

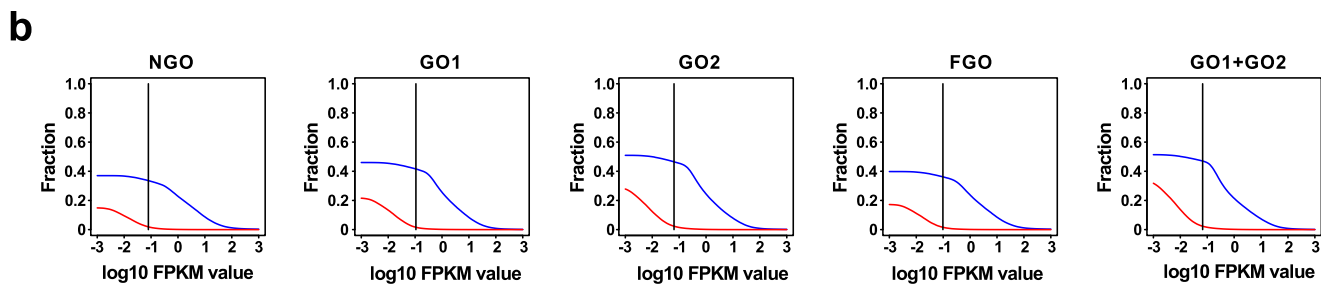
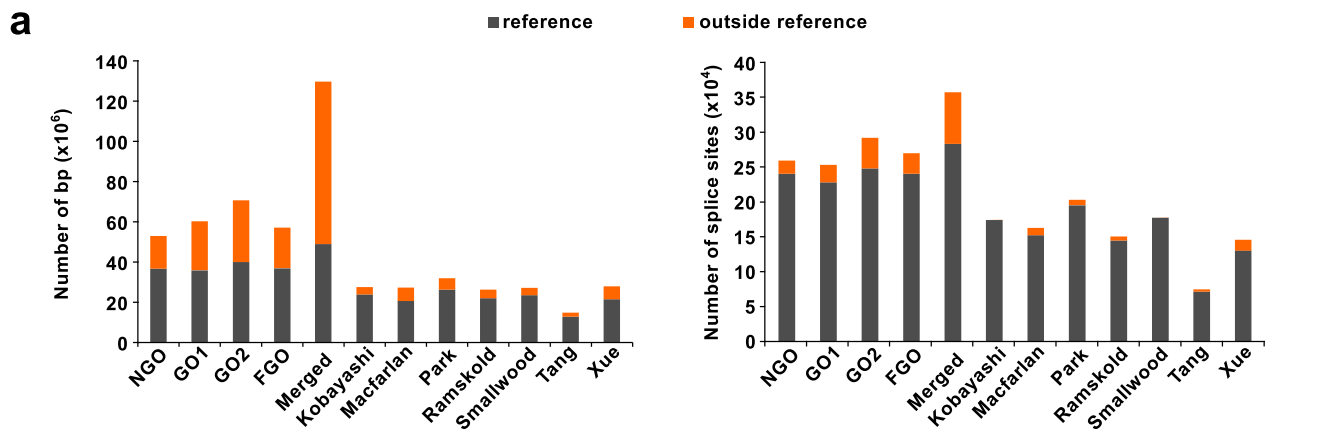
ADDITIONAL FILE1: Supplementary Figures

Deep sequencing and *de novo* assembly of the mouse oocyte transcriptome define the contribution of transcription to the DNA methylation landscape.

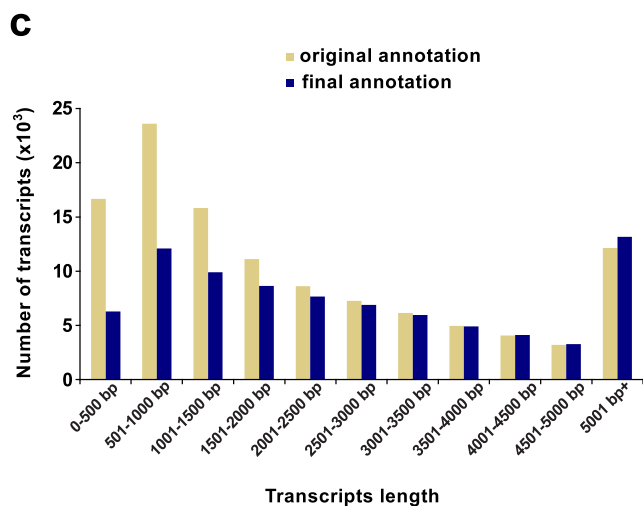
Veselovska L^{*1}, Smallwood SA^{*1}, Saadeh H^{1,2}, Stewart K¹, Krueger F², Maupetit-Méhouas S³, Arnaud P³, Tomizawa S-I⁴, Andrews S², Kelsey G^{1,5}.

