exon methylation, there were differences in how  
12 methylation levels changed across the gene (Figure 3C). For example, methylation was largely  
13 confined to the first three exons of *P. citri* and *N. vespilloides*, while methylation in *B. germanica* is  
14 largely found from exon four onwards (Figure 3C). Together these data suggest that exon-enriched  
15 methylation was an ancestral property of arthropod methylomes which is largely conserved across the  
16 phylum.

### 17 **Independent acquisition of promoter methylation in arthropod lineages**

18 In mammals, methylation of regions immediately upstream of genes, often at CpG islands, is  
19 associated with gene silencing. However, there is no evidence of promoter methylation in insects  
20 (Lyko *et al.*, 2010; Xiang *et al.*, 2010; Bonasio *et al.*, 2012). To examine promoter methylation  
21 associated with gene silencing across arthropods, we extracted 1kb upstream of genes for all  
22 species. In most species there was little difference in upstream methylation between high and low  
23 expression genes; however, low expression genes in *P. citri* and *S. maritima* had significantly higher  
24 upstream methylation than high expression genes (Figure 4A). In *S. maritima* only genes with very  
25 high upstream methylation showed clearly reduced gene expression ( $p=1e-15$ , Kruskal Wallis test),  
26 whilst in *P. citri* there was a positive correlation between upstream methylation and gene expression  
27 across a wider range of upstream methylation levels (Figure 4B). The different relationship between  
28 upstream methylation and gene expression between *S. maritima* and *P. citri* and the lack of a similar



1 relationship in other arthropod species suggests that promoter methylation associated with gene  
2 silencing may have evolved independently in these two species.

### 3 **Methylated genes are conserved and have moderate to high expression**

4 Our results suggest that the most highly conserved feature of arthropod methylomes is enrichment of  
5 methylation at the exons of a subset of genes. Across species, we asked whether there was any  
6 tendency for orthologous genes to be methylated in different species. We ranked orthologous genes  
7 by relative methylation levels across species and observed that there was a clear tendency for  
8 orthologs to have similar levels of methylation in different species (Figure 5A). The observation that  
9 the same genes are methylated in different species raised the question of what determines which  
10 genes acquire methylation. We used comparative analysis to investigate this across the phylum.

11 Methylation has been shown to be enriched at alternatively spliced genes in some insects (Lyko *et al.*,  
12 2010; Bonasio *et al.*, 2012). To test for a link between methylation and splicing across arthropods, we  
13 compared the level of methylation between genes with one exon (which cannot undergo splicing) and  
14 genes with two or more exons (which may undergo splicing). We found no clear difference in any  
15 species (Figure S5), suggesting that splicing does not explain the propensity of genes to acquire  
16 methylation across arthropods.

17 Previously, methylation of genes in individual insect species has been correlated to higher levels of  
18 expression (Xiang *et al.*, 2010; Bonasio *et al.*, 2012). We find a statistically significant tendency for  
19 genes with high methylation to have higher expression across most species (Table S2). However,  
20 many highly expressed genes are not methylated. Instead a more prominent trend is for methylated  
21 genes to have more focussed levels of gene expression such that genes with very low expression  
22 levels are rarely methylated (Figure 5B,C; Figure 5 supplement 2). Curiously, this pattern is reversed  
23 in *P. citri*, where the exons of methylated genes tend to have low expression (Figure 5 supplement 2).

24 Previously it has been noted that methylated genes are more likely to perform conserved  
25 “housekeeping” functions (Hunt *et al.*, 2013). We clustered genes into orthologous groups across  
26 species and examined genes that were conserved across all species compared to species-specific  
27 genes. Across all species carrying gene body methylation, conserved genes with moderate to high  
28 expression were more likely to be methylated (Figure 5C; Figure S5). Nevertheless many conserved

1 and highly expressed genes lacked methylation suggesting that neither conservation nor expression  
2 is sufficient to explain gene body methylation.

### 3 **Nucleosome positioning influences DNA methylation levels across arthropods**

4 In order to investigate molecular mechanisms that might be responsible for influencing DNA  
5 methylation we examined how the correlation in methylation between pairs of CpGs varied with  
6 increasing separation. In many species with exon-enriched methylation the correlation coefficient  
7 between methylation levels of individual CpGs oscillated periodically (Figure 6A,B). Fourier analysis  
8 showed that the period of oscillation was ~160 nucleotides, roughly corresponding to the average  
9 nucleosome repeat length (Figure 6A,B; Figure S6-1). We quantified this nucleosome-length  
10 periodicity within exons across all species. While the majority of species with exonic methylation  
11 displayed a nucleosome periodicity signal, its magnitude varied greatly – for example *H. melpomene*  
12 has gene methylation but less apparent periodicity (Figure 6B). Interestingly a clear signal of  
13 periodicity was also seen for TE methylation in *S. maritima* and *P. citri*, both of which have high levels  
14 of TE methylation (Figure S6-1).

15 We wondered whether the periodicity in correlation between methylated DNA might reflect an  
16 influence of nucleosome positioning on DNA methylation, as has been shown in plants (Chodavarapu  
17 *et al.*, 2010) and inferred from analysis of mammalian DNA methylation profiles. In the absence of  
18 genome-wide nucleosome positioning data for the majority of species, we investigated nucleosome  
19 positioning from *Drosophila* (Ho *et al.*, 2014), examining orthologues of genes either enriched or  
20 depleted for DNA methylation across arthropods. The promoters of methylated genes possessed high  
21 nucleosome occupancy overall and strongly positioned nucleosomes just upstream (-1) and  
22 downstream (+1) of the transcription start site (TSS) (Figure 6C). The promoters of unmethylated  
23 genes showed lower nucleosome occupancy overall and demonstrated weaker positioning of the -1  
24 and +1 nucleosome. Previous analyses of promoter types across eukaryotes have indicated that  
25 promoters with strong positioning of nucleosomes lead to initiation of transcription across a broad  
26 region (broad TSS) whilst promoters with weaker nucleosome positioning tend to have a much  
27 narrower TSS focussed around a dominant initiation site (Haberle and Lenhard, 2016). Using cap  
28 analysis of gene expression (CAGE) data from *D. melanogaster* we found that the TSS of *D.*

1 *melanogaster* orthologs of methylated genes was broader than the TSS of orthologs unmethylated  
2 genes (Figure 6C).

3 Further evidence for a connection between nucleosome occupancy and a periodic signal in the  
4 correlation between methylation sites comes from a comparison of exons and introns. Exons are  
5 known to have much higher nucleosome occupancy than introns and accordingly the periodic signal  
6 of methylation correlation is markedly weaker in introns than in exons (Figure S6-2). Together this  
7 supports a potential role for nucleosome occupancy in shaping CpG methylation patterns in  
8 arthropods.

9 The alternative patterns of nucleosome occupancy and transcription initiation corresponded to  
10 previous analyses across organisms demonstrating that housekeeping genes tend to have well  
11 positioned nucleosomes just downstream of promoters and broad TSS whereas tissue-specific genes  
12 tend to have less well-defined nucleosome positions at promoters and narrow TSS (Carninci *et al.*,  
13 2006; Hoskins *et al.*, 2011; Lenhard, Sandelin and Carninci, 2012; Haberle *et al.*, 2014). We therefore  
14 tested whether methylated genes were more likely to have tissue-specific or global gene expression  
15 using RNAseq data from different tissue types. In every species with gene body methylation, we  
16 found that methylated genes tended to have less variable expression across different tissues (Figure  
17 6D). Altogether this suggests that across arthropods conserved genes with strongly positioned  
18 nucleosomes, broad TSS and housekeeping functions are targeted for methylation whilst tissue-  
19 specific genes with opposite patterns of nucleosome occupancy and TSS width tend to be depleted of  
20 methylation.

21

## 22 **Discussion**

23 Molecular pathways involved in epigenetic gene regulation evolve surprisingly rapidly and DNA  
24 methylation is no exception. Our work adds to the complex picture of how DNA methylation patterns  
25 change across evolutionary time and offers new insight into potential factors influencing the  
26 distribution of DNA methylation within genomes.

27

28

## 1 **Plasticity of DNA methylation landscapes**

2 Prior to this study, DNA methylation had been characterised across insects (Bewick *et al.*, 2017) but  
3 only isolated species from more basal arthropod clades had been studied (Falckenhayn *et al.*, 2013;  
4 Kao *et al.*, 2016; Kvist *et al.*, 2018; de Mendoza, Pflueger and Lister, 2019; Liu *et al.*, 2019). By  
5 examining a phylogenetically broad range of arthropod methylomes we reconstructed the trajectory of  
6 DNA methylation patterns across the phylum. Our data show that ancestral arthropods likely had  
7 moderate genome-wide methylation including methylation of a small number of transposable  
8 elements. Methylation of genes was also prominent and was enriched in exons over introns; however,  
9 the magnitude of the difference between exonic and intronic methylation was not as striking as in  
10 insects such as *A. mellifera* reflecting the presence of a higher background genomic methylation.  
11 Crucially our data also show that changes in methylation patterns can evolve rapidly within individual  
12 lineages. Most strikingly, we find strong enrichment of TE methylation evolved independently in the  
13 centipede *S. maritima* and the mealybug *P. citri*, which very likely occurred independently. This  
14 enrichment does not correlate to any obvious change in genome structure such as increased TE  
15 proportion or genome size, however it is interesting that a recent paper reported acquisition of a  
16 relatively recent TE family in *S. maritima* that acquires high levels of methylation (de Mendoza,  
17 Pflueger and Lister, 2019), which may underpin gain of TE methylation in that species.

18 It is intriguing that the two species with high TE methylation had independently acquired methylation  
19 of promoters of silent genes, whilst the exons of these genes are devoid of methylation. Gene  
20 regulation by promoter methylation is also found in mammals and was likely acquired independently  
21 in the sponge *Amphimedon queenslandica* (de Mendoza *et al.*, 2019). In all these cases TE  
22 methylation is also prominent so it is possible that the two are linked, perhaps relating to a  
23 requirement to control TE-derived promoter regions; however testing this hypothesis would require  
24 experimental manipulation of methylation in *P. citri* or *S. maritima* which is currently not possible.

25 It is curious that repeated acquisition of similar types of DNA methylation occurs across phylogenies.  
26 This may indicate that targeting of DNA methylation to new regions can be achieved with very few  
27 genetic changes. In vertebrates, a possible example is the KRAB-Zinc finger proteins which can  
28 recruit DNA methylation to TEs through sequence-specific binding (Quenneville *et al.*, 2012). Further

1 work to identify potential “pioneer” factors that recruit DNMTs to specific regions and underlie the  
2 divergence of methylation patterns between species will be of great interest.

