

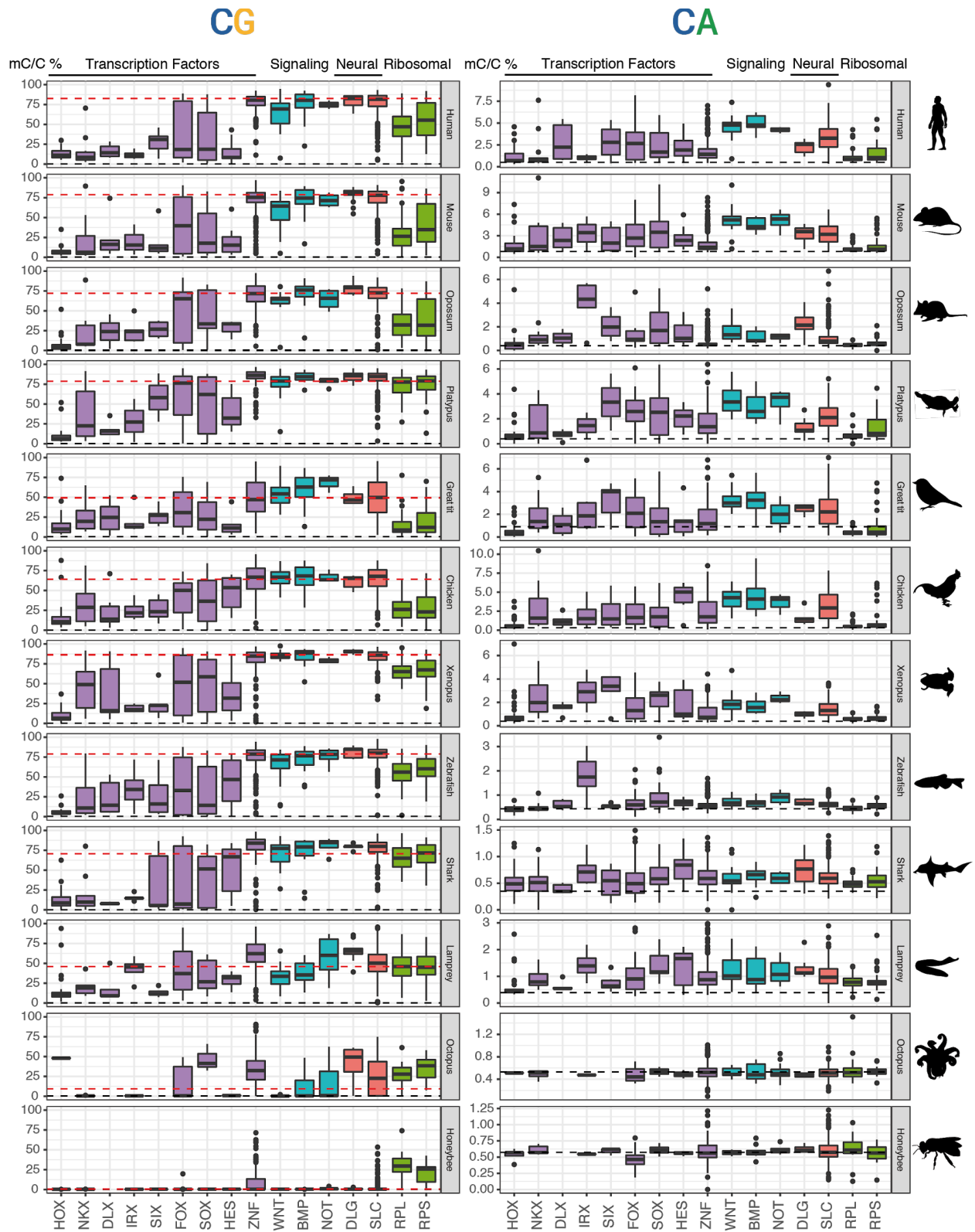
Supplementary material

The emergence of the brain non-CpG methylation system in vertebrates

Alex de Mendoza, Daniel Poppe, Sam Buckberry, Jahnvi Pflueger, Caroline B. Albertin, Tasman Daish, Stephanie Bertrand, Elisa de la Calle Mustienes, Jose Luis Gomez-Skarmeta, Joseph R. Nery, Joseph R. Ecker, Boris Baer, Clifton W. Ragsdale, Frank Grützner, Hector Escriva, Byrappa Venkatesh, Ozren Bogdanovic, Ryan Lister