Figure S1. Oocyte RNA-seq datasets and transcriptome assembly.

(a) Fraction of the genome and splice sites covered by at least 5 unique reads in our datasets (individual and merged) and individual published oocyte RNA-seq datasets (see **Additional file 2, Table S2** for the full list of datasets). **(b)** Quantification of threshold FPKM values for filtering non-transcribed transcripts from Cufflinks RABT assembler. The threshold value is visualised as the point of maximum difference (black line) between the cumulative distribution of FPKM values of transcripts in the assembly (blue) and of random size-matched intergenic regions (red). Threshold FPKMs determined using this strategy are displayed in the table below, together with, for each oocyte dataset, the proportion of transcripts exceeding this threshold (expressed and kept in the annotation), and the proportion of intergenic regions exceeding this threshold (not expressed transcripts) representing the percentage of false positives. **(c)** Length distribution of transcripts in the original oocyte annotation (beige), and after our annotation curating (blue). **(d)** Numbers of transcripts and genes in our oocyte Cufflinks annotation compared to Ensembl, RefSeq and UCSC reference annotations. Numbers of genes are as defined in Cuffcompare output.



Dataset	log FPKM threshold	FPKM threshold	Proportion of Transcripts kept	Proportion of not expressed transcripts
NGO	-1.1	0.08	33.58%	1.88%
GO1	-0.97	0.11	41.59%	1.60%
GO2	-1.19	0.07	46.48%	2.33%
FGO	-1.01	0.1	36.19%	1.60%
GO1+GO2	-1.17	0.07	47.02%	2.32%



d

	Oocyte Transcriptome	Ensembl	RefSeq	UCSC
Total No. transcripts	82,939	94,084	32,808	58,314
Toal No. genes	39,099	37,319	23,666	30,606
Multi-exonic transcripts	67,112	79,782	28,906	48,713
Mono-exonic transcripts	15,827	14,302	3,902	9,601
Multi-exonic genes	24,104	24,586	20,221	21,784
Mono-exonic genes	14,995	12,733	3,445	8,822

Figure S1

Figure S2. Novel genes and transcripts in the oocyte transcriptome.

(a) Expression of the novel transcripts identified in our assembly in the merged published oocyte RNA-seq datasets. Novel transcripts were divided into five categories according to their expression from the lowest (cat. 1) to the highest (cat. 5) level in our merged RNA-seq datasets. **(b)** Top: validation by RT-PCR of a random selection of novel multi-exonic (L1-L7) and mono-exonic (M1-M7) transcripts. For mono-exonic genes, -/+ indicates controls without reverse transcriptase. Positions of primers relative to the annotation of multi-exonic genes are indicated in the methods. DNA ladders are 100bp (hyperladder IV, Bioline) and 1kbp (hyperladder I, Bioline). Bottom: visualisation of a selection of novel mono-exonic (M4 marked by an asterisk) and multi-exonic genes and their expression in the GO2 RNA-Seq dataset. **(c)** The remaining five clusters with the smallest number of genes from Figure 2c hierarchical clustering.

