# Supplemental Material

## Supplemental Tables

| reads | counts |
|---|---:|
| sequenced (PE1+PE2) | 28,607,102 |
| aligned | 10,551,289 |
| uniquely aligned | 5,656,887 |
| multiply aligned | 4,894,402 |
| not aligned | 18,055,813 |
| (aligned pairs | 3,269,861) |

**Table S1. Summary of sequencing and alignment.**

| chromosome | unique | multiple, primary | multiple, secondary | CT strand unique+multiple | GA strand unique+multiple | ambiguous unique+multiple | inconsistent unique+multiple |
|---|---:|---:|---:|---:|---:|---:|---:|
| chr2L | 777,176 | 300,491 | 473,536 | 704,331 | 778,144 | 55,110 | 13,618 |
| chr2LHet | 2,353 | 4,484 | 8,702 | 6,973 | 7,557 | 1,007 | 2 |
| chr2R | 918,103 | 376,437 | 590,496 | 888,397 | 924,331 | 58,060 | 14,248 |
| chr2RHet | 22,844 | 44,843 | 86,204 | 73,632 | 72,494 | 7,691 | 74 |
| chr3L | 895,666 | 380,247 | 608,774 | 844,929 | 924,050 | 97,564 | 18,144 |
| chr3LHet | 20,230 | 41,584 | 79,222 | 69,894 | 66,991 | 4,107 | 44 |
| chr3R | 1,128,338 | 481,935 | 770,886 | 1,095,641 | 1,137,840 | 121,472 | 26,206 |
| chr3RHet | 20,661 | 38,654 | 71,630 | 60,010 | 62,820 | 8,063 | 52 |
| chr4 | 12,482 | 17,669 | 32,282 | 22,702 | 35,870 | 3,847 | 14 |
| chrX | 1,783,287 | 1,247,351 | 1,812,379 | 2,440,229 | 2,113,552 | 263,102 | 26,134 |
| chrXHet | 2,266 | 4,051 | 7,953 | 6,524 | 7,311 | 433 | 2 |
| chrYHet | 1,232 | 3,472 | 6,571 | 5,522 | 5,162 | 591 | 0 |
| chrU | 72,236 | 1,953,063 | 4,036,585 | 458,907 | 993,522 | 4,607,289 | 2,166 |
| chrM | 13 | 121 | 247 | | | | |
| TOTAL | 5,656,887 | 4,894,402 | 8,585,467 | | | | |
| TOTAL (-chrM) | 5,656,874 | 4,894,281 | 8,585,220 | 6,677,691 | 7,129,644 | 5,228,336 | 100,704 |

**Table S2. Alignments by chromosome.**
Paired-read alignments are assigned to either the CT of the GA strand; only reads mapping to a strand are retained. If a read can be mapped with equal probability to either strand, it is labeled "ambiguous". If the two reads of a pair do not map to the same strand, they are labeled "inconsistent". Ambiguous and inconsistent reads are not analyzed further.

|  |  | methylated regions | | resequenced regions | |
|---|---|---|---|---|---|
| Condition |  | CT strand | GA strand | CT strand | GA strand |
| 2 methylated reads/position 3 methylated positions/region | positive | 879 | 953 | 40 | 17 |
|  | negative | 321 | 268 | 2 | 1 |
|  | undetermined | 11337 | 11650 | 5 | 1 |
| 2 methylated reads/position 7 methylated positions/region | positive | 383 | 388 | 31 | 11 |
|  | negative | 215 | 196 | 2 | 1 |
|  | undetermined | 11939 | 12287 | 14 | 7 |
|  | total number of regions | 12537 | 12871 | 47 | 19 |

**Table S3. Comparison with whole genome bisulfite data.**
Methylated regions identified in this study were compared with data generated by whole genome bisulfite sequencing by Raddatz *et al.* "Positive" regions are methylated regions identified in this study that have supporting evidence in the Raddatz *et al.* data; we illustrate two conditions of different stringency (see Methods). "Negative" regions are regions identified in this study that have sufficient coverage in the Raddatz data to reveal methylation if present, but lack support in that data. "Undetermined" regions are not positive, and lack sufficient coverage to provide confidence that they are truly negative. The vast bulk of our methylated regions are positive or undetermined in the Raddatz data, and when positive they usually remain positive under the more stringent condition used to determine the status of a region. The coverage threshold we used to call regions as negative (100x) would allow detection of some regions methylated on ~1% of alleles. Many of the validated regions found in our study are methylated on <5% of alleles (Figure 3 and Figure S5). **Methods for comparison with whole-genome bisulfite sequencing data from Zemach *et al*. and Raddatz *et al*.:** The whole-genome bisulfite sequencing data for Stage 5 *Drosophila* embryo generated by Zemach *et al.* (Zemach et al. 2010) (accession #: GSM497255) and Raddatz *et al.* (Raddatz et al. 2013) (accession #: GSM983094) were downloaded from the GEO database. Sequence reads were aligned with Novoalign v2.07.11 with default options in bisulfite mode with the "b2" (directional) option and reporting only unique alignments, against a reference consisting of either the sequences that we validated by bisulfite PCR or the sequences of the 25,497 methylated regions identified in this study. Reads were assigned to the CT or GA strands according to Novoalign mappings. Reads aligning to the strand that was not amplified were discarded, and reads that are potential PCR duplicates were removed with MarkDuplicates from the Picard suite (picard.sourceforge.net). We determined sequence coverage and the percentage of methylation at each cytosine using the output of mpileup from the samtools suite (Li et al. 2009). We define a cytosine as methylated if at least two reads are unconverted at that position. To evaluate the agreement between Raddatz *et al.* and our data, the reference regions described above were divided into three groups according to the evidence for methylation in Raddatz *et al.*'s data: "positive", if the region contains at least three methylated cytosines (or seven in the more stringent condition); "negative", if the region contains less than 3 methylated cytosines and at least three (or seven) cytosines with coverage greater than 100 reads; "undetermined", if the region is neither positive nor negative.
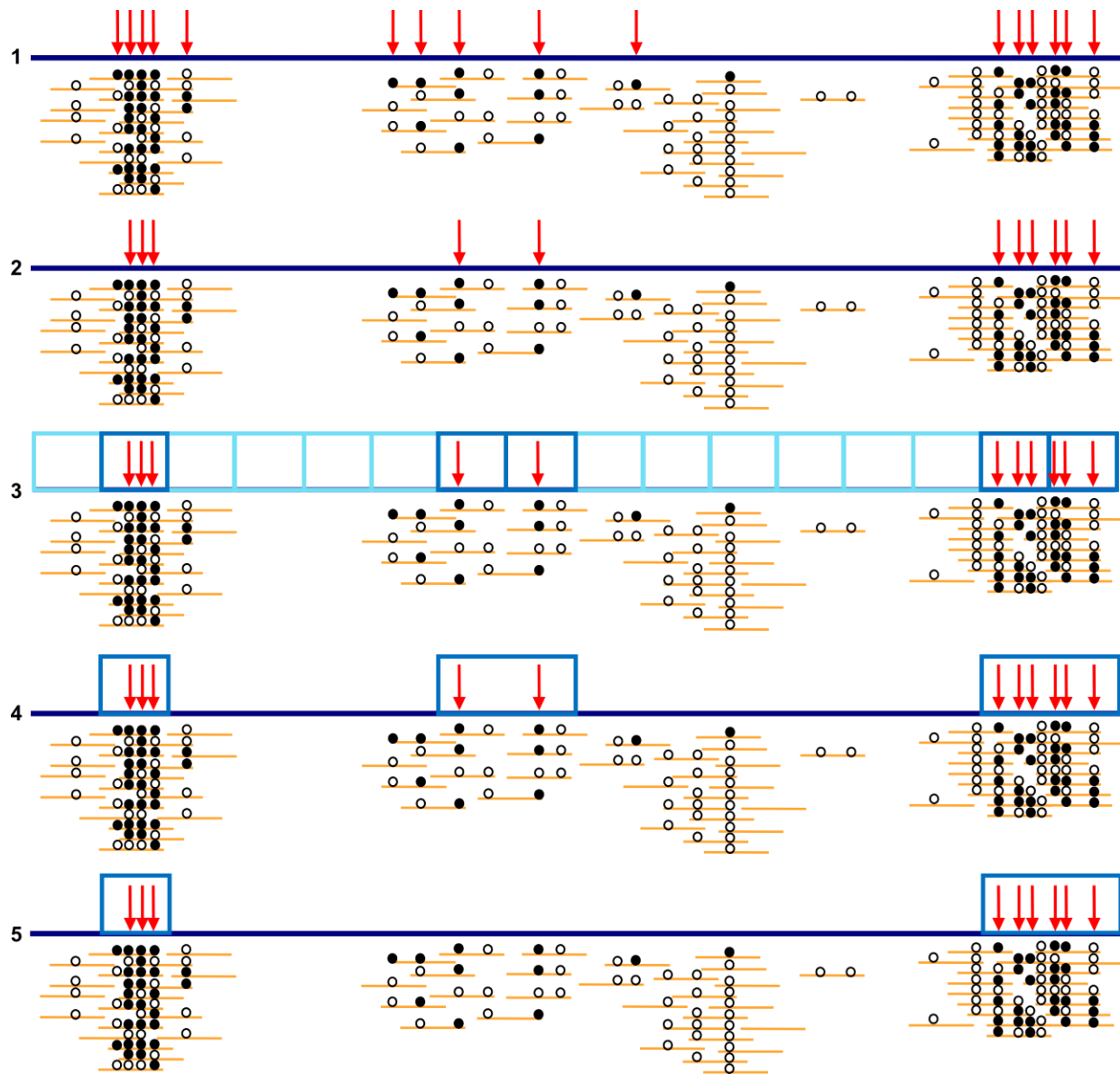
**Figure S1. Methylated regions in the genome of Stage 5 *Drosophila* embryos**. Procedure for identification of methylated cytosines and methylated regions from MeDIP-Bseq data. The steps illustrated here led to the identification of the 25,497 methylated regions discussed in the text. Red arrows identify cytosine positions that pass the filter at each step; the steps are applied sequentially to the output of the preceding step.