### 3 **Potential factors influencing methylation of genes**

4 Our study confirms earlier speculation that the most widely conserved feature of arthropod  
5 methylomes is methylation of genes, biased towards exon methylation (Keller, Han and Yi, 2016).  
6 Additionally, we confirm insights from insects that broadly expressed, housekeeping genes are more  
7 likely to be targeted for methylation than tissue-specific genes (Hunt *et al.*, 2013). This is strikingly  
8 similar to observations in plants and other animal groups, suggesting an ancient evolutionary origin  
9 (Bewick and Schmitz, 2017; Zilberman, 2017). Exactly what the function of this modification is  
10 remains to be elucidated. It is clearly dispensable under some circumstances as, in addition to the  
11 complete loss of DNA methylation in *Drosophila*, we found that DNA methylation at genes has been  
12 lost in both the crustaceans we examined, suggesting that even in species where DNA methylation is  
13 present in the genome, enrichment of DNA methylation at exons is not essential for viability.

14 Whilst we cannot decipher the function of exon-enriched DNA methylation, our analyses potentially  
15 offer new insights into the molecular mechanisms whereby DNA methylation might be deposited. We  
16 identify a remarkable methylation pattern across many arthropods such that methylation levels vary  
17 periodically with the nucleosome-repeat length. This striking genome-wide pattern that we observe in  
18 some species, in particular *S. maritima*, has not been observed to our knowledge in any species  
19 previously. However, there are specific regions within the human genome that display apparently  
20 nucleosome length periodicity in the correlation between adjacent sites (Zhang *et al.*, 2017);  
21 furthermore the influence of nucleosomes on methylation by DNMT3B was observed in human and  
22 yeast cells (Baubec *et al.*, 2015; Morselli *et al.*, 2015). Moreover, DNA methylation levels show a 10bp  
23 periodicity in *Arabidopsis*, corresponding to methylation targeting nucleotides on the same face of the  
24 nucleosome (Chodavarapu *et al.*, 2010). Together these observations reflect a positive correlation  
25 between nucleosome occupancy and DNA methylation in *Arabidopsis* and mammals (Chodavarapu *et al.*  
26 *et al.*, 2010). Exons are known to have better positioned nucleosomes than introns (Schwartz, Meshorer  
27 and Ast, 2009; Tilgner *et al.*, 2009) which might explain why exons are enriched in methylation across  
28 species. We also find that promoters of genes with high levels of methylation tend to carry a clear  
29 nucleosome positioning pattern, typical of housekeeping genes, where nucleosome occupancy is high

1 upstream and just downstream of the TSS with a nucleosome-free region between the two (Lenhard,  
2 Sandelin and Carninci, 2012; Haberle *et al.*, 2014). Both nucleosome positioning and DNA  
3 methylation could be linked to transcription. Since tissue-specific genes are highly expressed in only a  
4 few cell types, this might explain why they do not appear methylated in whole animal bisulphite  
5 sequencing. This would also explain why across all species genes with very low expression are  
6 depleted of methylation (Figure 4D). Alternatively, nucleosomes themselves could dictate where DNA  
7 methylation takes place. Supporting this point there is little periodicity in DNA methylation in introns  
8 compared to exons (Figure S6-2), suggesting that transcription itself is insufficient to account for this  
9 effect.

10 Importantly, the fact that we see these patterns based on nucleosome positioning in *Drosophila* where  
11 DNA is not methylated suggests that nucleosome positioning may cause differences in DNA  
12 methylation. Thus, we suggest that nucleosome positioning may be a primary determinant of variation  
13 in DNA methylation across arthropod genomes. Our analyses may therefore prompt a search for how  
14 nucleosome occupancy might determine methylation patterns across eukaryotes.

15

16

## 17 **Methods**

### 18 **DNMT identification**

19 To identify species that have retained or lost the DNA methylation pathway, we searched for  
20 homologues of DNMT. For each species, we used DIAMOND (Buchfink, Xie and Huson, 2015) to  
21 perform BLASTp searches against all annotated proteins, with *A. mellifera* DNMT1 (NM001171051),  
22 DNMT2 (XM006562945) and DNMT3 (NM001190421) as query sequences. We used InterProScan to  
23 screen out hits that lacked the C-5 cytosine-specific DNA methylase domain, and NCBI BLASTP to  
24 screen out bacterial contaminants (i.e. hits that were more similar to bacterial DNMTs than eukaryotic  
25 DNMTs). To classify DNMTs into subclades (DNMT1, 2 & 3) we aligned all homologues with MAFFT,  
26 screened out badly-aligned regions with Gblocks (Castresana, 2000), and inferred a neighbour-joining  
27 phylogenetic tree under the Jukes-Cantor model using Geneious v10.1.3 (<https://www.geneious.com>).

28

## 1 **Genome annotation**

2 To annotate exons in each genome we used existing annotations, excluding genes that were split  
3 across multiple contigs. To annotate regions which may contain promoters or enhancers, we took  
4 1,000 bases upstream of each gene, excluding genes where this exceeded the contig start or end  
5 point. We annotated introns based on the position of exons, excluding genes that were split across  
6 multiple contigs (using `intron_finder.py` script available at  
7 <https://github.com/SamuelHLewis/BStoolkit/>). To annotate TEs, we used RepeatModeller v1.0.8 to  
8 generate a model of TEs for each genome separately, and then RepeatMasker v4.0.6 to annotate  
9 TEs based on the model for that genome. Within each TE, we used interproscan (Jones *et al.*, 2014)  
10 to search for the following TE-associated domains: PF03184, PF02914, PF13358, PF03732,  
11 PF00665 & PF00077.

12 To annotate rRNA, we either used existing annotations or RNAmmer v1.2 (Lagesen *et al.*, 2007). To  
13 annotate tRNA, we either used existing annotations or tRNAscan-SE v1.3.1 (Lowe and Eddy, 1997).  
14 To avoid ambiguous results caused by overlapping features, we screened out any TE annotations  
15 that overlapped any rRNA, tRNA or exon, and any upstream regions which overlapped any TE, rRNA,  
16 tRNA or exon.

## 17 **Whole genome bisulphite sequencing**

18 To measure DNA methylation on a genome-wide scale, we carried out whole-genome bisulphite  
19 sequencing. We used the DNeasy Blood and Tissue kit (QIAGEN) according to the manufacturer's  
20 protocol to extract DNA from adult somatic tissues of the following species: *L. polyphemus*, *P.*  
21 *tepidariorum*, *S. maritima*, *A. vulgare*, *B. germanica*, *A. pisum*, *B. terrestris*, *N. vespilloides*, *H.*  
22 *melpomene* and *D. melanogaster*. For *I. scapularis*, we used the same method to extract DNA from  
23 the IDE2, IDE8 and ISE18 cell culture. To estimate bisulphite conversion efficiency, we added a  
24 spike-in of unmethylated DNA (P-1025-1, EpiGentek) equal to 0.01% of the sample DNA mass to  
25 each sample. We then prepared whole-genome bisulphite sequencing libraries from each DNA  
26 sample using the Pico Methyl-Seq Library Prep Kit (Zymo Research), according to the manufacturer's  
27 protocol (see Supplementary Table 1 for detailed sample metadata and sequence accession codes).  
28 We sequenced these libraries on an Illumina HiSeq 2500 instrument to generate 100bp paired-end

1 reads. We used pre-existing whole-genome bisulphite sequencing datasets for *P. hawaiiensis*  
2 (SRR3618947, (Kao *et al.*, 2016)) and *A. mellifera* (SRR1790690, (Galbraith *et al.*, 2015)).

3 To generate bisulphite sequencing data for *P. citri*, we extracted DNA from adult males using the  
4 DNeasy Blood and Tissue kit (QIAGEN) according to the manufacturer's protocol. To estimate  
5 bisulphite conversion efficiency, we included a spike-in of non-methylated *Escherichia coli* lambda  
6 DNA (isolated from a heat-inducible lysogenic *E. coli* W3110 strain, provided by Beijing Genomics  
7 Institute (BGI), GenBank/EMBL accession numbers J02459, M17233, M24325, V00636, X00906).  
8 Sequencing of bisulphite libraries was carried out by BGI on an Illumina HiSeq 4000 instrument to  
9 generate 150bp paired-end reads.

## 10 **Bisulphite sequencing data analysis**

11 Before mapping reads to the genome, we trimmed sequencing adapters from each read, and then  
12 trimmed 10 bases from the 5' and 3' end of each read (using the script  
13 <https://github.com/SamuelHLewis/BStoolkit/blob/master/BStrim.sh>). We aligned bisulphite sequencing  
14 reads to each genome using Bismark v0.19.0 (Krueger and Andrews, 2011) in --non\_directional mode  
15 with default settings. We used MethylExtract v1.9.1 (Barturen *et al.*, 2014) to estimate the level of  
16 methylation at each CpG site, calculated as the number of reads in which the cytosine is methylated  
17 divided by the total number of reads covering the cytosine, excluding sites covered by fewer than 10  
18 reads on each strand. Due to the large number of contigs in their genome assemblies exceeding the  
19 memory limit for MethylExtract, we split the genomes of *I. scapularis*, *L. polyphemus* and *P.*  
20 *hawaiiensis* into individual contigs, ran MethylExtract on each contig separately, and concatenated the  
21 resulting output files into one file for each genome.

22 To estimate the genome-wide background level of CpG methylation, we calculated the mean  
23 methylation for all CpGs outside annotated features (exon, intron, upstream region, TE, rRNA &  
24 tRNA). To gain an accurate estimate of the methylation level of each feature, we calculated the mean  
25 methylation level of all CpGs within that feature, excluding any feature with fewer than 3 sufficiently-  
26 covered CpGs (only CpGs covered by >10 reads are analysed). We estimated 95% confidence  
27 intervals for the mean methylation of genes and TEs within each species using 1000 nonparametric  
28 bootstrap replicates (i.e. genes or TEs were resampled with replacement 1000 times to generate an  
29 empirical distribution of the mean).



## 1 **Phylogenetics and ancestral state reconstruction**

2 To infer the ancestral levels of genome-wide methylation across 29 species of arthropods with newly-  
3 produced or publicly-available methylation data (Figure 1), we obtained a time-scaled species tree  
4 from TimeTree ([www.timetree.org](http://www.timetree.org), accessed 12.03.2019). We then used a maximum-likelihood  
5 approach to infer the genome-wide methylation level at all internal nodes of this tree based on the  
6 levels at the tips, using the fastAnc function within phytools (Revell, 2012).