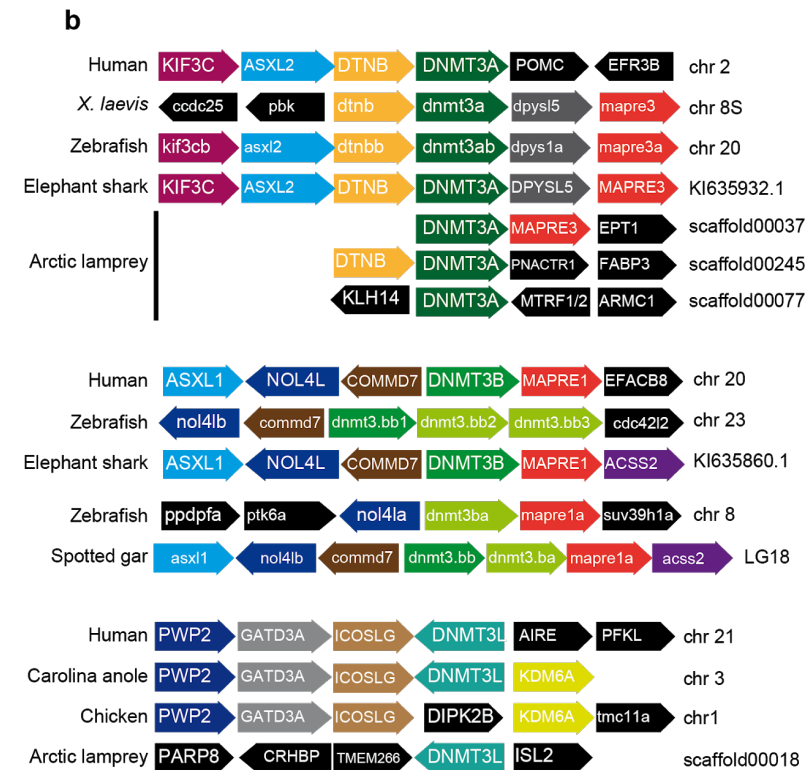
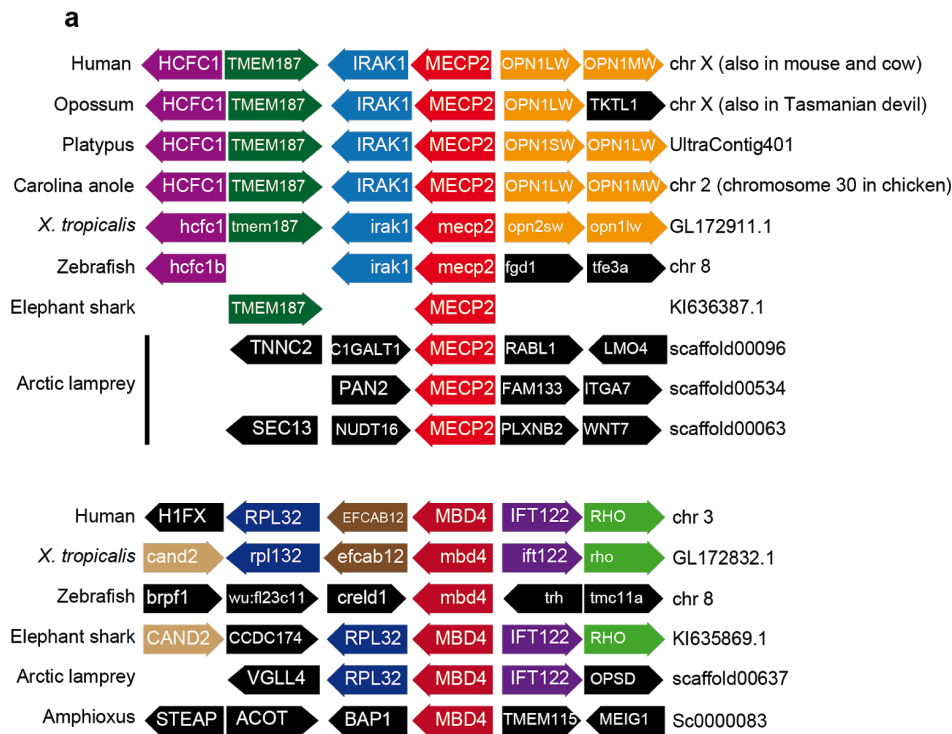
Correspondence to: ryan.lister@uwa.edu.au

Supplementary Figures 1-3.