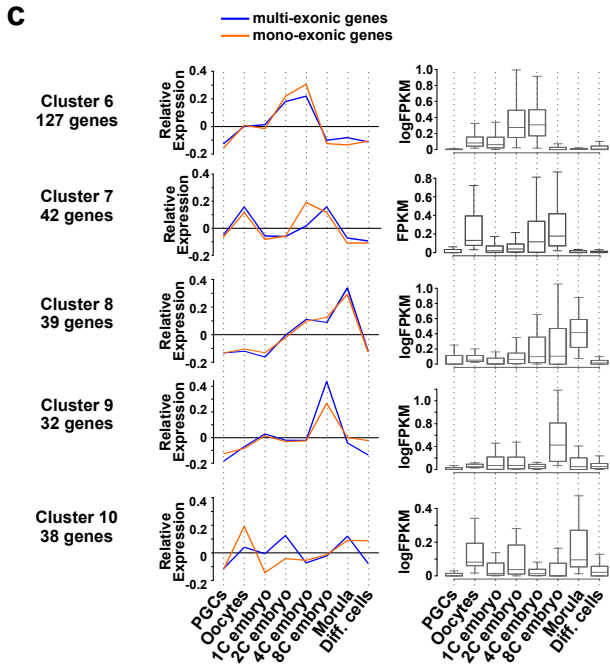
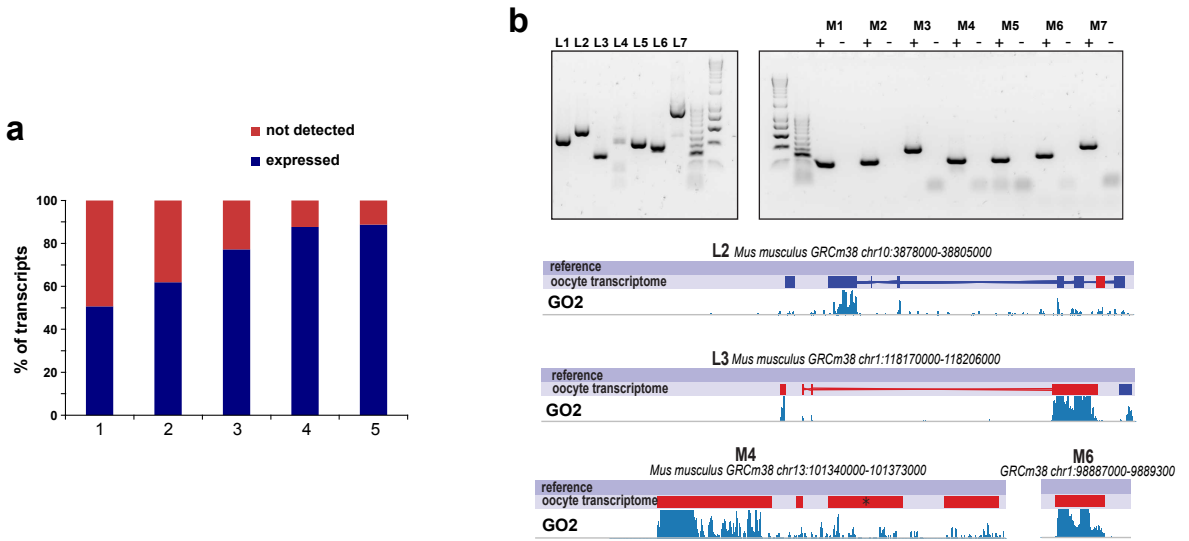


Figure S2

Figure S3. Novel upstream TSSs of reference genes in the oocyte transcriptome.

(a) Expression of novel upstream TSSs of reference genes in merged published oocyte RNA-seq datasets. Novel transcripts were divided into five categories according to their expression from the lowest (cat. 1) to the highest (cat. 5) level in our merged RNA-seq datasets. **(b)** Top: validation of a random selection of novel upstream TSSs of reference genes (T1-T12) by RT-PCR. Primers were designed to span at least one intron. Positions of primers relative to the annotation of multi-exonic genes are indicated in the methods. DNA ladders are 100bp (hyperladder IV, Bioline) and 1kbp (hyperladder I, Bioline). Bottom: visualisation of a selection of novel upstream TSSs and their expression in the GO2 dataset. **(c)** Expression of analysed TE families in our oocyte datasets, assessed by mapping RNA-Seq reads to custom repeat genomes of individual repeat families. **(d)** Expected and observed frequencies of TE families overlapping (+/- 1bp) TSSs on the same strand. Expected frequencies were calculated for the whole genome and for intergenic regions only, as TEs are more frequently found in intergenic than gene regions. Observed frequencies were calculated for (from top in the legend) all reference TSSs (reference all), the most upstream TSSs of reference genes (reference most upstream), TSSs of oocyte genes matching reference genes (oocyte matching reference), all oocyte TSSs (oocyte all), the most upstream TSSs of the oocyte genes (oocyte most upstream), all TSSs of oocyte genes after removal of transcribed independent TEs (oocyte filtered all), and most upstream TSSs of oocyte genes after removal of transcribed independent TEs (oocyte filtered most upstream). Statistically significant enrichment of TEs acting as TSSs after removal of expressed independent TEs compared to the expected frequency by chance in intergenic regions is marked by asterisk ($p < 0.0001$, chi-squared test). **(e)** Boxplot representation of the expression levels of genes with a TSS associated with TEs, not associated with TEs and associated specifically with MaLR, ERVK and other TEs, shown with and without outliers. Asterisks mark significant differences (p -value < 0.0001 , chi-squared test). **(f)** Profiles of TSS activity of reference genes with a novel upstream TSS identified in our transcriptome assembly (1560 in total), in PGCs, early embryos and somatic tissues. Percentage values on the right represent the 4 top categories, and the weight of the lines indicates the proportion of genes with each profile.