7 To infer the ancestral levels of gene-body and TE methylation for the 14 focal species, we constructed  
8 a Bayesian time-scaled species tree for 14 focal species (Figure 2 & 3). We first identified 236  
9 proteins present as 1:1:1 orthologues across our species set, concatenated the protein sequences  
10 together, and aligned them using MAFFT v7.271 (Katoh and Standley, 2013) with default settings. We  
11 then screened out poorly-aligned regions using Gblocks (Castresana, 2000) with least stringent  
12 settings. Using this alignment, we constructed a phylogenetic tree using BEAST v1.8.4 (Drummond *et*  
13 *al.*, 2012) to infer branch lengths. We specified a strict molecular clock, gamma-distributed rate  
14 variation, no invariant sites, and a birth-death speciation process. We fixed the topology and set prior  
15 distributions on key internal node dates (Arthropoda  $\tau = 568 \pm 29$ , Insecta–  
16 Crustacea  $\tau = 555 \pm 33$ , Insecta  $\tau = 386 \pm 27$ , Hymenoptera–Coleoptera–Lepidoptera–  
17 Diptera  $\tau = 345 \pm 27$ , Coleoptera–Lepidoptera–Diptera  $\tau = 327 \pm 26$ ), deriving these values from  
18 an existing phylogenetic analysis of arthropods (Misof *et al.*, 2014). We ran the analysis for 10 million  
19 generations, and used TreeAnnotator (Drummond *et al.*, 2012) to generate a maximum clade  
20 credibility tree. We then used a maximum-likelihood approach to infer the gene-body and TE  
21 methylation levels (separately) at all internal nodes of this tree, using the fastAnc function within  
22 phytools (Revell, 2012).

23 To test whether genome-wide methylation levels differ between species with and without ALKB2, we  
24 fitted a phylogenetic mixed model using MCMCglmm (Hadfield, 2010). To account for phylogenetic  
25 non-independence caused by sampling species with different levels of relatedness, we used the  
26 branch lengths of the time-scaled (ultrametric) species tree (see above) to calculate a genetic  
27 distance matrix, and included this in the model as a random factor. We ran the analysis for 6 million  
28 iterations, with a burn-in of 1 million iterations and thinning of 500 generations.

29

1

## 2 **RNA-Seq data analysis**

3 To investigate the link between DNA methylation and transcription, we used RNA-Seq data generated  
4 previously for arthropod somatic tissue (NCBI PRJNA386859, (Lewis *et al.*, 2018) and the *I.*  
5 *scapularis* IDE-8 cell line (SRR1756347, Arthropod Cell Line RNA Seq initiative, Broad Institute,  
6 broadinstitute.org). To measure the expression of each feature, we trimmed adaptors and low-quality  
7 ends using Trim Galore with default settings, and mapped RNA-Seq reads to the genome of each  
8 species using TopHat2 v2.1.1 (Kim *et al.*, 2013) with default settings for strand-specific libraries (--  
9 library-type fr-firststrand mode). We counted the number of reads overlapping each feature using  
10 BEDTools coverage v2.25.0 in strand-specific mode, and divided the number of reads by the feature  
11 length to generate expression level estimates in fragments per per kilobase million (FPKM).

12 To test whether variation in tissue-specific expression differs between highly- and lowly-methylated  
13 genes, we calculated the coefficient of variation for expression of each gene in each species with  
14 RNA-Seq data (i.e. excluding *B. germanica*, *I. scapularis* & *P. hawaiiensis*). For *S. maritima* we used  
15 RNA-Seq data for fat body and nerve chord; for *P. citri* & *A. pisum* we used RNA-Seq data for female  
16 soma and germline; and for all other species we used RNA-Seq data for female and male soma and  
17 germline.

## 18 **Periodic correlation in methylation levels**

19 To obtain an estimate of how the correlation between the methylation levels of sites varied with  
20 distance between the sites, we collected all pairs of sites separated by  $d$  nucleotides where  $d$  could  
21 vary between 3 and 500 nucleotides within the same exon. For each separate  $d$  we then computed  
22 the correlation coefficient across all the pairs. To quantify the periodic component of the signal we  
23 subtracted any gradual change in correlation across the entire window by calculating the residuals of  
24 a linear model. This signal was subjected to Fourier analysis using the fast Fourier transform  
25 algorithm implemented in R. A linear model was used to subtract the baseline across the 500bp and  
26 the residuals were used as a time series for input into the algorithm, with 50000 0 values ended on to  
27 the end of the series to increase the resolution of the algorithm. The total intensity of the components  
28 between 140 and 200 base pairs was calculated to give the nucleosome periodicity for each species.

## 1 **Nucleosome positioning analysis**

2 The genomic coordinates of the *D. melanogaster* members of orthogroups conserved across all  
3 species were extracted and the top 20% (high methylation) and bottom 20% methylation (low  
4 methylation) levels selected. Nucleosome positioning data from the *D. melanogaster* S2 cell line was  
5 downloaded from Modencode (Ho *et al.*, 2014). The average signal was computed across 200bp  
6 windows spanning 2kb either side of the annotated transcription start site for each gene. The mean  
7 signal was computed within the high methylation and low methylation sets separately and a loess fit  
8 performed. To obtain confidence intervals, the mean signal was computed on 100 random samples  
9 containing 90% of the data and a loess fit calculated on the lowest and highest values obtained for  
10 each 200bp window.

## 11 **CAGE data analysis**

12 Total body RNA was extracted from L3 *Drosophila melanogaster* ( $w^{1118}$ ) larvae using the Qiagen  
13 RNeasy kit. CAGE library preparation was performed using the nAnT-iCAGE protocol (Murata *et al.*,  
14 2014). Two biological replicates were prepared from 5 ug of total RNA each. The libraries were  
15 sequenced in single-end 50 bp-pair mode. CAGE tags (47 bp) were mapped to the reference *D.*  
16 *melanogaster* genome (assembly Release 6) using Bowtie2 (Langmead and Salzberg, 2012) with  
17 default parameters. Uniquely mapped reads were imported into R (<http://www.R-project.org/>) as bam  
18 files using the standard workflow within the CAGER package (Haberle *et al.*, 2015). The 5' ends of  
19 reads are CAGE-supported transcription start sites (CTSSs) and the number of tags for each CTSS  
20 reflects expression levels. Raw tags were normalised using a referent power-law distribution and  
21 expressed as normalized tags per million (TPMs). Biological replicates were highly correlated ( $r^2 =$   
22 0.99) and were therefore merged prior to downstream analyses using standard Bioconductor  
23 packages (<http://www.bioconductor.org/>) and custom scripts.

24 CTSSs were clustered together into tag clusters, a single functional transcriptional unit, using  
25 distance-based clustering, with the maximum distance allowed between adjacent CTSSs being 20 bp.  
26 For each tag cluster, the interquartile width was calculated as the distance between CTSSs at the  
27 10<sup>th</sup> and 90<sup>th</sup> quartile of the cumulative distribution of expression across the cluster. The interquartile  
28 range of each gene within the top 20% and bottom 20% of methylation levels was extracted and  
29 compared.

1

## 2 **Availability of scripts and data**

3 Sequence data that was newly-generated for this project have been deposited in the NCBI Short  
4 Read Archive under the BioProject accession code PRJNA589724. The source code, input data and  
5 newly-identified DNMT & ALKB2 gene sequences are available from the Cambridge Data Repository  
6 (<https://doi.org/10.17863/CAM.45964>).

7

## 8 **Acknowledgements**

9 We thank L. Bell-Sakyi, the Tick Cell Biobank, A. McGregor, R. Jenner, M. Akam, A. McLean, D.  
10 Collins, R. Kilner, A. Pinharanda and C. Jiggins for providing arthropod samples. This research was  
11 supported by a Leverhulme Research Project Grant (RPG-2016-210 to F.M.J., E.A.M. and P.S.) and  
12 the Medical Research Council (MC-A652-5PY80). Sequencing of bisulphite and CAGE libraries was  
13 carried out by the LMS Genomics Facility.

14

15

16

17

18

19

20

21

22

23

24

1

2 **References**

- 3 Barturen, G. *et al.* (2014) 'MethylExtract: High-Quality methylation maps and SNV calling from  
4 whole genome bisulfite sequencing data', *F1000Research*. F1000 Research Limited, 2, p. 217.  
5 doi: 10.12688/f1000research.2-217.v2.
- 6 Baubec, T. *et al.* (2015) 'Genomic profiling of DNA methyltransferases reveals a role for  
7 DNMT3B in genic methylation', *Nature*, 520(7546), pp. 243–247. doi: 10.1038/nature14176.
- 8 Bewick, A. J. *et al.* (2017) 'Evolution of DNA methylation across insects', *Molecular Biology and  
9 Evolution*, 34(3), pp. 654–665. doi: 10.1093/molbev/msw264.
- 10 Bewick, A. J. *et al.* (2019) 'Diversity of cytosine methylation across the fungal tree of life', *Nature  
11 Ecology & Evolution*. Nature Publishing Group, 3(3), pp. 479–490. doi: 10.1038/s41559-019-  
12 0810-9.
- 13 Bewick, A. J. and Schmitz, R. J. (2017) 'Gene body DNA methylation in plants', *Current Opinion in  
14 Plant Biology*. doi: 10.1016/j.pbi.2016.12.007.
- 15 Bird, A. (2002) 'DNA methylation patterns and epigenetic memory.', *Genes & development*. Cold  
16 Spring Harbor Laboratory Press, 16(1), pp. 6–21. doi: 10.1101/gad.947102.
- 17 Bonasio, R. *et al.* (2012) 'Genome-wide and caste-specific DNA methylomes of the ants  
18 camponotus floridanus and harpegnathos saltator', *Current Biology*, 22(19), pp. 1755–1764. doi:  
19 10.1016/j.cub.2012.07.042.
- 20 Buchfink, B., Xie, C. and Huson, D. H. (2015) 'Fast and sensitive protein alignment using  
21 DIAMOND', *Nature Methods*. Nature Publishing Group, 12(1), pp. 59–60. doi:  
22 10.1038/nmeth.3176.
- 23 Carninci, P. *et al.* (2006) 'Genome-wide analysis of mammalian promoter architecture and  
24 evolution', *Nature Genetics*. Nature Publishing Group, 38(6), pp. 626–635. doi: 10.1038/ng1789.
- 25 Casadesús, J. and Low, D. (2006) 'Epigenetic gene regulation in the bacterial world.',  
26 *Microbiology and molecular biology reviews*: MMBR. American Society for Microbiology, 70(3),  
27 pp. 830–56. doi: 10.1128/MMBR.00016-06.
- 28 Castresana, J. (2000) 'Selection of Conserved Blocks from Multiple Alignments for Their Use in  
29 Phylogenetic Analysis', *Molecular Biology and Evolution*. Narnia, 17(4), pp. 540–552. doi:  
30 10.1093/oxfordjournals.molbev.a026334.
- 31 Chodavarapu, R. K. *et al.* (2010) 'Relationship between nucleosome positioning and DNA  
32 methylation', *Nature*. Nature Publishing Group, 466(7304), pp. 388–392. doi:  
33 10.1038/nature09147.
- 34 Drummond, A. J. *et al.* (2012) 'Bayesian Phylogenetics with BEAUti and the BEAST 1.7',  
35 *Molecular Biology and Evolution*. Narnia, 29(8), pp. 1969–1973. doi: 10.1093/molbev/mss075.
- 36 Falckenhayn, C. *et al.* (2013) 'Characterization of genome methylation patterns in the desert  
37 locust *Schistocerca gregaria*.', *The Journal of experimental biology*. The Company of Biologists  
38 Ltd, 216(Pt 8), pp. 1423–9. doi: 10.1242/jeb.080754.
- 39 Feng, S. *et al.* (2010) 'Conservation and divergence of methylation patterning in plants and  
40 animals', *Proceedings of the National Academy of Sciences*, 107(19), pp. 8689–8694. doi:  
41 10.1073/pnas.1002720107.