- o Step 1 identifies cytosine positions at which the ratio of C-containing alignments (unconverted, i.e. methylated) over the sum of C- and T-containing alignments was greater than 0.1 (methylated cytosine: closed circles; unmethylated: open circles).
- o Step 2 identifies cytosine positions at which at least three alignments contain a methylated cytosine.
- o Step 3 divides the genome into contiguous 25-base segments.
- o Step 4 removes segments that do not contain cytosines passing Step 2, and merges contiguous segments.
- o Step 5 retains only those segments in which the alignments contained at least 25 methylated cytosines; these are the methylated regions.
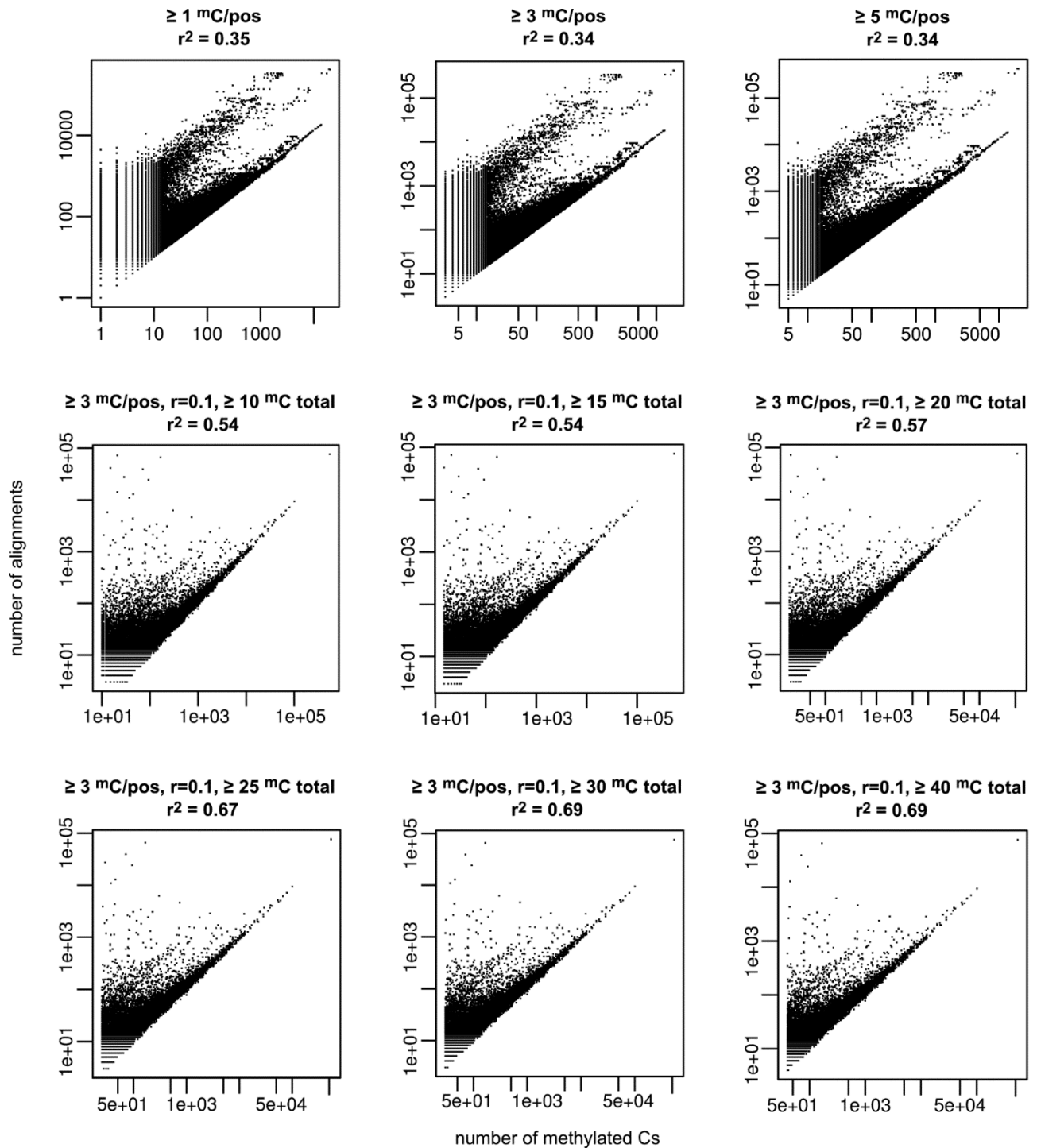
**Figure S2. Effects of different parameter choices on the identification of methylated cytosines in the MeDIP-Bseq data.** In all of the plots displayed, the number of methylated cytosines in the reads aligning to a given cytosine in the reference sequence is shown on the x-axis, and the number of sequence reads containing that position is shown on the y-axis. Only those cytosines meeting the condition shown at the top of each plot are displayed (see Methods for a detailed description of the conditions). The procedure described in the text and Figure S1 removes cytosine positions with weakly supported methylation states. The top row is derived from step 2 in Figure S1: it illustrates the effect of requiring more methylated reads supporting the status of the position. The middle and bottom rows are derived from the step shown in Figure S1, step 5: only those cytosines included in one of the regions meeting the parameters are shown. For each set of parameters (denoted at the top of each plot), the correlation coefficient between the number of methylated cytosines and the number of alignments within a region was calculated. We chose the set of parameters that optimizes the correlation coefficient at the lowest cost in discarded methylated regions.
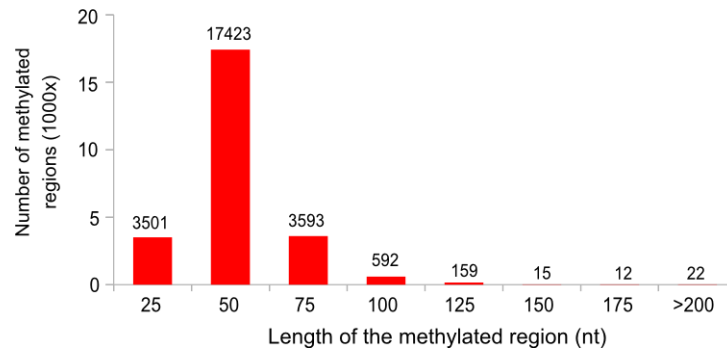
**Figure S3. Length distribution of methylated regions.** The number of methylated regions of a given length range is shown at the top of each column. 97% of the regions have a length of 75 bases or less.
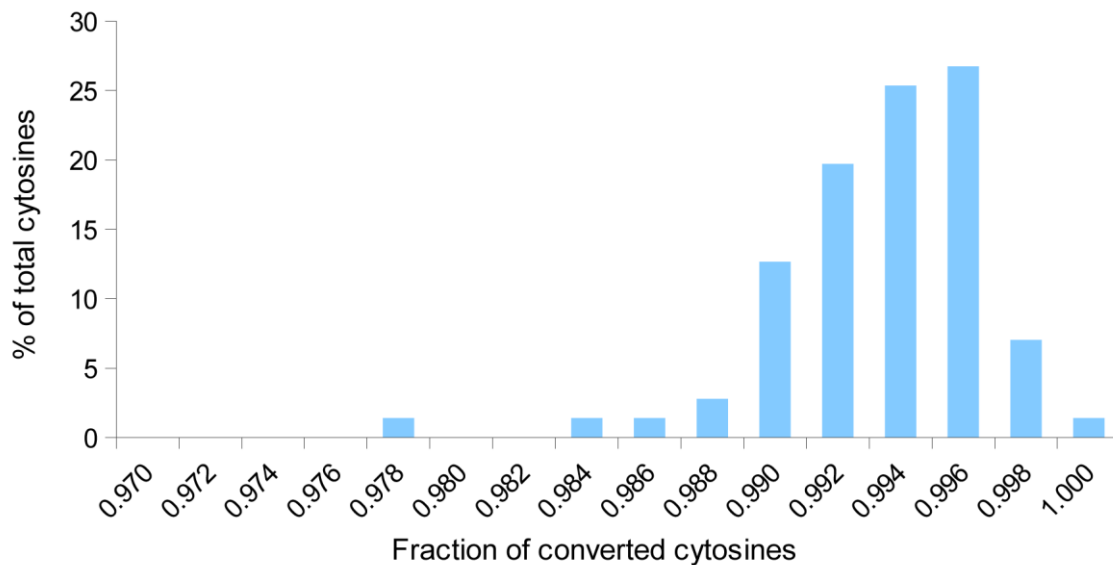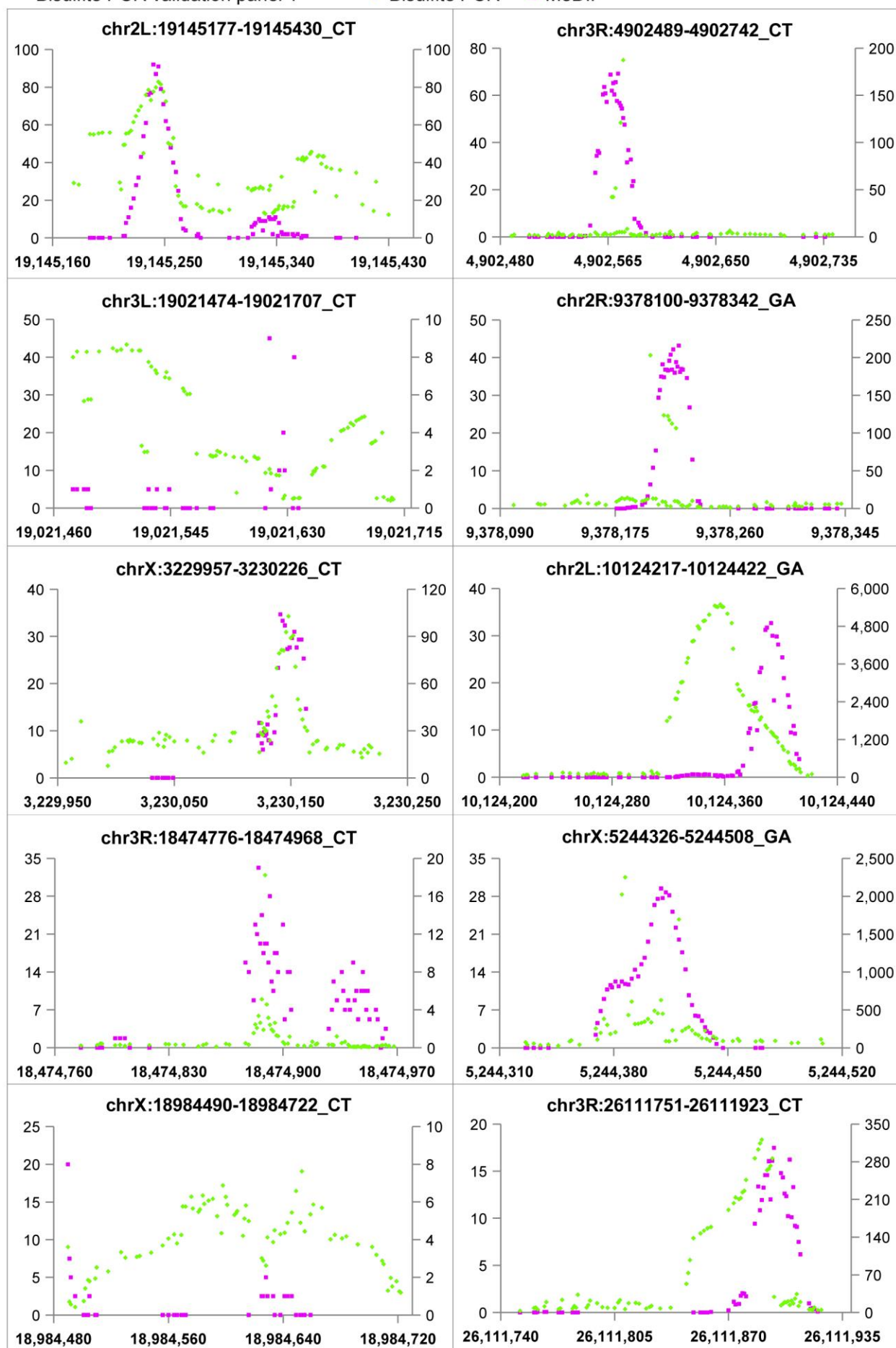


**Figure S4. Efficiency of bisulfite conversion as determined by bisulfite-PCR sequencing of lambda phage DNA.** Lambda phage DNA, grown in an *E. coli* strain deficient for methylation, was bisulfite converted in the same reaction as the *Drosophila* DNA used for the validation in Figures 3, S5, S6 and S8. Two segments of the lambda genome were amplified and sequenced to a median coverage of 109, 207 reads (range 82,602-161,273). The x-axis shows the rate of conversion from C to T, determined as the ration of T over C +T at each cytosine position. The y-axis shows the percent of total cytosines with a given conversion rate. Average conversion is >99%, and 96% of all cytosine positions had a conversion rate ≥ 0.988.