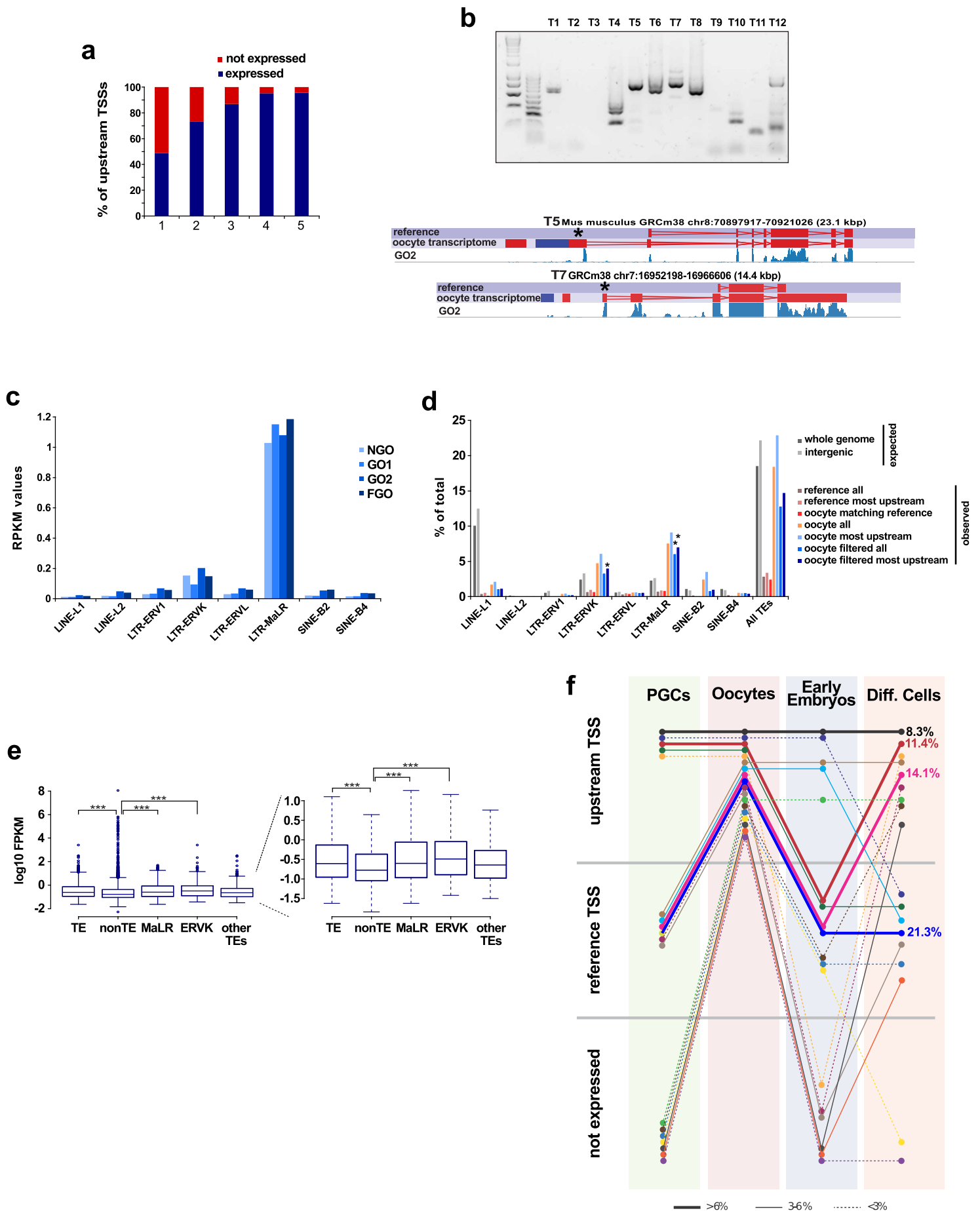


Figure S3

Figure S4. Definition of HyperDs and HypoDs.

(a) Detailed flowchart of the strategy and parameters used for the definition of HyperDs and HypoDs. The numbers in red indicate the key steps for final definition of domains. **(b)** Length of the domains after each step of domain definition as described in (a) (red numbers).