- 1 Galbraith, D. A. *et al.* (2015) 'Parallel Epigenomic and Transcriptomic Responses to Viral  
2 Infection in Honey Bees (*Apis mellifera*)', *PLoS Pathogens*. Edited by D. S. Schneider. Public  
3 Library of Science, 11(3), p. e1004713. doi: 10.1371/journal.ppat.1004713.
- 4 Goll, M. G. *et al.* (2006) 'Methylation of tRNA<sup>Asp</sup> by the DNA methyltransferase homolog  
5 Dnmt2.', *Science (New York, N.Y.)*. American Association for the Advancement of Science,  
6 311(5759), pp. 395–8. doi: 10.1126/science.1120976.
- 7 Haberle, V. *et al.* (2014) 'Two independent transcription initiation codes overlap on vertebrate  
8 core promoters', *Nature*. Nature Publishing Group, 507(7492), pp. 381–385. doi:  
9 10.1038/nature12974.
- 10 Haberle, V. *et al.* (2015) 'CAGER: Precise TSS data retrieval and high-resolution promoterome  
11 mining for integrative analyses', *Nucleic Acids Research*. doi: 10.1093/nar/gkv054.
- 12 Haberle, V. and Lenhard, B. (2016) 'Promoter architectures and developmental gene regulation',  
13 *Seminars in Cell and Developmental Biology*. doi: 10.1016/j.semcdb.2016.01.014.
- 14 Hadfield, J. D. (2010) 'MCMC Methods for Multi-Response Generalized Linear Mixed Models: The  
15 **MCMCglmm** R Package', *Journal of Statistical Software*, 33(2), pp. 1–22. doi:  
16 10.18637/jss.v033.i02.
- 17 Ho, J. W. K. *et al.* (2014) 'Comparative analysis of metazoan chromatin organization', *Nature*.  
18 Nature Publishing Group, 512(7515), pp. 449–452. doi: 10.1038/nature13415.
- 19 Holliday, R. *et al.* (1987) 'The inheritance of epigenetic defects.', *Science (New York, N.Y.)*.  
20 American Association for the Advancement of Science, 238(4824), pp. 163–70. doi:  
21 10.1126/science.3310230.
- 22 Holliday, R. (2006) 'Epigenetics: A historical overview', *Epigenetics*, pp. 76–80. doi:  
23 10.4161/epi.1.2.2762.
- 24 Hoskins, R. A. *et al.* (2011) 'Genome-wide analysis of promoter architecture in *Drosophila*  
25 *melanogaster*.', *Genome research*. Cold Spring Harbor Laboratory Press, 21(2), pp. 182–92. doi:  
26 10.1101/gr.112466.110.
- 27 Hunt, B. G. *et al.* (2013) 'Patterning and Regulatory Associations of DNA Methylation Are  
28 Mirrored by Histone Modifications in Insects', *Genome Biology and Evolution*. Narnia, 5(3), pp.  
29 591–598. doi: 10.1093/gbe/evt030.
- 30 Jeltsch, A. (2006) 'Molecular enzymology of mammalian DNA methyltransferases', in *Current*  
31 *Topics in Microbiology and Immunology*, pp. 203–225. doi: 10.1007/3-540-31390-7\_7.
- 32 Jones, P. *et al.* (2014) 'InterProScan 5: genome-scale protein function classification',  
33 *Bioinformatics*. Narnia, 30(9), pp. 1236–1240. doi: 10.1093/bioinformatics/btu031.
- 34 Jurkowski, T. P. and Jeltsch, A. (2011) 'On the evolutionary origin of eukaryotic DNA  
35 methyltransferases and Dnmt2', *PLoS ONE*, 6(11). doi: 10.1371/journal.pone.0028104.
- 36 Kao, D. *et al.* (2016) 'The genome of the crustacean parhyale hawaiiensis, a model for animal  
37 development, regeneration, immunity and lignocellulose digestion', *eLife*, 5(November 2016).  
38 doi: 10.7554/eLife.20062.001.
- 39 Katoh, K. and Standley, D. M. (2013) 'MAFFT multiple sequence alignment software version 7:  
40 improvements in performance and usability', *Mol Biol Evol*, 30. doi: 10.1093/molbev/mst010.
- 41 Keller, T. E., Han, P. and Yi, S. V. (2016) 'Evolutionary transition of promoter and gene body DNA  
42 methylation across invertebrate-vertebrate boundary', *Molecular Biology and Evolution*, 33(4),  
43 pp. 1019–1028. doi: 10.1093/molbev/msv345.

- 1 Kim, D. *et al.* (2013) 'TopHat2: accurate alignment of transcriptomes in the presence of  
2 insertions, deletions and gene fusions', *Genome Biology*. BioMed Central, 14(4), p. R36. doi:  
3 10.1186/gb-2013-14-4-r36.
- 4 Krueger, F. and Andrews, S. R. (2011) 'Bismark: a flexible aligner and methylation caller for  
5 Bisulfite-Seq applications', *Bioinformatics*. Narnia, 27(11), pp. 1571–1572. doi:  
6 10.1093/bioinformatics/btr167.
- 7 Kvist, J. *et al.* (2018) 'Pattern of DNA Methylation in *Daphnia*: Evolutionary Perspective', *Genome  
8 Biology and Evolution*. Narnia, 10(8), pp. 1988–2007. doi: 10.1093/gbe/evy155.
- 9 Lagesen, K. *et al.* (2007) 'RNAmmer: consistent and rapid annotation of ribosomal RNA genes',  
10 *Nucleic Acids Research*. Narnia, 35(9), pp. 3100–3108. doi: 10.1093/nar/gkm160.
- 11 Langmead, B. and Salzberg, S. L. (2012) 'Fast gapped-read alignment with Bowtie 2', *Nat Meth.*  
12 Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved., 9(4),  
13 pp. 357–359. Available at: <http://dx.doi.org/10.1038/nmeth.1923>.
- 14 Law, J. A. and Jacobsen, S. E. (2010) 'Establishing, maintaining and modifying DNA methylation  
15 patterns in plants and animals', *Nature Reviews Genetics*, pp. 204–220. doi: 10.1038/nrg2719.
- 16 Lenhard, B., Sandelin, A. and Carninci, P. (2012) 'Metazoan promoters: emerging characteristics  
17 and insights into transcriptional regulation', *Nature Reviews Genetics*. Nature Publishing Group,  
18 13(4), pp. 233–245. doi: 10.1038/nrg3163.
- 19 Lewis, S. H. *et al.* (2018) 'Pan-arthropod analysis reveals somatic piRNAs as an ancestral defence  
20 against transposable elements', *Nature Ecology and Evolution*, 2(1), pp. 174–181. doi:  
21 10.1038/s41559-017-0403-4.
- 22 Liu, S. *et al.* (2019) 'DNA Methylation Patterns in the Social Spider, *Stegodyphus dumicola*',  
23 *Genes*. Multidisciplinary Digital Publishing Institute, 10(2), p. 137. doi:  
24 10.3390/genes10020137.
- 25 Lowe, T. M. and Eddy, S. R. (1997) 'tRNAscan-SE: A Program for Improved Detection of Transfer  
26 RNA Genes in Genomic Sequence', *Nucleic Acids Research*. Narnia, 25(5), pp. 955–964. doi:  
27 10.1093/nar/25.5.955.
- 28 Lyko, F. *et al.* (2010) 'The Honey Bee Epigenomes: Differential Methylation of Brain DNA in  
29 Queens and Workers', *PLoS Biology*. Edited by L. Keller. Public Library of Science, 8(11), p.  
30 e1000506. doi: 10.1371/journal.pbio.1000506.
- 31 de Mendoza, A. *et al.* (2019) 'Convergent evolution of a vertebrate-like methylome in a marine  
32 sponge', *Nature Ecology & Evolution*. Nature Publishing Group, 3(10), pp. 1464–1473. doi:  
33 10.1038/s41559-019-0983-2.
- 34 de Mendoza, A., Pflueger, J. and Lister, R. (2019) 'Capture of a functionally active methyl-CpG  
35 binding domain by an arthropod retrotransposon family', *Genome Research*, 29(8), pp. 1277–  
36 1286. doi: 10.1101/gr.243774.118.
- 37 Misof, B. *et al.* (2014) 'Phylogenomics resolves the timing and pattern of insect evolution.',  
38 *Science (New York, N.Y.)*. American Association for the Advancement of Science, 346(6210), pp.  
39 763–7. doi: 10.1126/science.1257570.
- 40 Morselli, M. *et al.* (2015) 'In vivo targeting of de novo DNA methylation by histone modifications  
41 in yeast and mouse', *eLife*, 4. doi: 10.7554/eLife.06205.
- 42 Murata, M. *et al.* (2014) 'Detecting expressed genes using CAGE', *Methods in Molecular Biology*.  
43 doi: 10.1007/978-1-4939-0805-9\_7.
- 44 Nashun, B., Hill, P. W. S. and Hajkova, P. (2015) 'Reprogramming of cell fate: epigenetic memory