Bisulfite PCR validation panel 1 — Bisulfite PCR (green diamonds), MeDIP (magenta squares). Bisulfite PCR (% methylation) on left axis; MeDIP (number of methyl cytosines) on right axis. Panels: chr2L:19145177-19145430_CT, chr3R:4902489-4902742_CT, chr3L:19021474-19021707_CT, chr2R:9378100-9378342_GA, chrX:3229957-3230226_CT, chr2L:10124217-10124422_GA, chr3R:18474776-18474968_CT, chrX:5244326-5244508_GA, chrX:18984490-18984722_CT, chr3R:26111751-26111923_CT.

Bisulfite PCR validation panel 2 — Bisulfite PCR (green diamonds), MeDIP (magenta squares). Panels: chr3L:10937398-10937633_CT, chr3R:12855269-12855564_CT, chr2L:6787664-6787866_CT, chr2R:1660580-1660794_CT, chr3R:15048426-15048672_GA, chrX:16425570-16425728_CT, chr3L:5844440-5844692_CT, chr2R:5346056-5346273_CT, chrX:5999750-5999959_GA, chrX:14609561-14609724_CT. Left axis: Bisulfite PCR (% methylation). Right axis: MeDIP (number of methyl cytosines).

Bisulfite PCR validation panel 3

Bisulfite PCR validation panel **4**   ◆ Bisulfite PCR   ■ MeDIP

Bisulfite PCR validation panel 5 ◆ Bisulfite PCR ■ MeDIP

Bisulfite PCR validation panel 6 — ♦ Bisulfite PCR ■ MeDIP

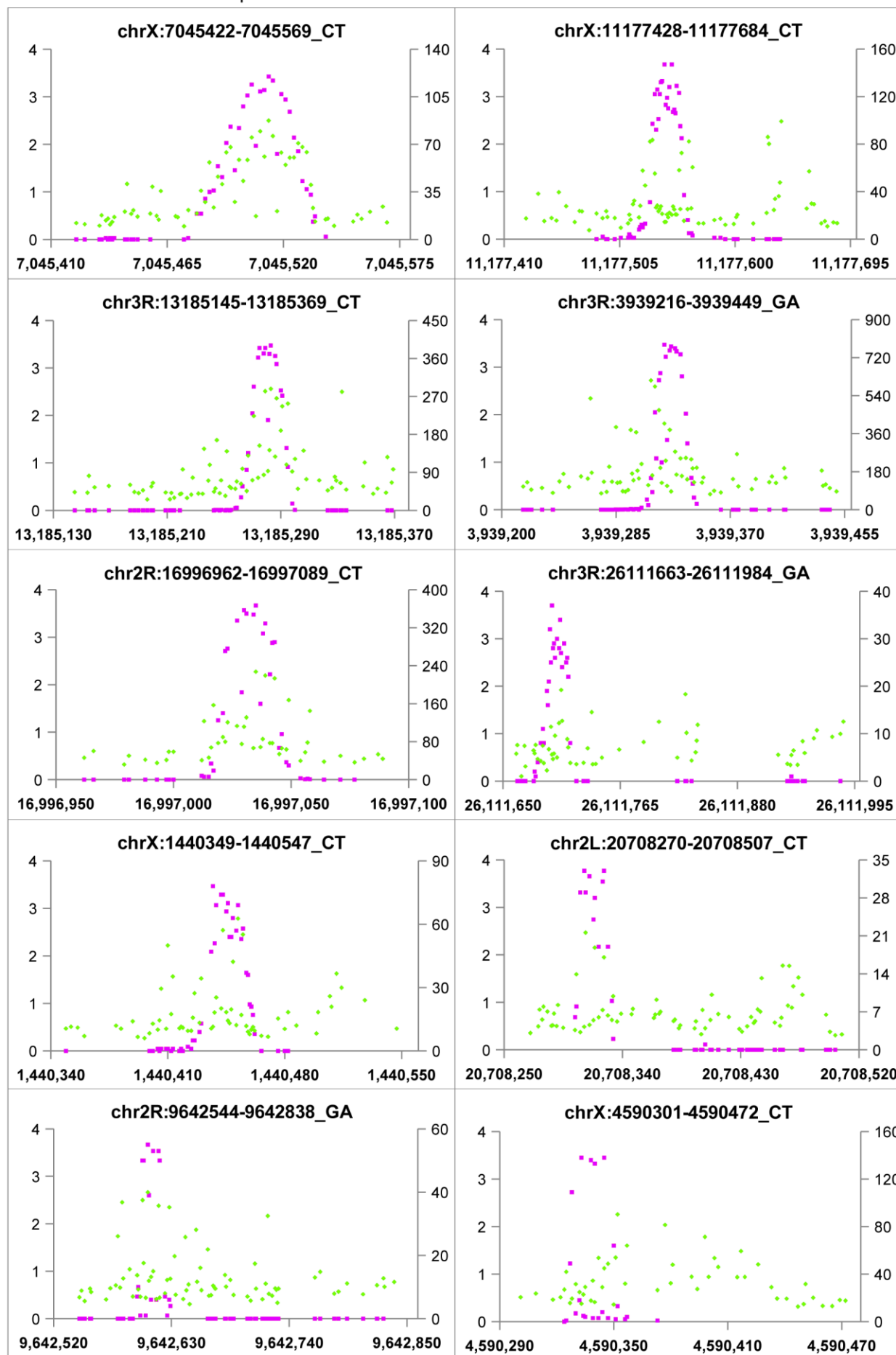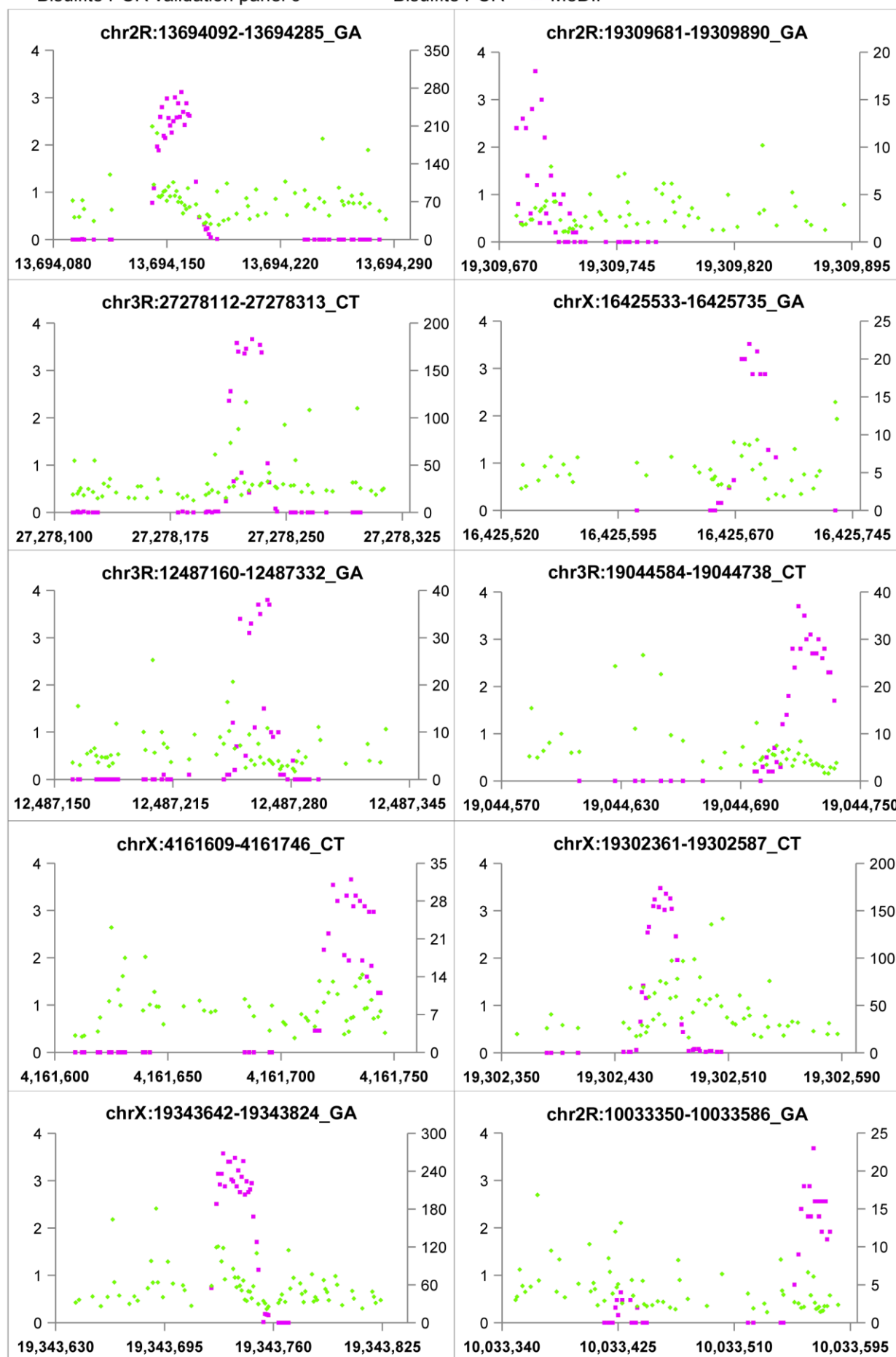Bisulfite PCR (% methylation)
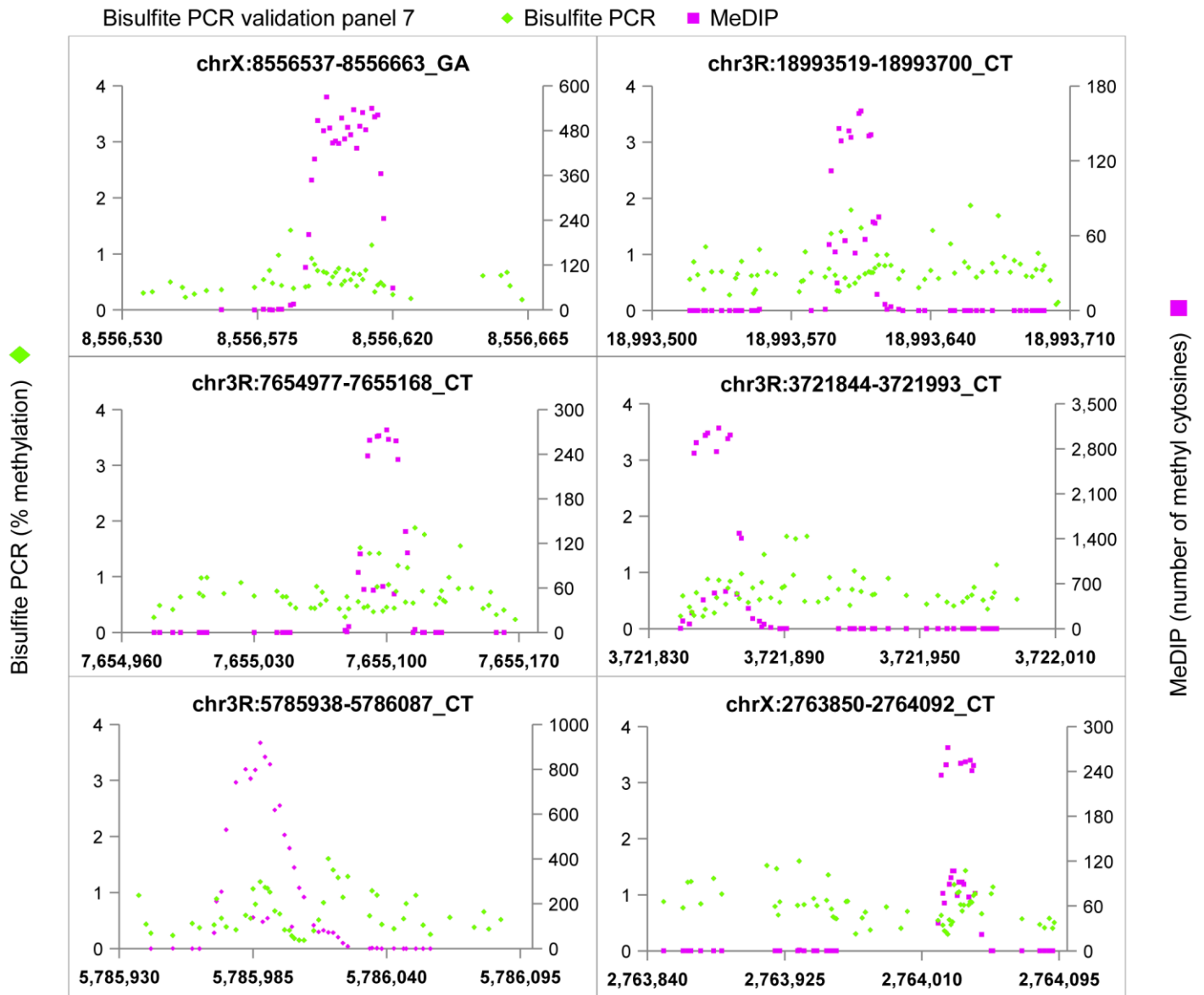
MeDIP (number of methyl cytosines)