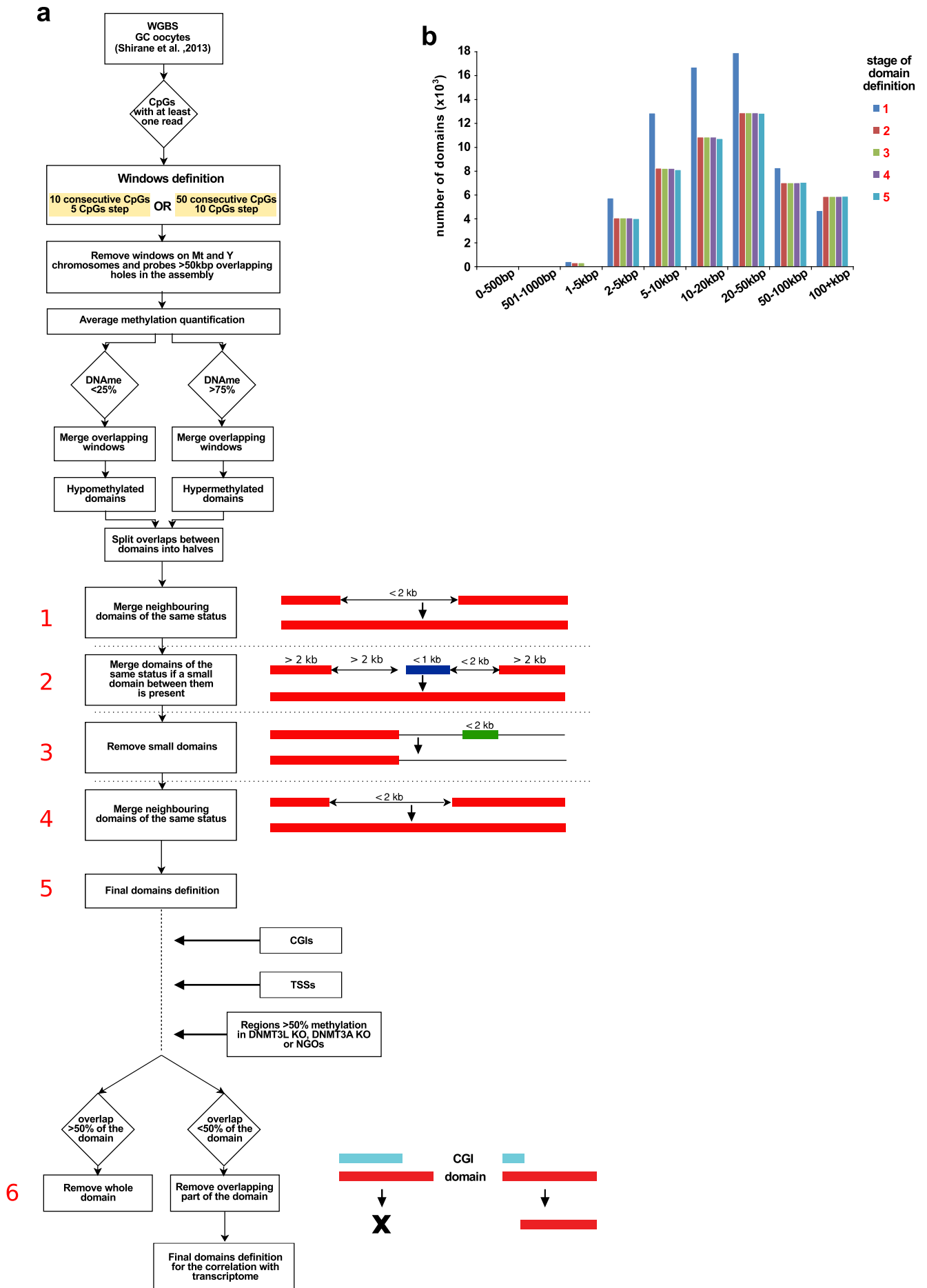


Figure S4

Figure S5. Characteristics of HyperDs and HypoDs and their correlation with transcription.

(a) Distribution of small (<10kbp) and large (>50kbp) HyperDs according their overlap with transcription in oocytes, based on the expressed reference genes (Ref. FPKM>0.001), our transcriptome assembly, our assembly combined with read contigs, our assembly/contig combined with transcribed regions of partial DNAm (>25%) in DNMT KO oocytes and NGOs and overlap with ERVK. **(b)** Proportion of 1kbp windows within HyperDs overlapping with the following features: transcripts of different FPKM values, oocyte contigs, regions within 2 kb downstream of the TTSs of the oocyte transcripts, regions with 25-50% DNA methylation in DNMT3A-deficient, DNMT3L-deficient and oocytes early NGOs (partially methylated regions) and ERVK elements. **(c)** Distribution of small (<10kbp) and large (>50kbp) HypoDs according their overlap with transcription in oocytes, based on the expressed reference genes (Ref. FPKM>0.001), our transcriptome assembly, our assembly excluding genes with FPKM \leq 0.5 alone or including also alternative TSSs. **(d)** Proportion of 1kbp windows within HypoDs overlapping with the following features: transcripts of different FPKM values and alternative TSSs. **(e)** Metagene analysis of DNAm level and distribution across gene bodies +/-5kbp as defined in our oocyte assembly, according to their expression levels. Genes were divided into 5 categories of equal numbers based on their FPKM values, from the lowest (0-20%) to the highest (80-100%).