- 1 and the erasure of memories past.', *The EMBO journal*. EMBO Press, 34(10), pp. 1296–308. doi:  
2 10.15252/embj.201490649.
- 3 Ponger, L. and Li, W. H. (2005) 'Evolutionary diversification of DNA methyltransferases in  
4 eukaryotic genomes', *Molecular Biology and Evolution*, 22(4), pp. 1119–1128. doi:  
5 10.1093/molbev/msi098.
- 6 Quenneville, S. *et al.* (2012) 'The KRAB-ZFP/KAP1 system contributes to the early embryonic  
7 establishment of site-specific DNA methylation patterns maintained during development.', *Cell*  
8 *reports*. Elsevier, 2(4), pp. 766–73. doi: 10.1016/j.celrep.2012.08.043.
- 9 Revell, L. J. (2012) 'phytools: an R package for phylogenetic comparative biology (and other  
10 things)', *Methods in Ecology and Evolution*. John Wiley & Sons, Ltd (10.1111), 3(2), pp. 217–223.  
11 doi: 10.1111/j.2041-210X.2011.00169.x.
- 12 Rošić, S. *et al.* (2018) 'Evolutionary analysis indicates that DNA alkylation damage is a  
13 byproduct of cytosine DNA methyltransferase activity', *Nature Genetics*, 50(3), pp. 452–459. doi:  
14 10.1038/s41588-018-0061-8.
- 15 Schwartz, S., Meshorer, E. and Ast, G. (2009) 'Chromatin organization marks exon-intron  
16 structure', *Nature Structural & Molecular Biology*. Nature Publishing Group, 16(9), pp. 990–995.  
17 doi: 10.1038/nsmb.1659.
- 18 Suzuki, M. M. and Bird, A. (2008) 'DNA methylation landscapes: provocative insights from  
19 epigenomics', *Nature Reviews Genetics*. Nature Publishing Group, 9(6), pp. 465–476. doi:  
20 10.1038/nrg2341.
- 21 Tilgner, H. *et al.* (2009) 'Nucleosome positioning as a determinant of exon recognition', *Nature*  
22 *Structural & Molecular Biology*. Nature Publishing Group, 16(9), pp. 996–1001. doi:  
23 10.1038/nsmb.1658.
- 24 Walsh, C. P., Chaillet, J. R. and Bestor, T. H. (1998) 'Transcription of IAP endogenous retroviruses  
25 is constrained by cytosine methylation', *Nature Genetics*. Nature Publishing Group, 20(2), pp.  
26 116–117. doi: 10.1038/2413.
- 27 Wang, X. *et al.* (2013) 'Function and Evolution of DNA Methylation in *Nasonia vitripennis*', *PLoS*  
28 *Genetics*, 9(10). doi: 10.1371/journal.pgen.1003872.
- 29 Wang, Y. *et al.* (2006) 'Functional CpG methylation system in a social insect', *Science*, 314(5799),  
30 pp. 645–647. doi: 10.1126/science.1135213.
- 31 Wu, P. *et al.* (2017) 'DNA methylation in silkworm genome may provide insights into epigenetic  
32 regulation of response to *Bombyx mori* cytopovirus infection', *Scientific Reports*. Nature  
33 Publishing Group, 7(1), p. 16013. doi: 10.1038/s41598-017-16357-7.
- 34 Xiang H. *et al.* (2010) 'Single base-resolution methylome of the silkworm reveals a sparse  
35 epigenomic map', *Nature Biotechnology*. Nature Publishing Group, 28(5), pp. 516–520. doi:  
36 10.1038/nbt.1626.
- 37 Zemach, A. *et al.* (2010) 'Genome-wide evolutionary analysis of eukaryotic DNA methylation',  
38 *Science*, 328(5980), pp. 916–919. doi: 10.1126/science.1186366.
- 39 Zhang, L. *et al.* (2017) 'DNA Methylation Landscape Reflects the Spatial Organization of  
40 Chromatin in Different Cells.', *Biophysical journal*. Elsevier, 113(7), pp. 1395–1404. doi:  
41 10.1016/j.bpj.2017.08.019.
- 42 Zilberman, D. (2017) 'An evolutionary case for functional gene body methylation in plants and  
43 animals', *Genome Biology*. BioMed Central, 18(1), p. 87. doi: 10.1186/s13059-017-1230-2.
- 44 ☐





1 **Figure Legends**

2 **Figure 1. Genome-wide CpG methylation across the arthropod phylogeny.** (A) A phylogeny of 29  
3 arthropod species that have publicly available or newly computed genome-wide methylation  
4 estimates, with branches coloured to show an ancestral state reconstruction of the percentage of CpG  
5 sites that are methylated in the genome. (B) The percentage of CpG sites that are methylated  
6 genome-wide. (C) The number of *DNMT* and *ALKB2* homologues in the genomes of each species.

7 **Figure 2. Methylation of transposable elements.** For 14 diverse arthropod species with annotated  
8 genomes, we explored methylation characteristics of genomic features. (A) Density plot of the mean  
9 % CpG methylation per gene and per TE. (B) Ancestral state reconstruction of the mean %  
10 methylation of CpGs within TEs. (C) Mean % methylation of CpGs within TEs with 95% bootstrap  
11 confidence intervals. Red points are CpGs >1kB from annotated regions of the genome. (D)  
12 Metagene plot of methylation within TEs (pink) and in flanking sequence for *S. maritima* and *P. citri*.

13 **Figure 3 Gene body methylation.** (A) Ancestral state reconstruction of the mean % methylation of  
14 CpGs within exons. (B) Mean % methylation of CpGs within exons with 95% bootstrap confidence  
15 intervals. Red points are CpGs >1kB from annotated regions of the genome. (C) Metagene plot of  
16 methylation across introns (white), exons (pink), UTRs (blue) and 1kB of flanking sequence (white).

17 **Figure 4 Promoter methylation.** (A) Methylation across upstream regions for highly expressed  
18 genes (top 20%) and lowly expressed genes (bottom 20%). *P. hawaiiensis* is omitted due to lack of  
19 gene expression data. Expression of genes across bins of decreasing upstream methylation in *S.*  
20 *maritima* (B) and *P. citri* (C).

21 **Figure 5 The expression and conservation of methylated genes.** (A) Methylation of orthologous  
22 genes in different species. Only genes with orthologs in all species are shown, and in species with  
23 multiple paralogs the mean % CpG methylation is shown. Genes are ranked by their mean  
24 methylation. (B) Histogram of gene expression estimated from RNAseq data for methylated and  
25 unmethylated genes in *L. polyphemous* (FPKM: fragments per kilobase million). (C) The relationship  
26 between gene expression and CpG methylation for genes that are conserved across all species and  
27 species-specific genes. To combine data across species, the methylation rate was normalised by  
28 taking the Z-score of methylation and expression of each gene within each species. Each point is a  
29 gene from a single species, and the colour represents the density of overlaid points.

1 **Figure 6. Nucleosome occupancy and DNA methylation.** The Pearson correlation coefficient in  
2 DNA methylation levels between pairs of CpG at different distances apart in (A) *S. maritima* and (B)  
3 *H. melpomene*. (C) Nucleosome occupancy in *D. melanogaster* orthologues of genes that are either  
4 highly methylated (grey) or unmethylated (red) in arthropods. Shaded area is a 95% bootstrap  
5 confidence interval. Across all species in the dataset, mean methylation levels were estimated for  
6 each group of orthologous genes using a general linear mixed model. The top and bottom 20% were  
7 classified as methylated and unmethylated respectively. Only genes with orthologs in all species are  
8 shown. (D) Interquartile range of the TSS window for the *D. melanogaster* orthologues of highly  
9 methylated orthogroups (top 20%) and lowly methylated orthogroups (bottom 20%). (E) The  
10 coefficient of variation in expression of genes with high (top 20%) and low (bottom 20%) methylation  
11 across different tissues estimated using RNAseq data. *P. hawaiiensis* is omitted because no tissue-  
12 specific data is available for this species.

13 **Figure 6 supplement 1: estimation of nucleosomal periodicity signal for methylation of exons**  
14 **and transposable elements across all species.** (A) workflow using *S. maritima* exons as an  
15 example for how baseline correction and fast fourier transform were used to obtain a nucleosome  
16 signal. (C) Nucleosome signal as a fraction of total signal for exonic methylation across arthropods.  
17 (D) Nucleosome signal as a fraction of total signal for TE methylation across arthropods.

18

19

20

21

22

23

24

25

26

1 **Supplementary Information**

2 **Figure 1-Supplement: ALKB2 DNA repair is associated with high levels of DNA methylation**  
3 **across arthropods**

4 Boxplot showing genome-wide methylation levels in 29 arthropod species with and without ALKB2.

5 **Figure 2 supplement: Metagene plot of methylation within TEs and in flanking sequence for all**  
6 **species.** TEs are shown in pink, flanking sequence in white.

7 **Figure 5 supplement 1: Expression patterns of methylated and unmethylated genes for all**  
8 **species (cf Figure 5B)**

9 **Figure 5 supplement 2: Methylation of single exon and multi-exon genes for all species in**  
10 **which we see gene body methylation**

11 **Figure 6 Supplement 2:** Intron periodicity is markedly less apparent than exon periodicity. *S.*  
12 *maritima* exons 1 to 4 (A) and introns 1 to 4 (B) are shown for comparison.

13

14 **Supplementary Table 1:** Details of the tissue type, sex, caste, BioSample Accession and SRA  
15 Accession of each sample that was newly-sequenced in this study.