**Figure S5. Direct amplification of bisulfite-converted DNA confirms methylation patterns**. The full set of 66 regions that were analyzed is shown. Methylated regions identified by MeDIP-bisulfite sequencing were PCR amplified from bisulfite converted DNA and Illumina sequenced to at least 10,000X coverage. Each dot represents one cytosine (green – bisulfite PCR; purple – MeDIP bisulfite). The y-axis at the left indicates the percent of methylated cytosines in the bisulfite PCR; the y-axis at the right indicates the number of methylated cytosines detected by MeDIP bisulfite. While the MeDIP bisulfite analysis is not quantitative, bisulfite PCR demonstrates the proportion of methylated cytosines at a given position, as well as the pattern of methylation of the amplified region. There is good agreement in the pattern of methylation detected by the two methods.
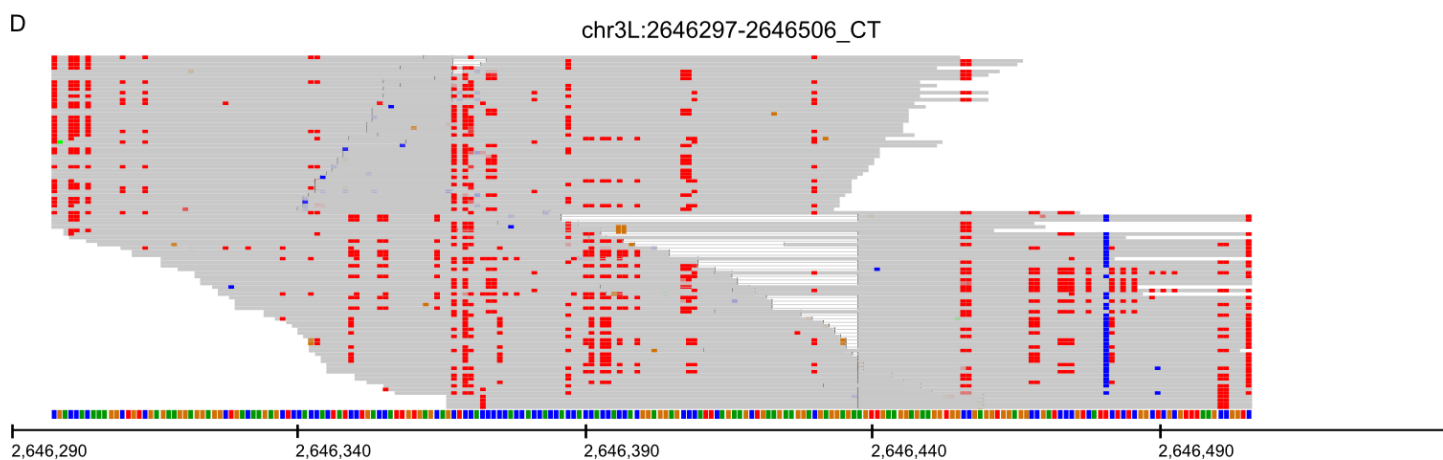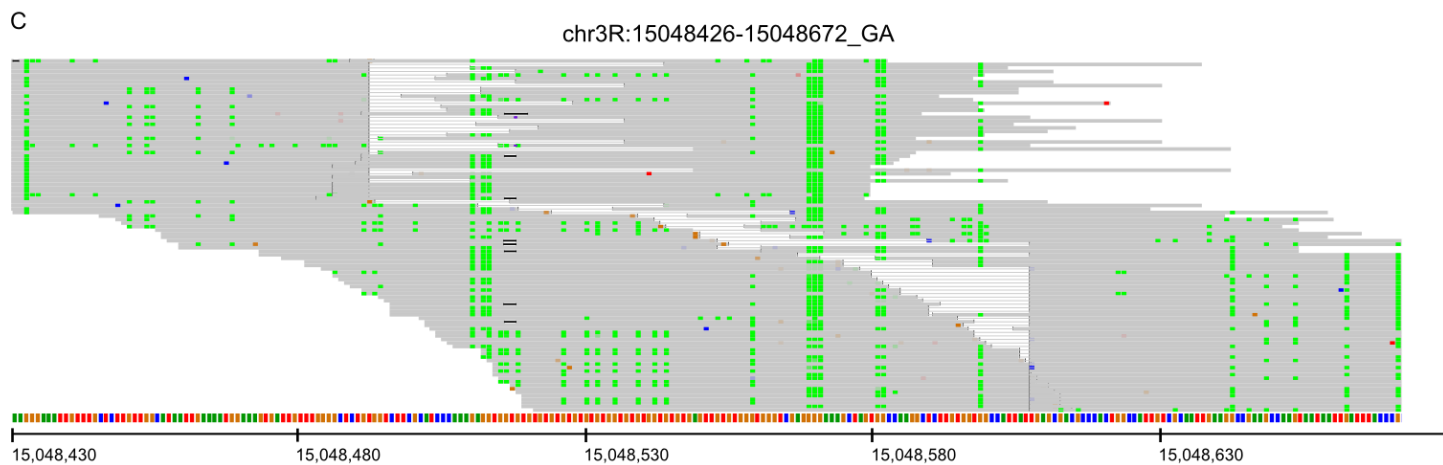
A — chrX:3229957-3230226_CT