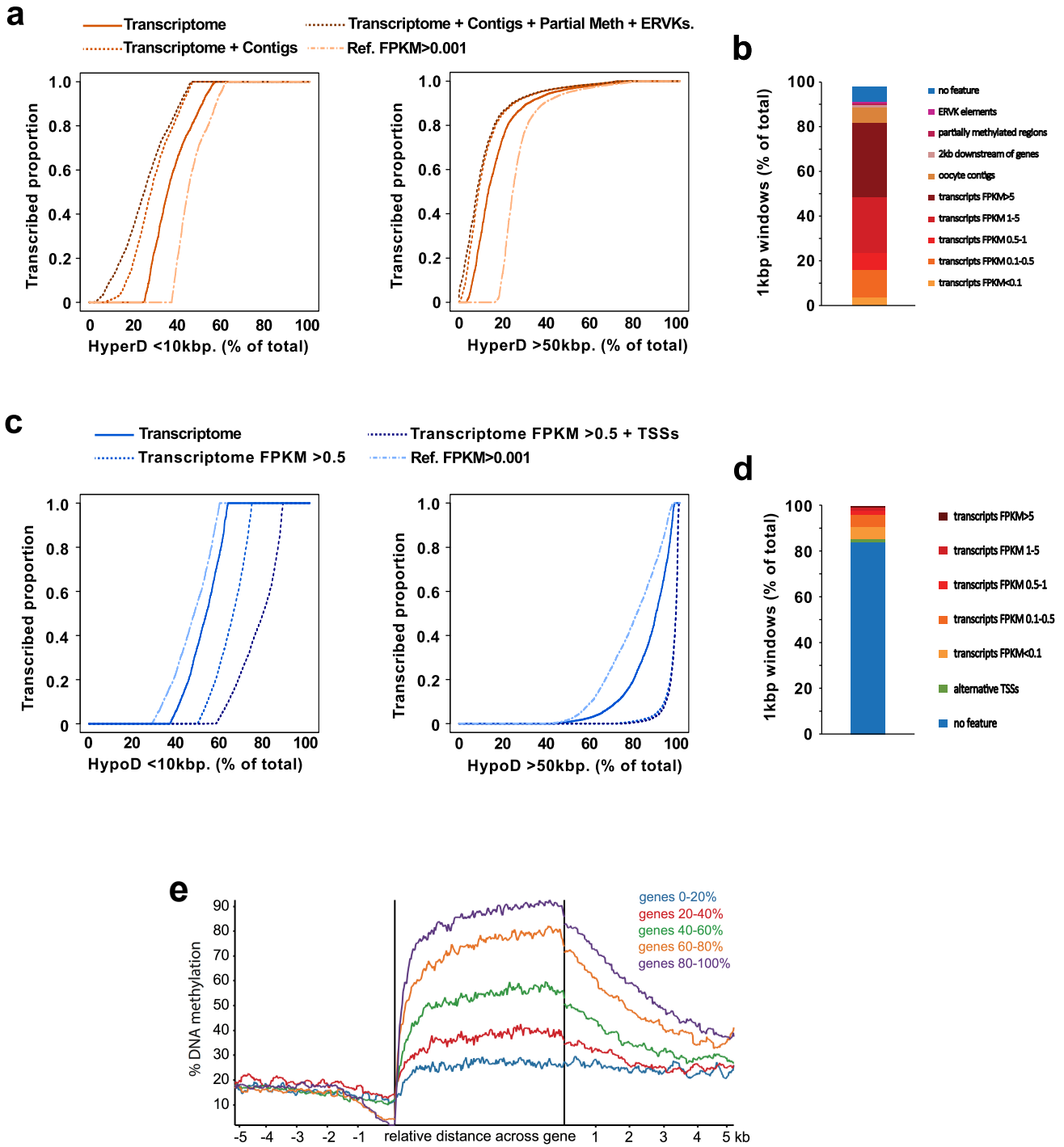


Figure S5

Figure S6. igDMRs newly identified as intragenic.

Snapshot of RNA-Seq (GO2 and FGO) and DNAm (1kbp windows, BS-Seq, FGO) of the imprinted loci for which igDMRs are newly classified as intragenic according to our transcriptome assembly. Colours of genes and RNA-seq reads represent their strand-specificity (red indicates the alignment to + strand, blue to - strand). For easier interpretation, the RNA-seq reads containing splice sites were visualised as if the read also covered the intron.