16

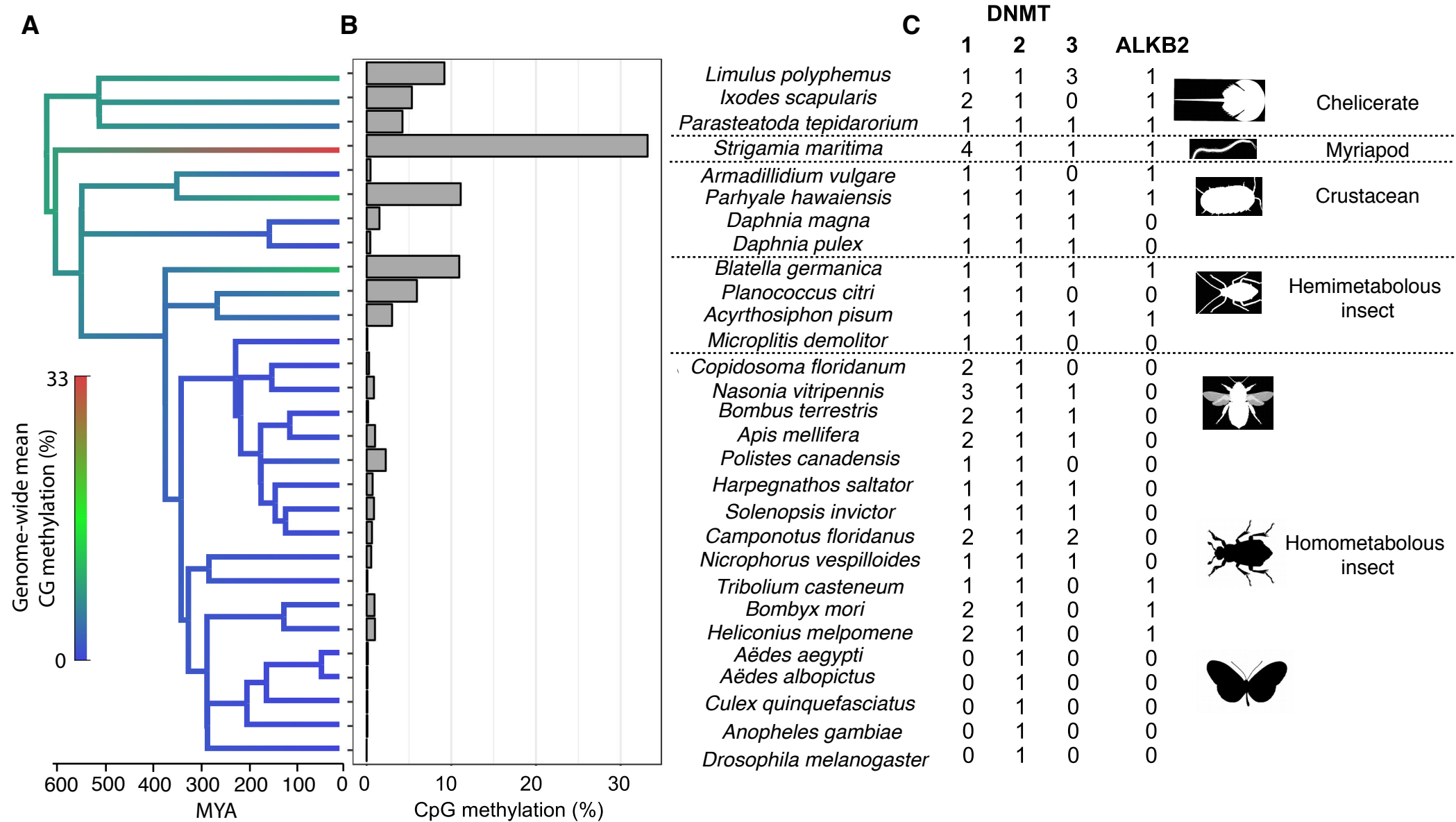


Figure 1

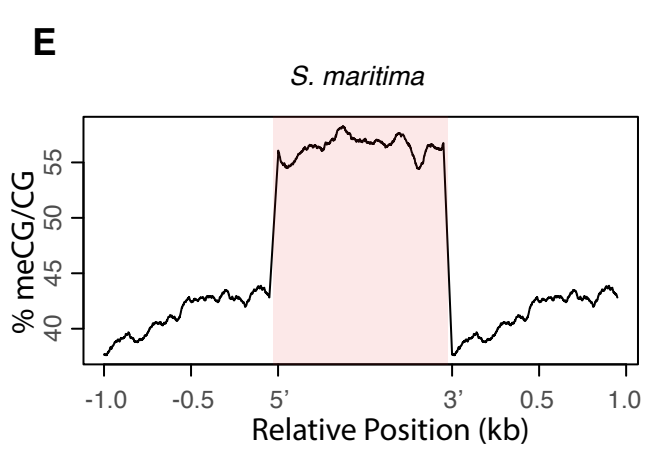
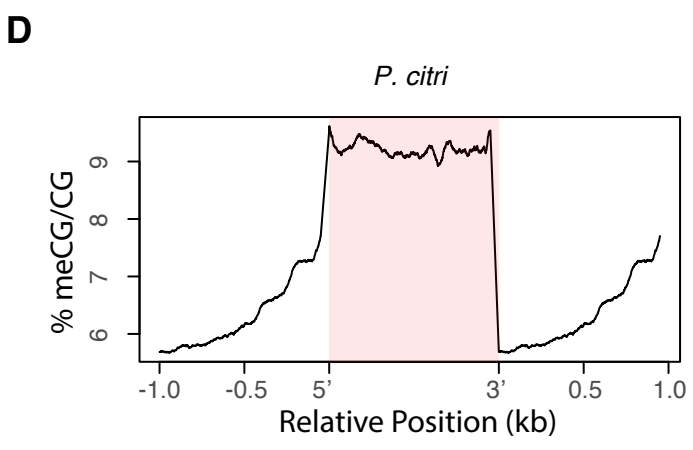
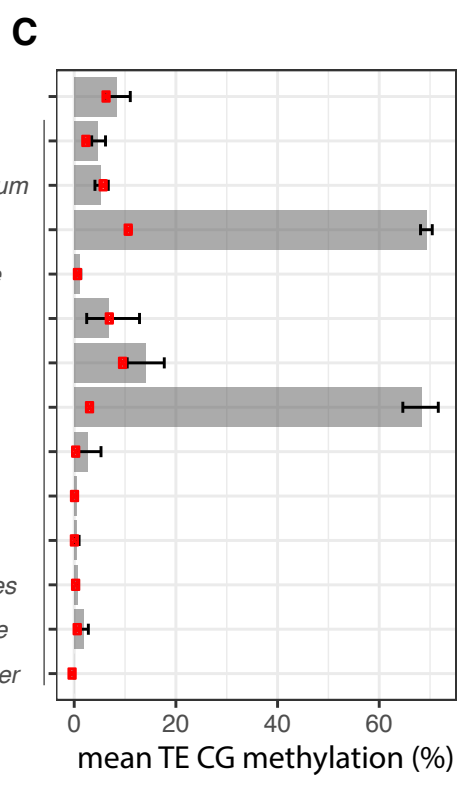
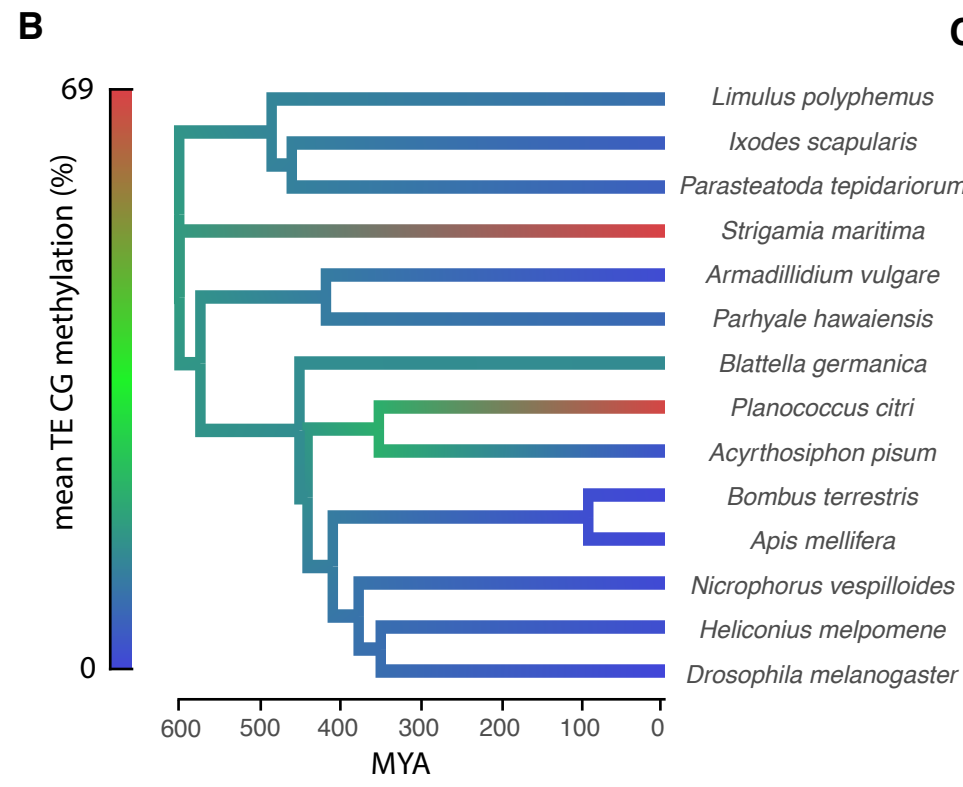
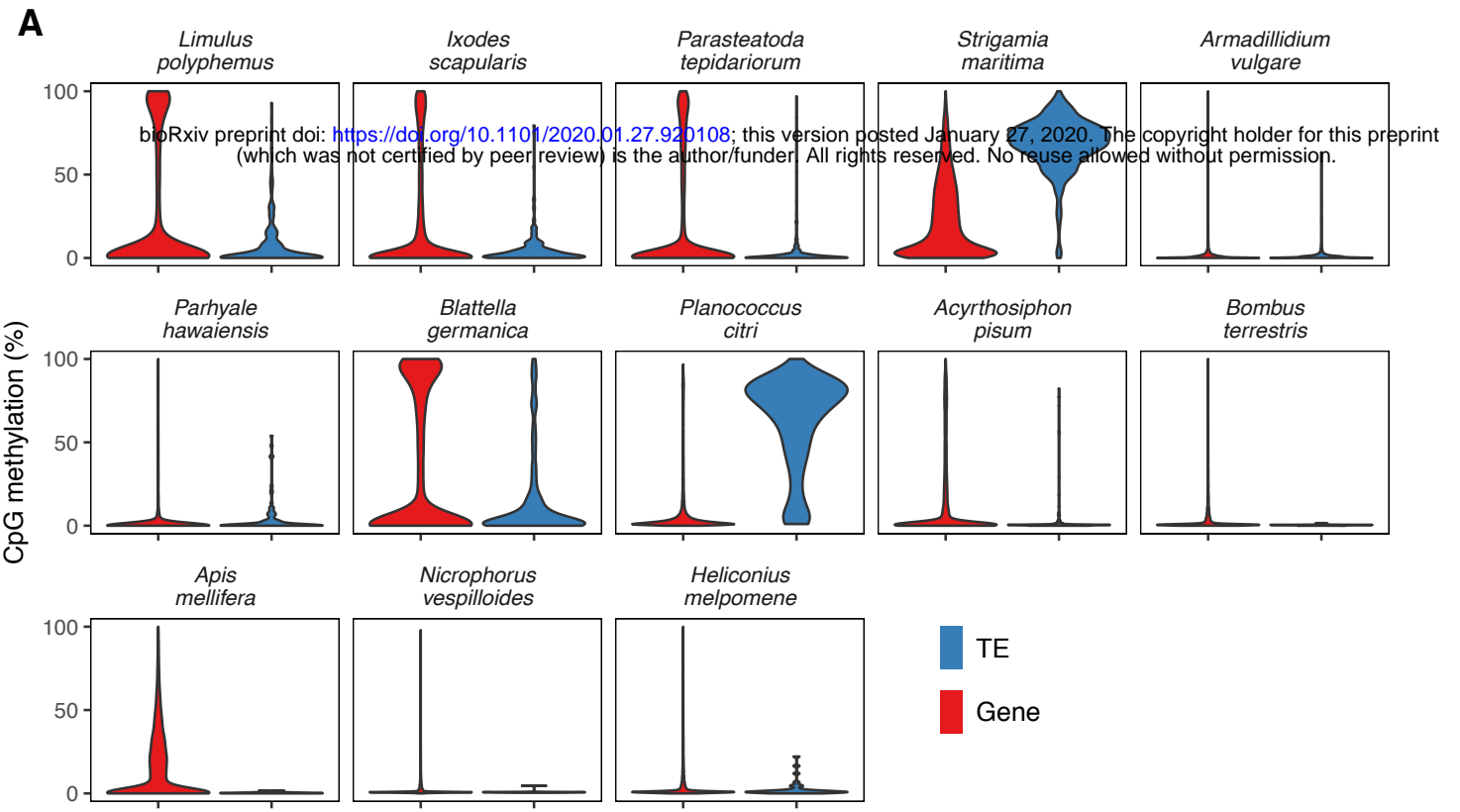