B — chr3L:10937398-10937633_CT

C — chr3R:15048426-15048672_GA
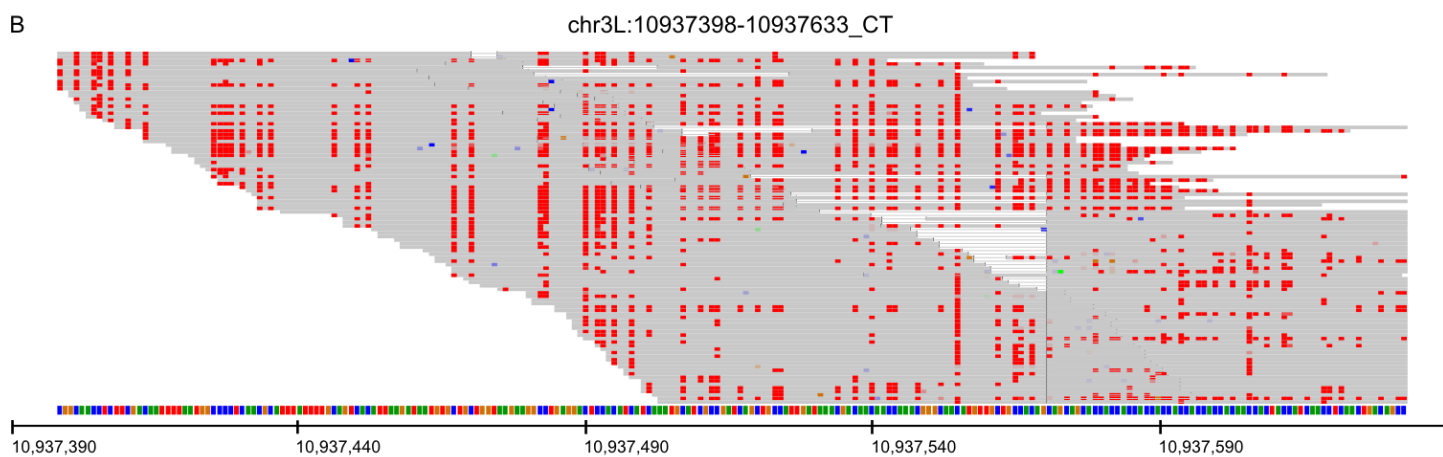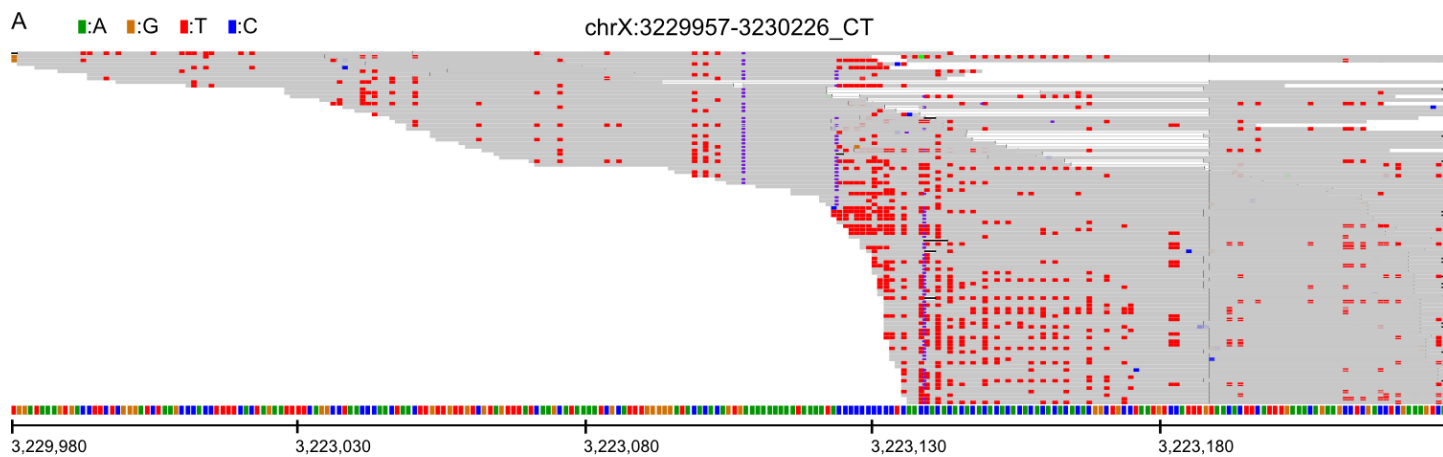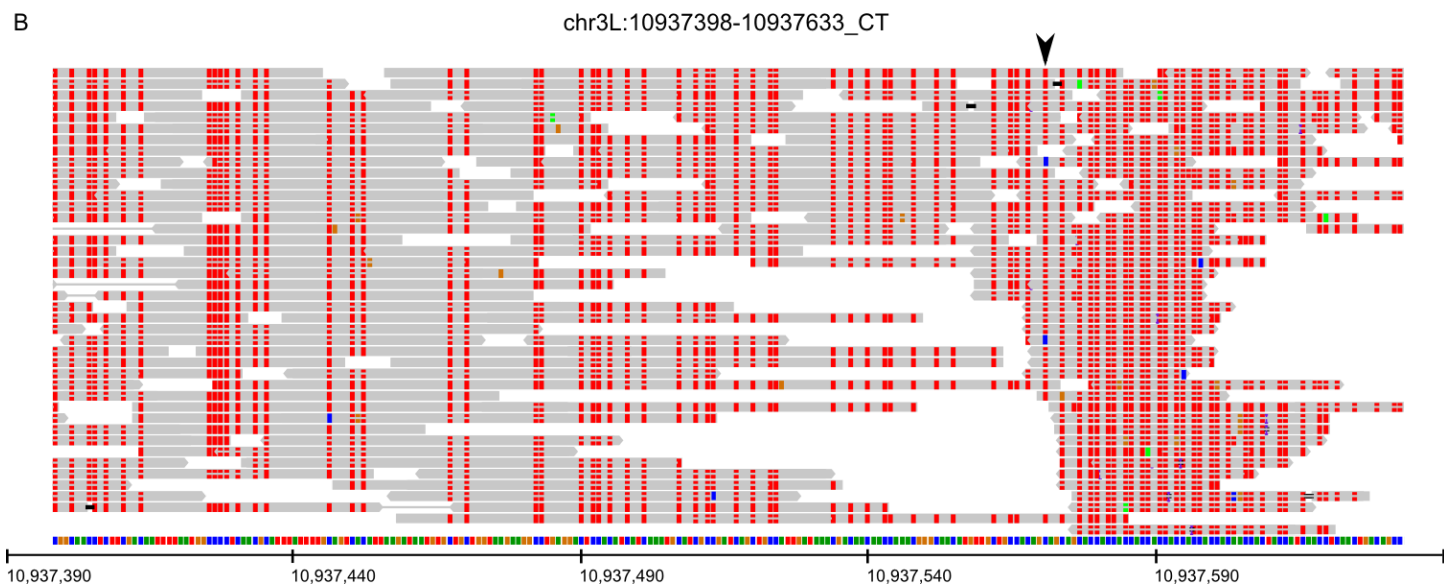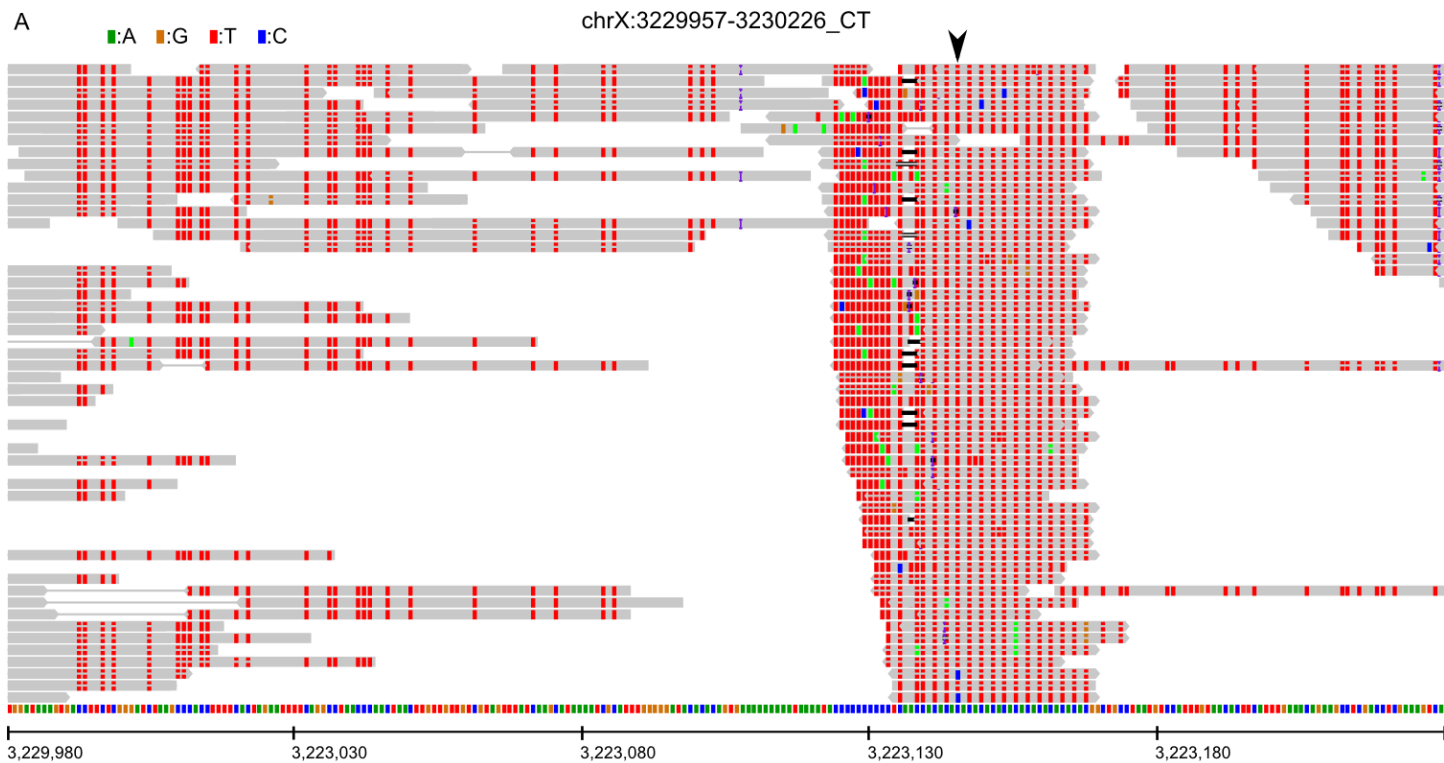
D — chr3L:2646297-2646506_CT