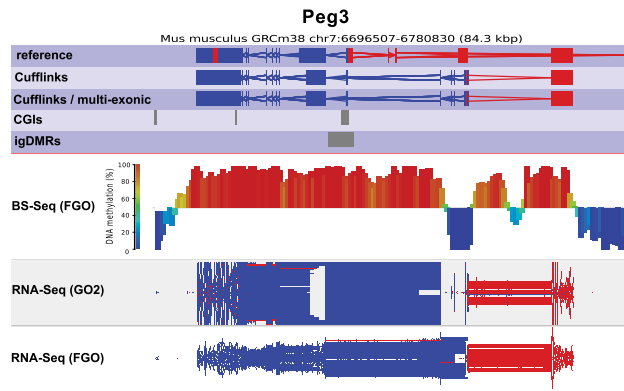
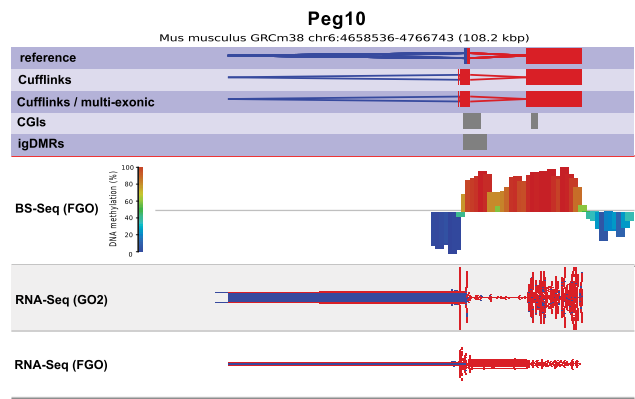
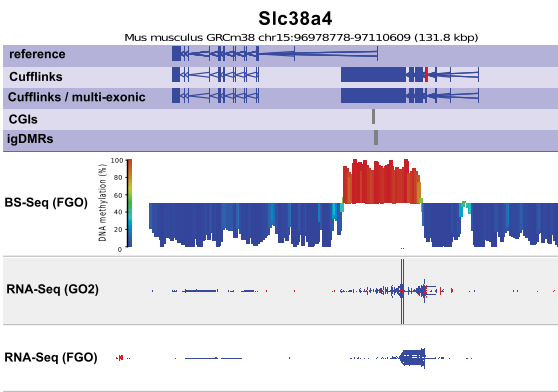
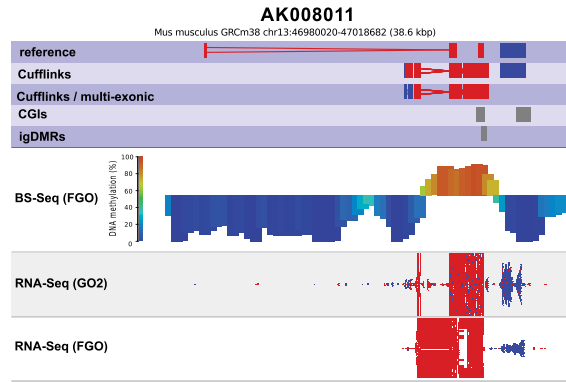
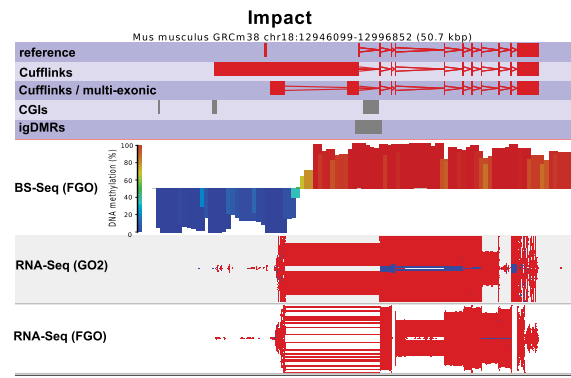


Figure S6

Figure S7. Functional test of the involvement of transcription in DNAm establishment.

(a) Scheme of the *Zac1o* conditional KO allele construct and strategy. 2 *loxP* sites were inserted 2.6kbp upstream of the *Zac1o* TSS and just downstream of *Zac1o* exon I, respectively, using homologous recombination in bacteria. The downstream *loxP* site was associated with the Neomycin selection cassette flanked by 3 *Frt* sites that was excised by crossing with Flpe mice. **(b)** Expression of *Zac1o*, *Zac1oAS* and *Ppia* (control) in *Zac1o^{fl/fl}* and *Zac1o^{fl/fl}* x Zp3-Cre (*Zac1o^{-/-}*) growing (GO) and fully-grown (FGO) oocytes. **(c)** DNAm at *Zac1* igDMR assessed by COBRA analysis in neonatal (P2) brain (Br) and neonatal (P2) pancreas (Pa) of *Zac1o^{+/+}* and *Zac1o^{+/-}* littermates from two different litters. **(d)** DNAm at *Zac1* igDMR assessed by COBRA analysis in neonatal brain (P2) of *Zac1o^{+/fl}* and *Zac1o^{+/fl}* x *Sox2-Cre* mice. **(e)** Enrichment of H3K4me3, H3K9me3, H3K9Ac and H3K20me3 at *Zac1* igDMR in *Zac1o^{+/+}* and *Zac1o^{+/-}* (LoM) mice, in neonatal brain. Bound/ input ratios were calculated and normalised to those for the imprinted *KvDMR* (***p*<0.001, Wilcoxon test) **(f)** ChIP-qPCR quantification of H3K4me2 and H3K36me3 enrichment in growing oocytes (dpp15) at *Zp3* (highly expressed), *Ppia* (highly expressed) and *Fam164b* (non-expressed) loci, in intergenic, promoter and intragenic regions, for *Zac1o^{+/+}* and *Zac1o^{-/-}* (LoM) mice (**p*<0.05, t-student test).