Figure 2

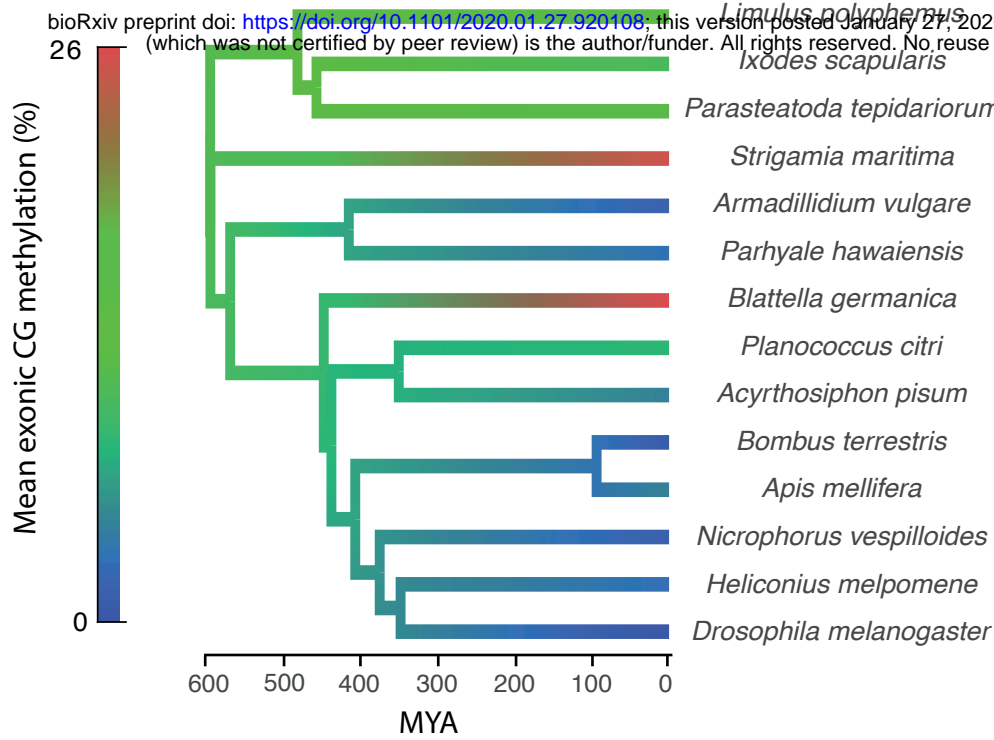
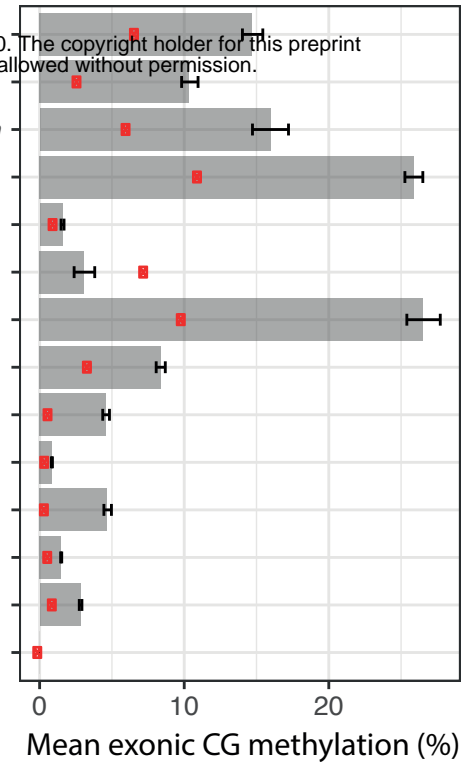
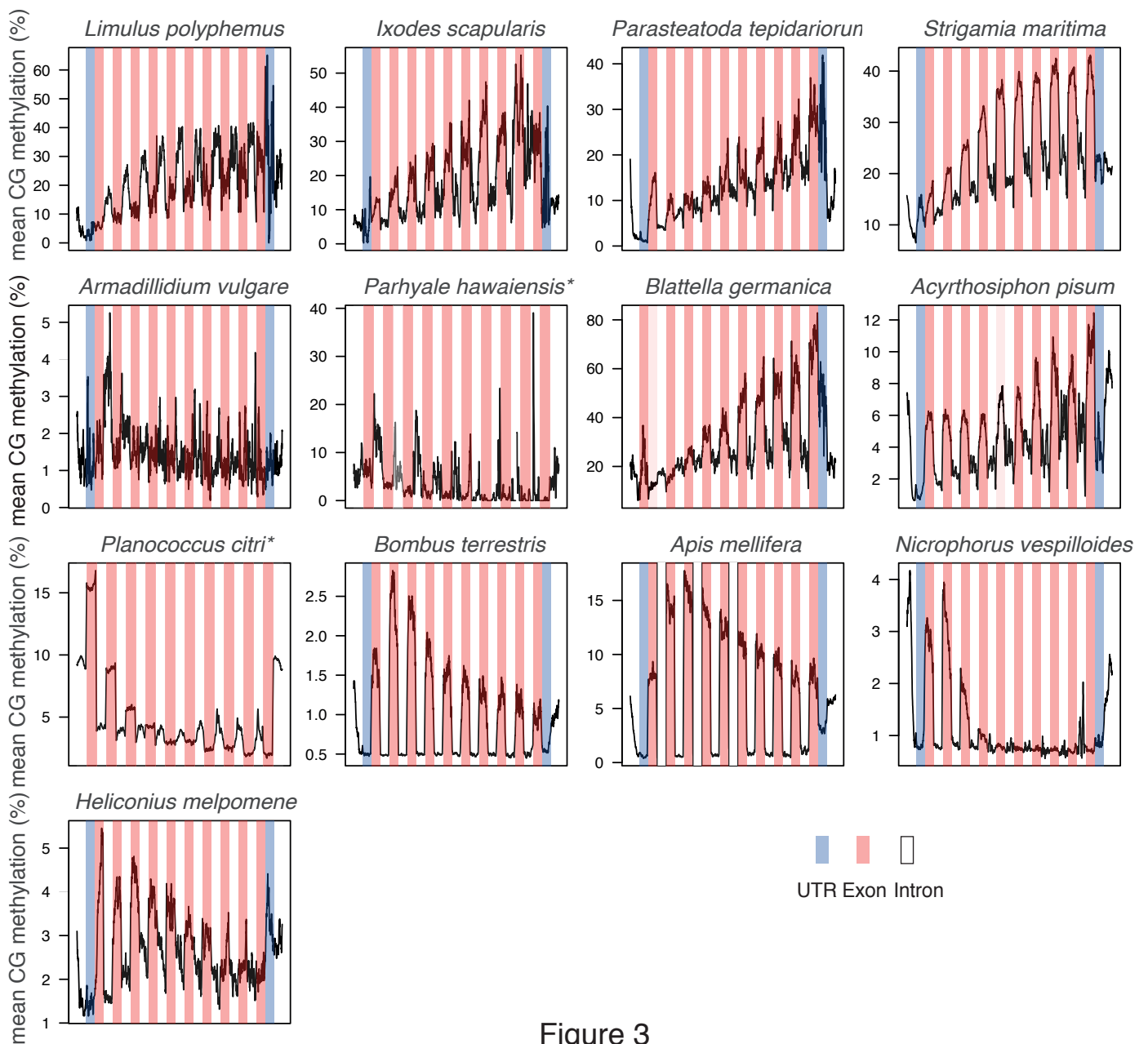
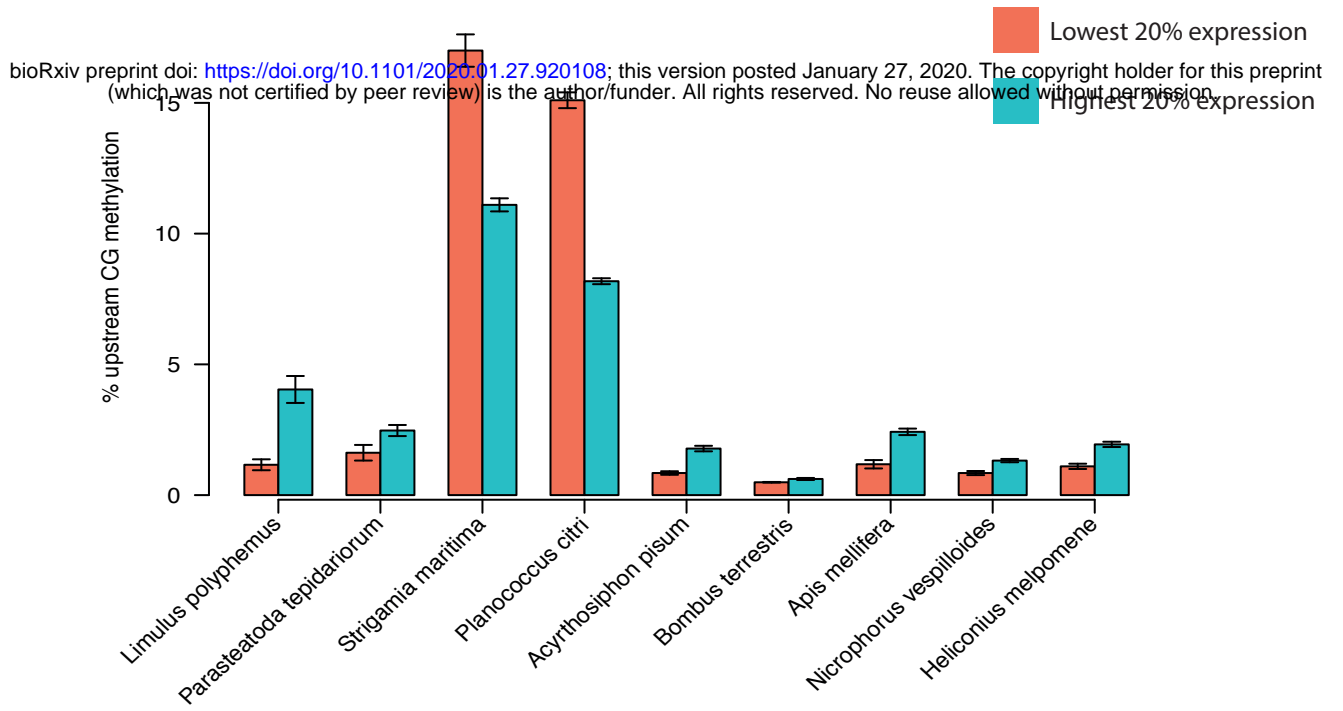
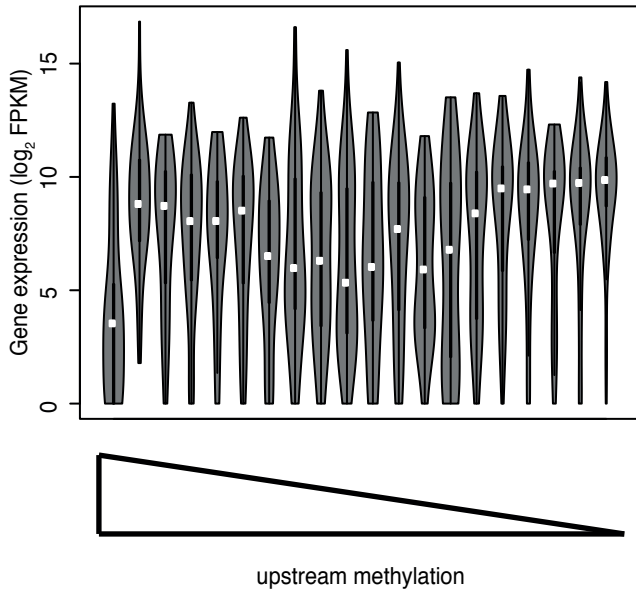
**A****B****C**