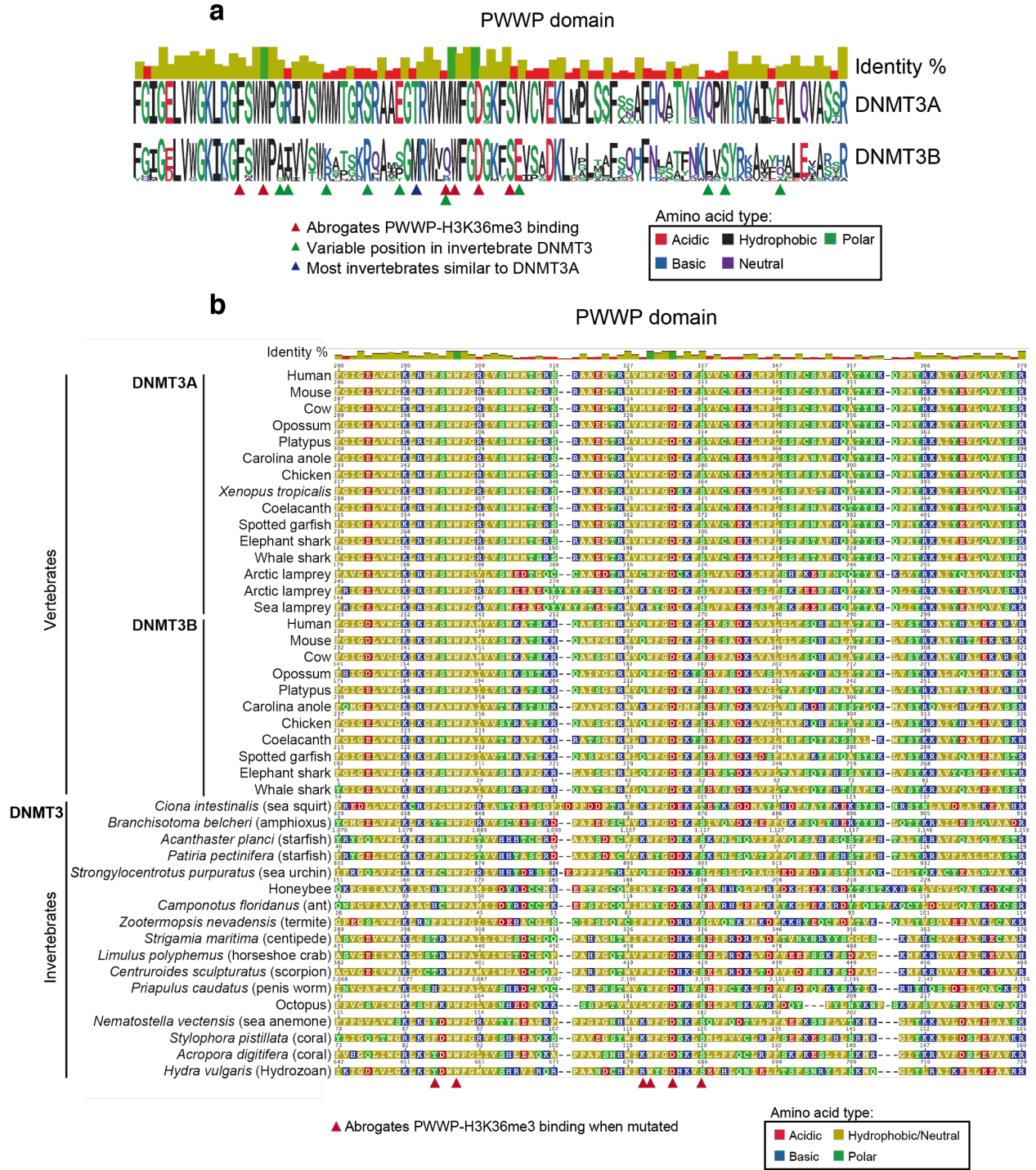
Supplementary Figure 1 | Gene body methylation on gene families/classes. Distribution of gene body methylation on genes belonging to specific transcription factor classes, signaling

molecules, neural related and ribosomal functions. Black dashed lines indicate the bisulfite non-conversion rate for each WGBS library, red dashed lines indicate the global CpG methylation levels in each species. HOX, NKX, DLX, IRX and SIX are Homeobox transcription factor (TF) families. FOX are forkhead domain TFs. SOX are HMG-box TFs belonging to the Sox family. HES are bHLH TFs belonging to the Hes gene family. ZNF are zinc finger C2H2 TFs. WNT, BMP and NOT (Notch) are signaling peptides and receptors. DLG (Discs Large Homologue) and SLC (Solute Carriers) are genes enriched in neural functions. RPL (L Ribosomal Proteins) and RPS (S Ribosomal Proteins) are genes with ribosomal functions. The median CpA methylation levels of invertebrate genes overlaps the non-conversion rate. Boxplot centre lines are medians, box limits are quartiles 1 (Q1) and 3 (Q3), whiskers are $1.5 \times$ interquartile range (IQR) and points are outliers.