**Figure S6. Visualization of alignments at four methylated regions, illustrating correct alignment of reads that support cytosine methylation**. Shown are 100 base (including a 6 base index) paired end reads aligning to the regions displayed in Figure 3; these reads derived from wild type *EP(2)GE 15695* flies. Only a subset of read pairs (100) is shown; reads were selected based on their content of unconverted cytosines, in order to illustrate the alignment of such reads. Each line shows the alignment of a read pair; when the two sequences of a pair of reads do not overlap, a thin grey line shows their connection. Alignments are displayed using the "collapsed" mode of the Integrative Genome Viewer; in this mode, the direction of the alignments is shown by a vertical grey line at the 3' end of the alignment. The color-coded reference sequence is at the bottom of each panel, with the color key shown at the top of the figure. A match between a read and the reference is shown in grey; a mismatch is shown with the color of the mismatched base. Unmethylated cytosines are sequenced as 'T' (red) on the CT strand (panels A, B, and D) and as 'A' (green) on the GA strand (panel C); thus converted (unmethylated) cytosines are shown in color, and any methylated cytosines are denoted by gray color at a position that is colored in other reads. The figure shows that alignments of reads containing unconverted (methylated) cytosines are unambiguous and extend well beyond the low complexity sequences where methylation is concentrated.
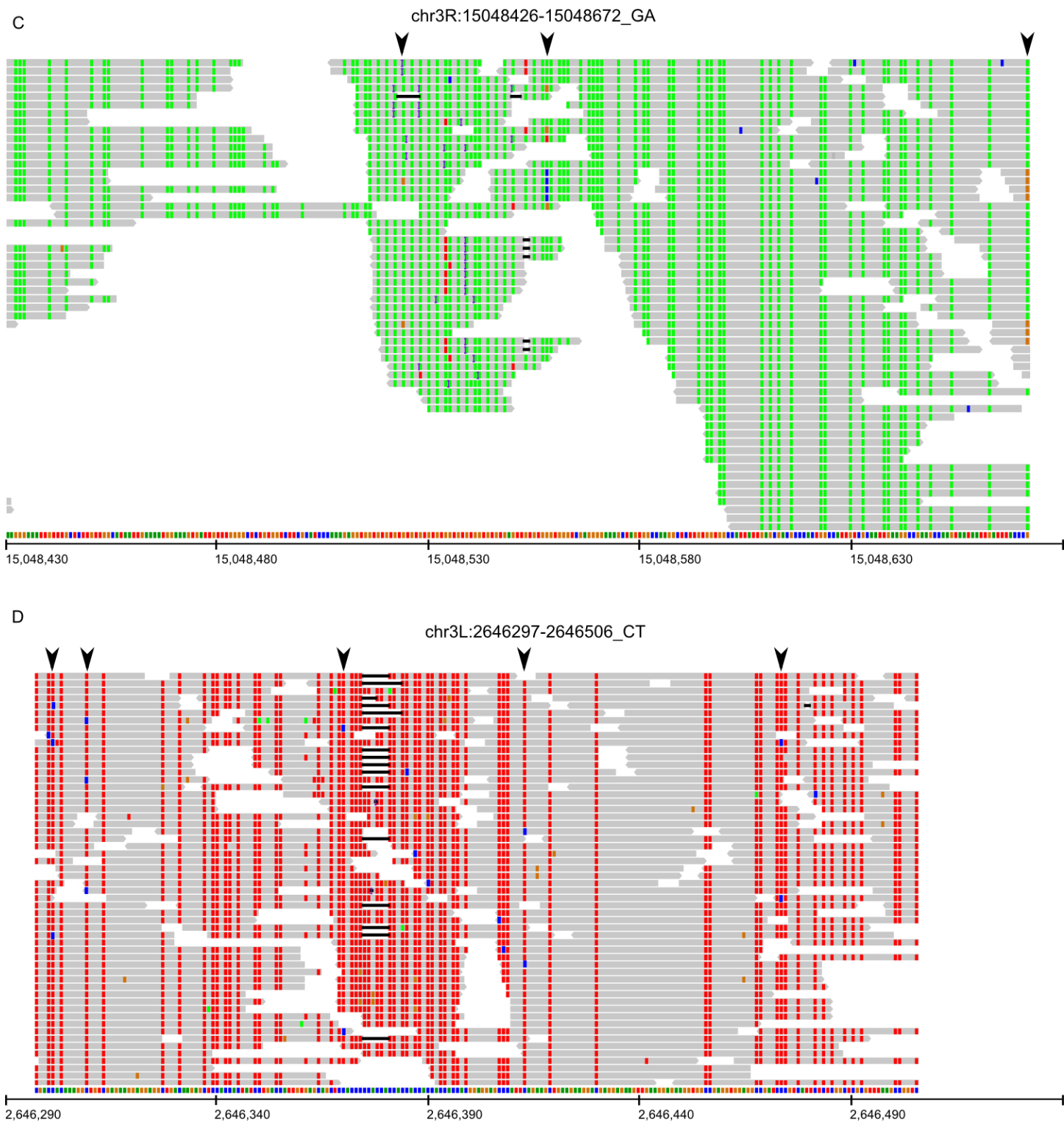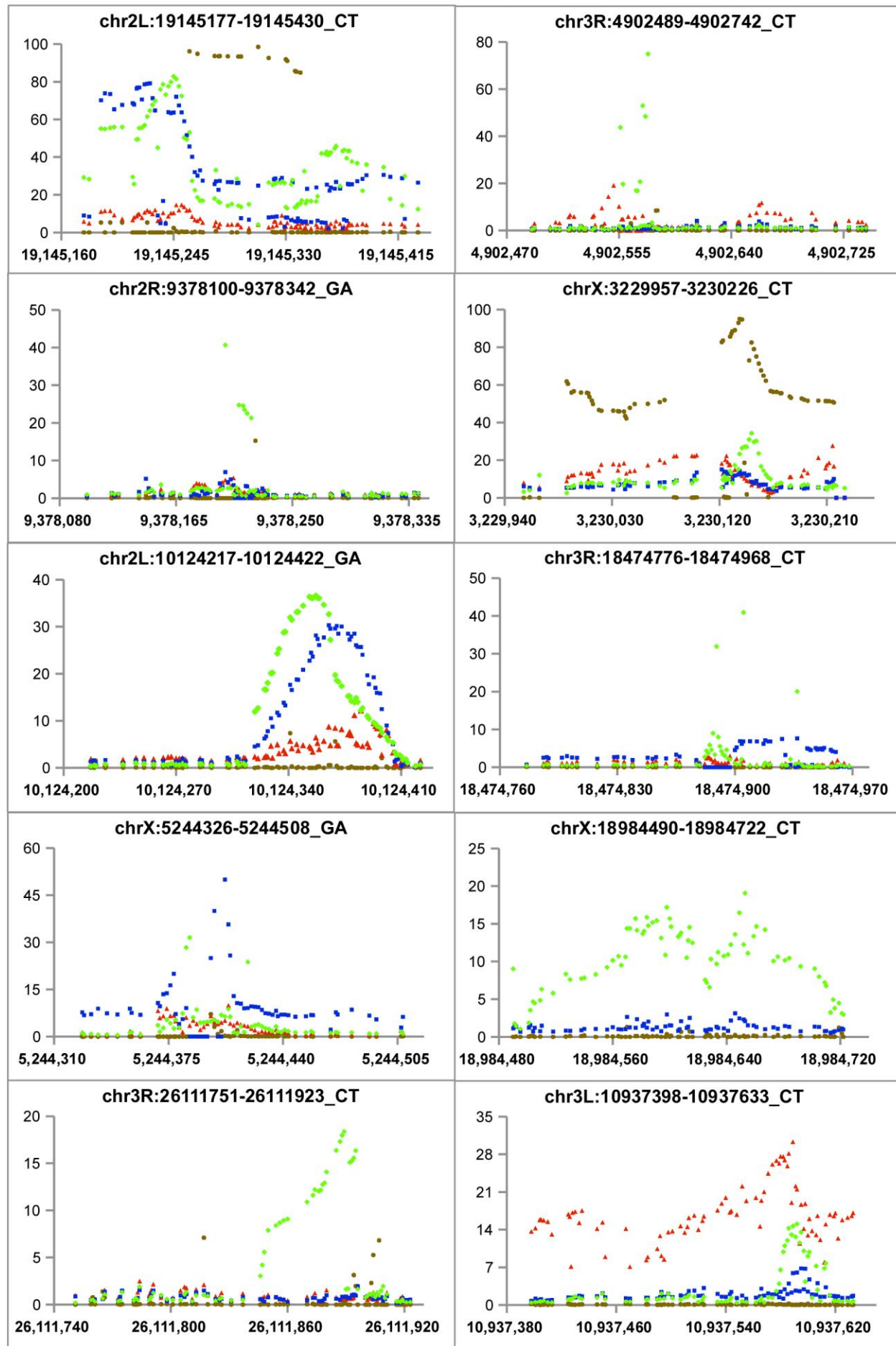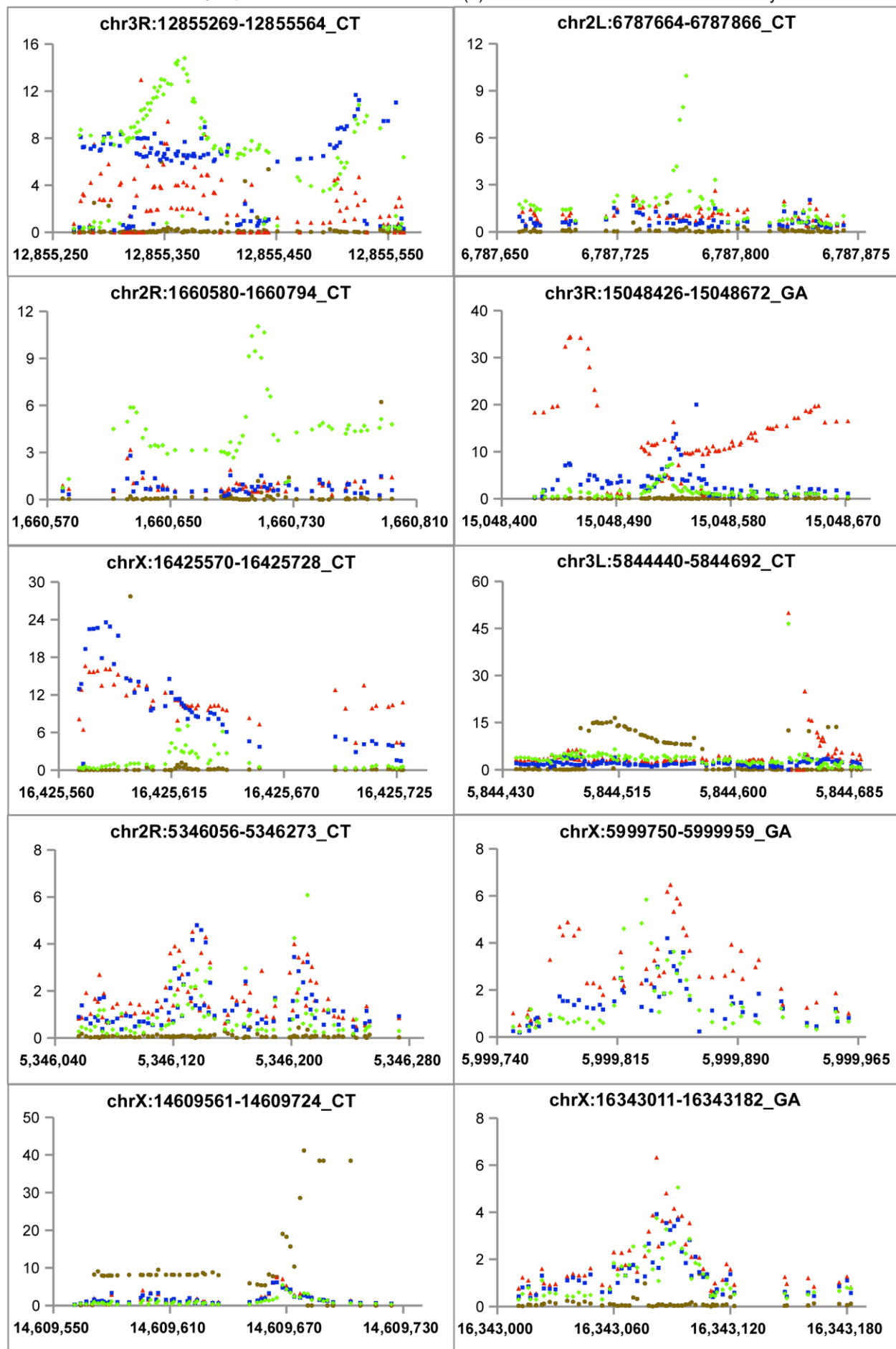
A

chrX:3229957-3230226_CT

■:A ■:G ■:T ■:C

3,229,980    3,223,030    3,223,080    3,223,130    3,223,180

B

chr3L:10937398-10937633_CT

10,937,390    10,937,440    10,937,490    10,937,540    10,937,590

**Figure S7. Reads from the dataset of Raddatz et al., aligned to the regions displayed in Figure 3.** As in Fig. S6, each line shows the alignment of a read pair; when the two sequences of a pair of reads do not overlap, a thin grey line shows their connection. The color-coded reference sequence is at the bottom of each panel, with the color key shown at the top of the figure. A match between a read and the reference is shown in grey; a mismatch is shown with the color of the mismatched base. Unmethylated cytosines are sequenced as 'T' (red) on the CT strand (panels A, B, and D) and as 'A' (green) on the GA strand (panel C). Any methylated cytosines are denoted by blue color in A, B, and D, and in brown in C. Arrows mark positions that contain at least two unconverted cytosines.