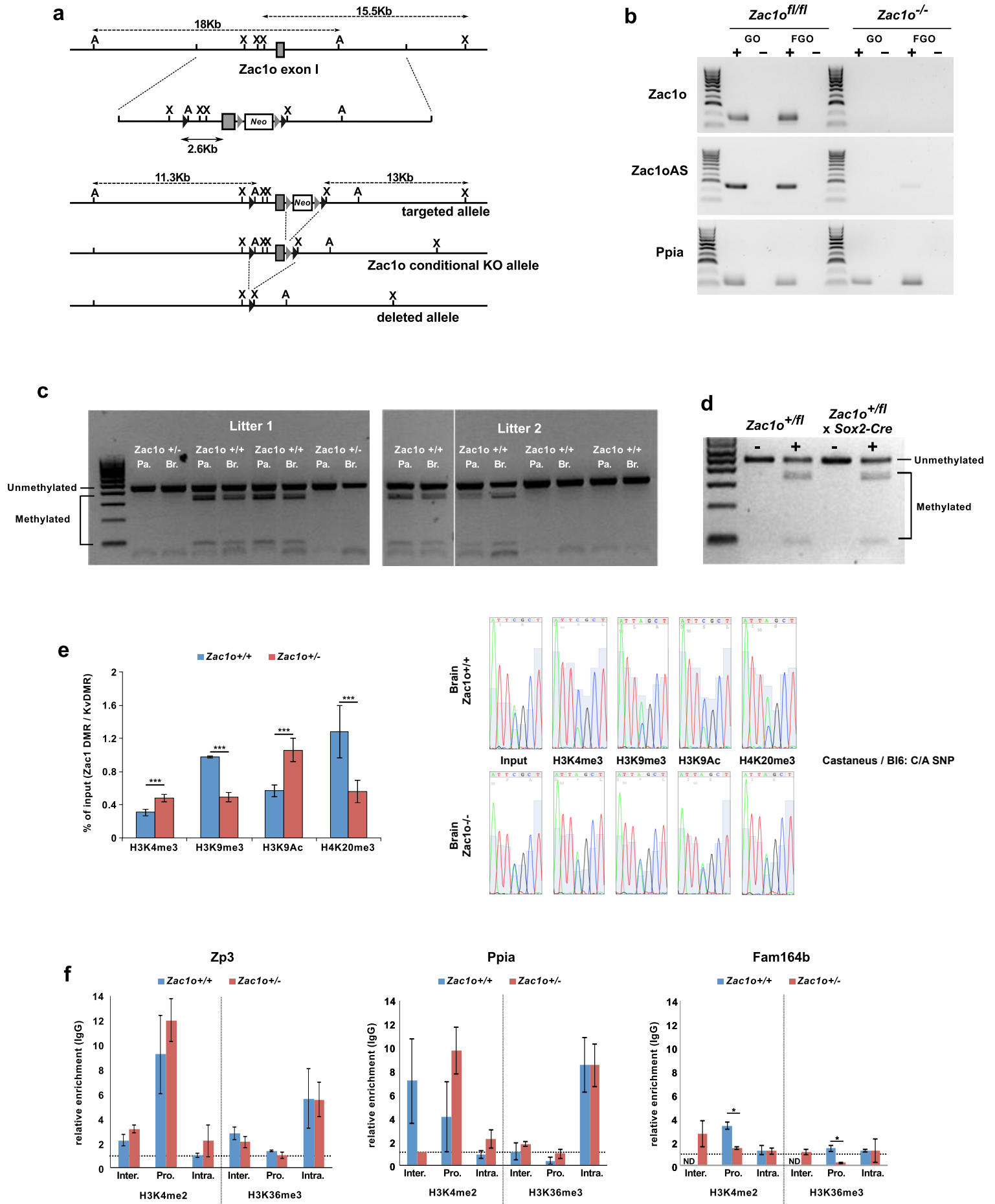


Figure S7