Supplementary Figure 2 | Syntenic regions surrounding MeCP2, MBD4, and DNMT3 genes.

(a) Gene order and orientation around the MECP2 and MBD4 loci in the genomes of selected species. Lampreys encode three copies of MECP2. Amphioxus MBD4/MECP2 locus belongs to the *Branchiostoma lanceolatum* genome, and is consistent with previous synteny conservation observed previously in *B. floridae*⁴⁷. (b) Gene order and orientation around the DNMT3 loci in selected species genomes. dnmt3ab and dnmt3bb are contiguous in the spotted gar genome assembly, whereas zebrafish has split paralogues in different scaffolds (dnmt3bb.1 is the conserved DNMT3B orthologue, whereas dnmt3ba, dnmt3bb.2 and dnmt3bb.3 are the teleost-specific duplication harbouring a Calponin homology domain). The DNMT3L syntenic region is conserved in bird genomes, however it does not contain the DNMT3L gene. Genes with the same colour are orthologous. Species with less than six genes are found in genomic scaffolds that do not contain more genes, either because it is too short or because it is truncated on one side. Gene names are based on UCSC genome browser annotation or best BLASTP hit against Uniprot in the case of the arctic lamprey.



Supplementary Figure 3 | PWWP domain conservation across DNMT3 genes. (a) Amino acid sequence motifs representing the PWWP domain of DNMT3A and DNMT3B orthologues across jawed vertebrates. Coloured triangles indicate specific positions that are known to be important for PWWP binding to H3K36me3 (red), differ between DNMT3A and DNMT3B but are variable in

invertebrate DNMT3 PWWP domain (green), and positions in which invertebrate DNMT3 are similar to DNMT3A. Amino acid letters are colour coded per type as indicated in the legend. Identity shows the total identity in a multi-sequence alignment including DNMT3A, DNMT3B and invertebrate DNMT3 sequences. **(b)** Complete amino acid multi-sequence alignment of the PWWP (Pro-Trp-Trp-Pro) motif from DNMT3A, DNMT3B, and invertebrate DNMT3 genes. Alignment visualised using Geneious software.

Differences between mammalian CpH methylation patterns in neural and stem cells are attributed to the distinct binding affinities of DNMT3A and DNMT3B⁴⁶. Stem cell CpH methylation and DNMT3B are preferentially localised on transcribed gene bodies and enriched in the CAG trinucleotide context. DNMT3B localisation to gene bodies is mediated by the PWWP domain binding of H3K36me3⁸⁸. In contrast, DNMT3A does not show such a preference in human embryonic stem cells⁸⁸, despite encoding a conserved PWWP domain that is able to bind H3K36me2/3 *in vitro*⁸⁹. Furthermore, neural CpH methylation is widespread across the genome, suggesting that DNMT3A does not exclusively bind to gene bodies¹⁰. Additionally, invertebrate CpG methylation is mostly restricted to gene bodies and has been shown to co-localise with H3K36me3⁹⁰, suggesting that DNMT3B might be more similar to the ancestral invertebrate DNMT3 enzymes. To test this idea, we constructed an alignment of DNMT3 enzymes (Supplementary Fig. 9). Counterintuitively, we found that the PWWP domains of invertebrate DNMT3 share higher identity with those of DNMT3A (44%) than DNMT3B (39.2%), highlighting that DNMT3A has accumulated less changes than DNMT3B. Five of the six mutations in the PWWP domain that have been shown to abrogate DNMT3A and DNMT3B preference for H3K36me3^{17,44,88} are highly conserved between the orthologues (Fig. 4d), and thus are not likely to mediate differences in histone tail modification affinity. However, glutamine in position 232 (Q232) of human DNMT3B is not shared with DNMT3A genes, which has a hydrophobic methionine in that position. Therefore, differences between PWWP domain key amino acid

positions are likely to explain the distinct binding affinities of DNMT3 family enzymes in vertebrates, likely fostering the specialisation of each ohnologue into distinct binding patterns and expression domains. Supporting the central role of PWWP in DNMT3 diversification, both DNMT3L and DNMT3C, a rodent-specific DNMT3B paralogue⁹¹, have lost the PWWP domain, which likely allowed novel functions for these genes to emerge. A similar process of specialisation must have occurred at the methyltransferase domain, given the known preference for CAT trinucleotide of DNMT3B, instead of the CAC that characterises DNMT3A and neural CpH methylation. In summary, the redundancy of DNMT3 enzymes is a product of the vertebrate ancestral WGD that was fundamental to allow neofunctionalization of the ohnologues, likely conferring neural-specific functions to DNMT3A at the origin of vertebrates.