Bisulfite PCR DNMT2 KO / oocyte panel 1
OR   EP(2)GE15695   Dnmt2[99]   OR Oocyte

Bisulfite PCR DNMT2 KO / oocyte panel 2

Legend: OR, EP(2)GE15695, Dnmt2[99], OR Oocyte

Y-axis: Bisulfite PCR (% methylation)

chr3R:12855269-12855564_CT
chr2L:6787664-6787866_CT
chr2R:1660580-1660794_CT
chr3R:15048426-15048672_GA
chrX:16425570-16425728_CT
chr3L:5844440-5844692_CT
chr2R:5346056-5346273_CT
chrX:5999750-5999959_GA
chrX:14609561-14609724_CT
chrX:16343011-16343182_GA

Bisulfite PCR DNMT2 KO / oocyte panel 3

◆ OR   ■ EP(2)GE15695   ▲ Dnmt2[99]   ● OR Oocyte

Bisulfite PCR DNMT2 KO / oocyte panel 4  ◆ OR  ■ EP(2)GE15695  ▲ Dnmt2[99]  ● OR Oocyte

Bisulfite PCR DNMT2 KO / oocyte panel 5

OR  EP(2)GE15695  Dnmt2$^{99}$  OR Oocyte

Bisulfite PCR (% methylation)

chrX:11177428-11177684_CT

chr3R:13185145-13185369_CT

chr3R:3939216-3939449_GA

chr2R:16996962-16997089_CT

chr3R:26111663-26111984_GA

chrX:1440349-1440547_CT

chr2L:20708270-20708507_CT

chr2R:9642544-9642838_GA

chrX:4590301-4590472_CT

chr2R:13694092-13694285_GA

Bisulfite PCR DNMT2 KO / oocyte panel 6

Legend: ◆ OR  ■ EP(2)GE15695  ▲ Dnmt2[99]  ● OR Oocyte

Y-axis (shared): Bisulfite PCR (% methylation)

chr2R:19309681-19309890_GA
chr3R:27278112-27278313_CT
chrX:16425533-16425735_GA
chr3R:12487160-12487332_GA
chr3R:19044584-19044738_CT
chrX:4161609-4161746_CT
chrX:19302361-19302587_CT
chrX:19343642-19343824_GA
chr2R:10033350-10033586_GA
chrX:8556537-8556663_GA

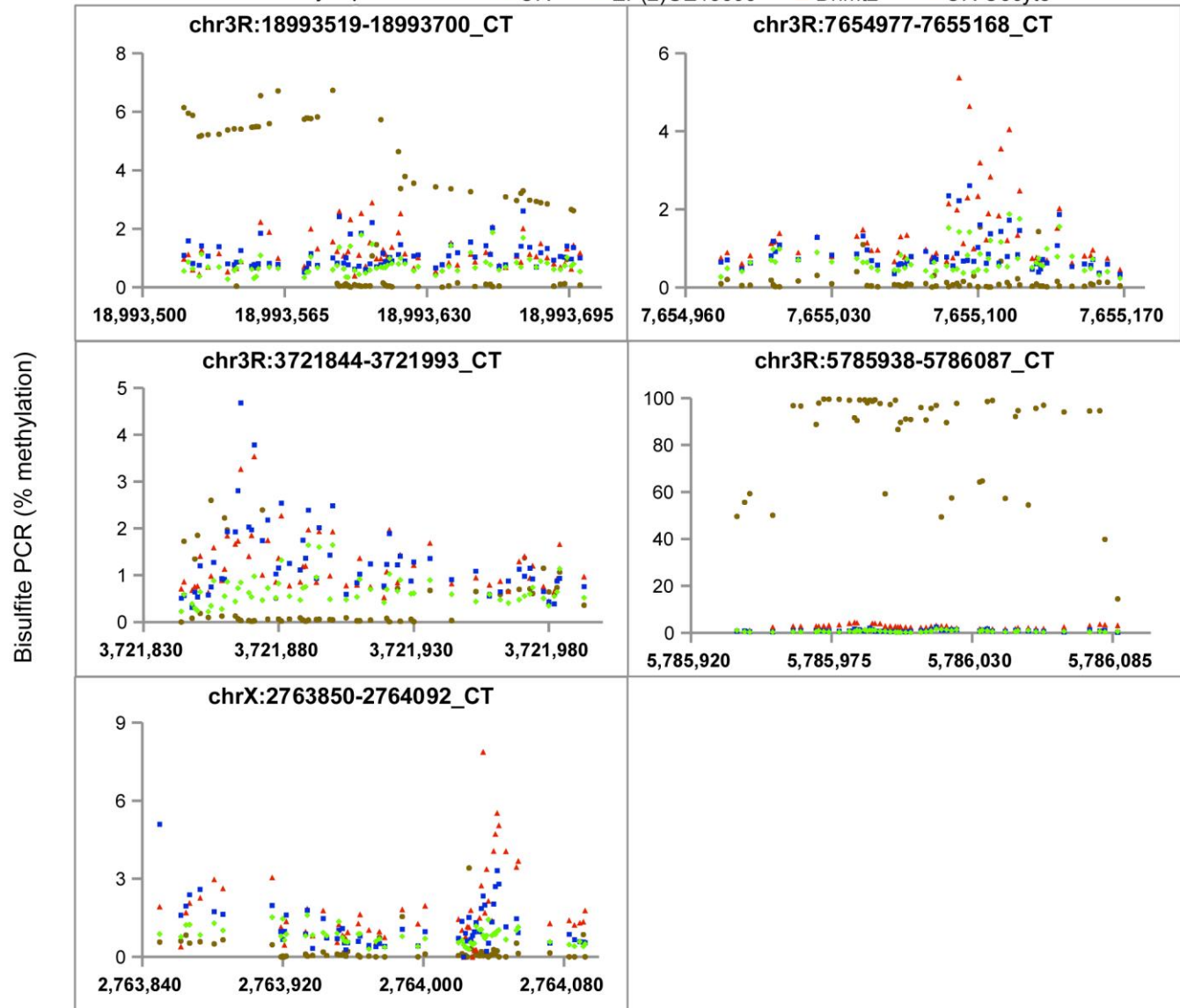**Figure S8. Methylation is present in flies deficient for the DNA methyltransferase MT2 and at some loci in unfertilized oocytes.** The full set of 66 regions that were analyzed is shown. Methylated regions identified by MeDIP-bisulfite sequencing were PCR amplified from bisulfite converted DNA and Illumina sequenced to at least 10,000X coverage. Each dot represents one cytosine: green – bisulfite PCR (same data as in Figure S5); brown – unfertilized oocyte; red – *Mt2* deficient; blue – EP(2)GE15695 (*Mt2* wild type). The y-axis indicates the percent of methylated cytosines in the bisulfite PCR.
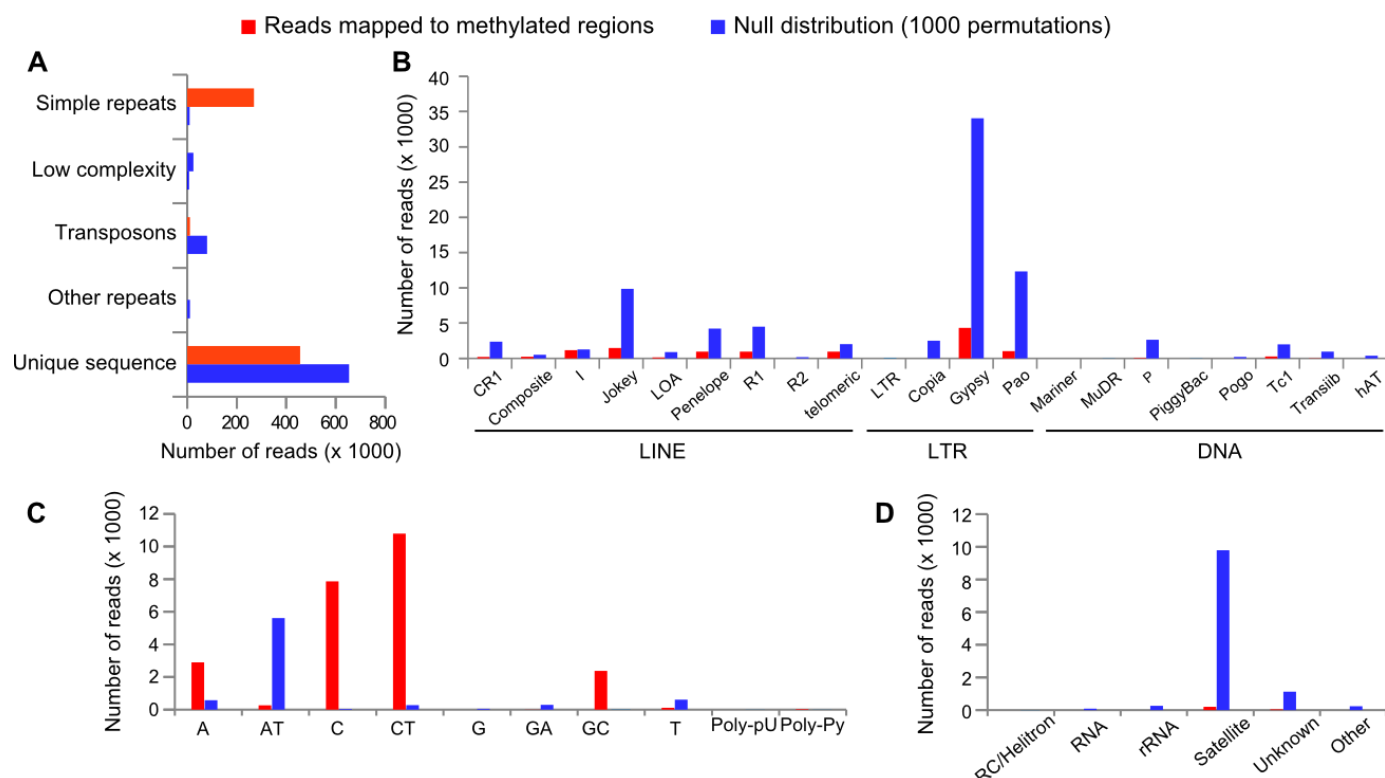
**Figure S9. Sequence properties of methylated regions**. The number of sequence reads within methylated regions that overlap with various sequence classes (red bars); blue bars represent the average and standard deviation of 1,000 randomized permutations of the same number of reads. **A.** Distribution of methylation between unique and repeat sequences. Methylation is much more likely to present in simple sequence repeats, and less likely to be present in transposons or unique sequences. **B-D.** Methylation of transposons and other repeat types. Methylation is depleted in all transposon families except the I element (B), enriched in some types of low-complexity sequence (C), and depleted from RNA, satellite, and other repeats (D). **Methods:** The repeat sequence annotation for the *D. melanogaster* dm3 assembly was downloaded from the UCSC Table Browser table:rmsk. The repeat annotation was intersected with the 762,655 primary alignments that align by at least 51% to the 25,497 methylated regions The intersection was obtained with intersectBed from the BEDTools suite (Quinlan and Hall 2010), run with -f0.51 option which requires that at least 51% of a read overlaps an annotated repeat. We used primary reads rather than methylated regions because of the difficulty in mapping a read to a specific repeat element. The results were compared to a random expectation distribution. We used shuffleBed from the BEDTools suite to randomly permute the locations of the 762,655 primary alignments. We used the –chrom option, which keeps the alignments on the same chromosome and only randomizes their location on the chromosome. A distribution of random annotations was generated by 1,000 repetitions of the permutation procedure, from which we calculated the mean and standard deviation.
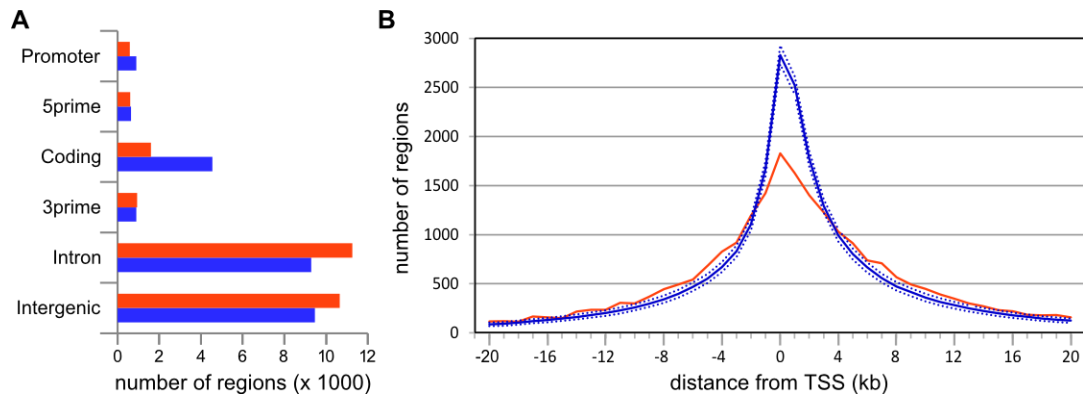
**Figure S10. Methylation of various classes of simple sequence repeats.** Red bars represent the number of sequence reads within methylated regions that overlap with a given simple sequence repeat; blue bars represent the average and standard deviation of 1,000 randomized permutations of the reads. Simple sequence repeats that lack a cytosine are not displayed. The scale on the y-axis of the top two panels is logarithmic. Few of the simple sequence repeats in which methylation is enriched contain Gs, but some of these are highly enriched.