Figure 3

A



B



C

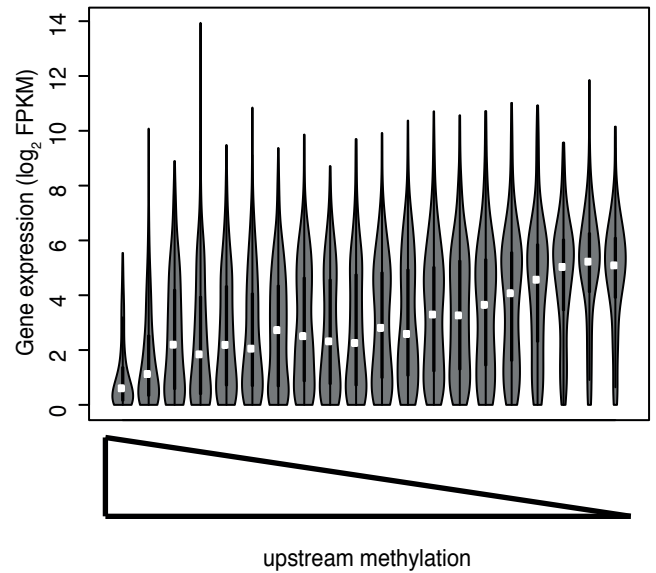


Figure 4



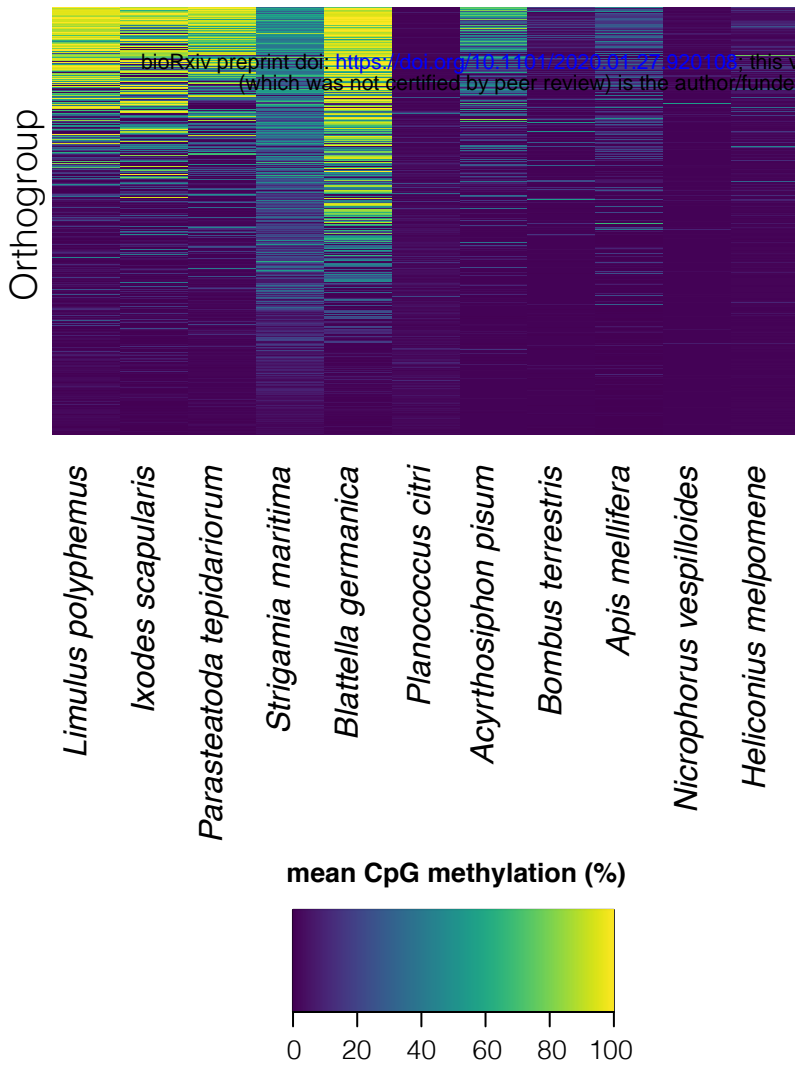
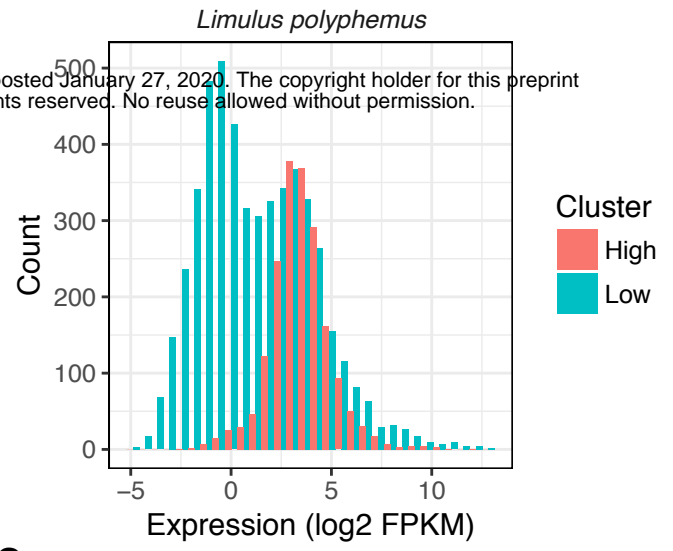
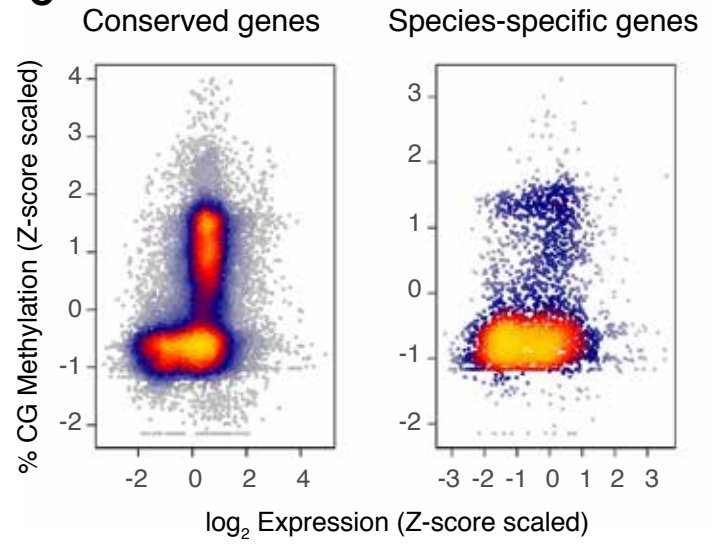
**A****B****C**

Figure 5

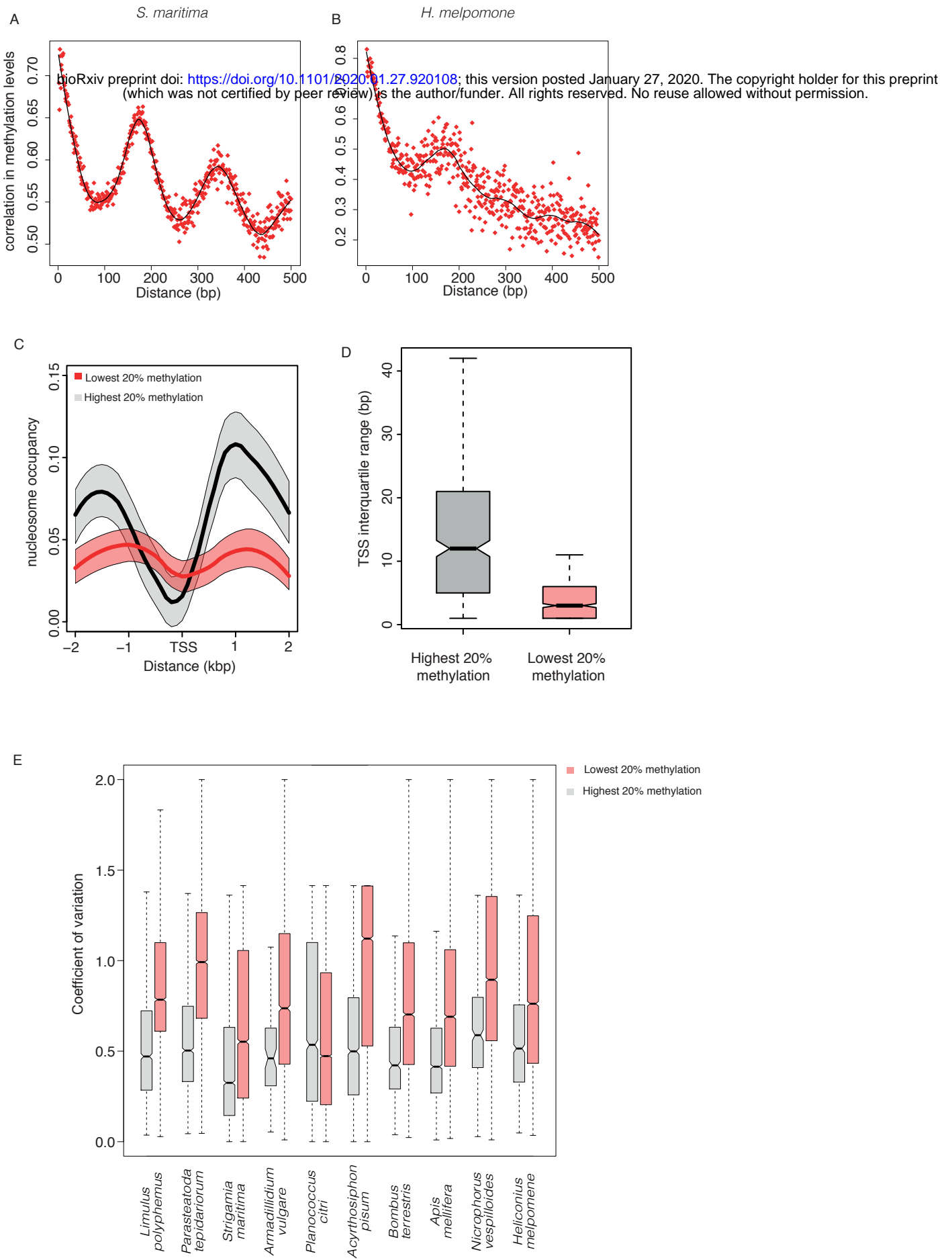


Figure 6