**Figure S11. Methylated regions and genic features. A**. The number of methylated regions overlapping with annotation features in the *Drosophila* genome (red bars). The blue bars represent the average and standard deviation of 1,000 randomized permutations of the 25, 319 methylated regions. Compared to the random selection of regions, methylated regions are more likely to be found in introns and intergenic regions, and less likely to be found at promoters and within coding regions. **B.** Distance of methylated regions from the nearest transcription start site. For each of the 25, 319 methylated regions, we calculated the distance to the nearest annotated transcription start site (TSS). The red line shows the number of methylated regions at a given distance from the nearest TSS. The solid blue line indicates the mean distance to the nearest TSS of 1,000 random permutations of the genomic locations of the 25, 319 methylated regions. The dotted blue lines denote the 95.6% confidence intervals. This analysis shows a depletion of methylated regions near TSSs. **Methods:** The gene annotation for the *D. melanogaster* dm3 assembly was downloaded in BED format from the UCSC Table Browser table:flyBaseGene. Non redundant files for the various gene annotation features (promoter, 5' UTR, coding exon, 3' UTR) were obtained by collapsing all features with overlapping coordinates; a promoter was defined as the sequence up to 300bp upstream of a transcription start sites. Regions that were annotated as more than one feature (e.g.: as 5' UTR and promoter) were retained independently. Introns were defined as the sequences within a gene that did not correspond to any exon. Intergenic regions were defined as the genome sequences that did not correspond to an intron or any other gene annotation feature. The degree of overlap between methylated regions and gene annotation features was determined with intersectBed from the BEDTools suite, run with -f0.51 option, which requires that at least 51% of a methylated region overlaps a gene annotation feature. The results were compared to a random expectation distribution. We used shuffleBed from the BEDTools suite to randomly permute the locations of the 25,497 methylated regions. We used the –chrom option, which keeps the regions on the same chromosome and only randomizes their location on the chromosome. The random expectation distribution was generated by 1,000 repetitions of the permutation procedure and by intersecting each repetition with the gene annotation features.
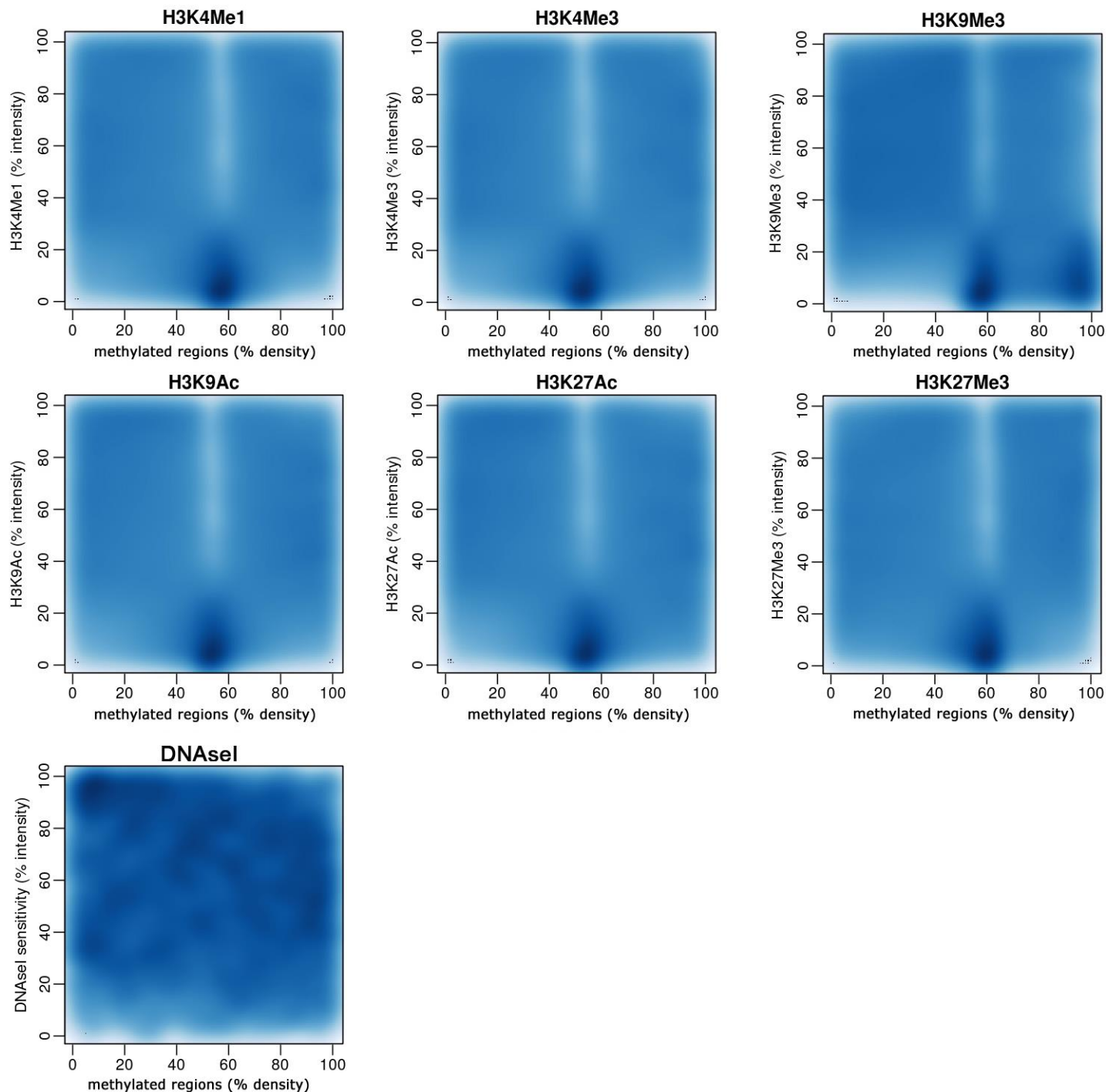
**Figure S12. Lack of correlation between methylation and chromatin features.** Scatterplots comparing the density of methylated region (x-axis) with the intensities of various histone tail modifications and of DNase I hypersensitivity (y-axis). The scatterplots illustrate a general lack of correlation. **Methods:** The density distribution of methylated regions over 100kb intervals was determined using fseq with the -l100000 -s100 options. Histone tail modification data for *D. melanogaster* (developmental stage: E0-4h) were downloaded as wiggle files from the GEO database with the accession numbers: GSM400656 (H3K4Me3), GSM401407 (H3K27Ac), GSM401408 (H3K9Ac), GSM401409 (H3K4Me1), GSM439448 (H3K27Me3), GSM439457 (H3K9Me3). DNase I sensitivity data for *D. melanogaster* (developmental stage 5) was downloaded in BED format from the UCSC Table Browser table:bdtnpDnaseAccS5, dm3 assembly, and converted to a density distribution using fseq with the -l100000 -s100 options. The intensities of the distributions of methylated regions and chromatin state features were percentile-normalized (histone tail modification data with a value of '0' were skipped during normalization) and compared with the normalized distribution of methylation density by scatterplot using the 'smoothscatter' package in R.