

## **DNA methylation programming and reprogramming in primate embryonic stem cells**

Netta Mendelson Cohen<sup>1</sup>, Vikas Dighe<sup>2\*</sup>, Gilad Landan<sup>1</sup>, Sigrun Reynisdottir<sup>3</sup>, Arnar Palsson<sup>3</sup>, Shoukhrat Mitalipov<sup>2</sup> and Amos Tanay<sup>1</sup>

1. Department of Computer Science and Applied Mathematics, The Weizmann Institute of Science, Rehovot 76100, Israel.

2. Oregon National Primate Research Center, Oregon Stem Cell Center and Departments of Obstetrics & Gynecology and Molecular & Medical Genetics Oregon Health & Science University, West Campus, 505 NW 185th Ave. Beaverton, OR 97006, USA

3. Institute of Biology and Department of Life and Environmental Sciences, School of Engineering and Natural Sciences, University of Iceland, Sturlugata 7, 101 Reykjavik, Iceland.

\* Present address: National Institute for Research in Reproductive Health (ICMR), Parel Mumbai, India.

Correspondence: Amos Tanay ([amos.tanay@weizmann.ac.il](mailto:amos.tanay@weizmann.ac.il)). Phone: 972-8-934-3579, Fax: 972-8-9348375

## ABSTRACT

DNA methylation is an important epigenetic mechanism, affecting normal development and playing a key role in reprogramming epigenomes during stem cell derivation. Here we report on DNA methylation patterns in native monkey embryonic stem (ES) cells, fibroblasts and ES cells generated through somatic cell nuclear transfer (SCNT), identifying and comparing epigenome programming and reprogramming. We characterize hundreds of regions that are hyper- or hypo-methylated in fibroblasts compared to native ES cells and show that these are conserved in human cells and tissues. Remarkably, the vast majority of these regions are reprogrammed in SCNT ES cells, leading to almost perfect correlation between the epigenomic profiles of the native and reprogrammed lines. At least 58% of these changes are correlated *in cis* to transcription changes, Polycomb Repressive Complex-2 occupancy, or binding by the CTCF insulator. We also show that while epigenomic reprogramming is extensive and globally accurate, the efficiency of adding and stripping DNA methylation during reprogramming is regionally variable. In several cases, this variability results in regions that remain methylated in a fibroblast-like pattern even after reprogramming.

## Introduction

DNA methylation is considered a key factor in the formation of cellular memory and identity, but due to experimental and conceptual limitations we still do not truly understand how the cell writes and erases DNA methylation marks in the course of normal cellular differentiation, and how these marks revert to their original embryonic stem cell (ESC)-like form following somatic cell nuclear transfer or iPS reprogramming (Reik 2007). Progress in the field was hampered for years by lack of quality methods for high throughput DNA methylation profiling, but recently several effective assays for profiling DNA methylation in large fractions of the mammalian genome were developed and applied successfully (Weber, Davies et al. 2005; Keshet, Schlesinger et al. 2006; Rollins, Haghghi et al. 2006; Cokus, Feng et al. 2008; Irizarry, Ladd-Acosta et al. 2008; Meissner, Mikkelsen et al. 2008). Another major source of confusion and difficulty in understanding the role of mammalian DNA methylation is the non-uniform CpG content of the genome, which led most of the experimental attention toward regions with the highest CpG content (CpG islands). Recent evidence suggests that classical CpG islands of high CpG content are almost never methylated in normal conditions, but much dynamic DNA methylation (manifested as differentially methylated regions (DMRs)) can be found in regions with intermediate CpG content, some of which are classically defined as CpG islands and some of which are not (Irizarry, Ladd-Acosta et al. 2009; Strausman, Nejman et al. 2009). Adding to these difficulties, multiple studies have shown that DNA methylation is stably acquired in culture, forming significant line-to-line variability (Allegrucci, Wu et al. 2007) and deterministic tissue culture effects (Mikkelsen, Hanna et al. 2008; Brunner, Johnson et al. 2009), all of which make the interpretation of the functional role of DNA methylation difficult to verify. Culture effects and variability are of particular importance when analyzing DNA methylation in stem cells and induced differentiation, since the compatibility of stem cell lines with various clinical applications may greatly depend on their epigenomic state. Taken together, the recent experiments in the field have completely changed the way by which DNA methylation is studied, yet left many challenges unresolved.

We sought to approach these challenges using a suite of rhesus monkey (*Macaca mulata*) stem cell lines that allowed us to compare epigenetic programming (reorganization of DNA methylation during normal differentiation) and reprogramming (reorganization of DNA methylation patterns following derivation of stem cells from somatic cells) (**Fig. 1A**). The comparison of these two reciprocal processes and the

use of ESC lines of markedly different developmental origin allowed for better control of ESC line and culture effects. Primate somatic nuclear transfer cells are currently unique to rhesus monkeys, and represent an opportunity to study DNA methylation patterns in a native reprogramming environment. We used a native in-vitro fertilization ESC line (ORMES-22), a primary XY fibroblast line and the CRES-2 line generated through somatic cell nuclear transfer from the fibroblast line (Byrne, Pedersen et al. 2007). We supplemented our panel with a homozygous parthenote (unfertilized, spontaneously diploid and active) ESC line (ORMES-9), which provided an additional control against ESC-line specific effects. We obtained gene expression data from each cell type, and performed MeDIP-chip (Keshet, Schlesinger et al. 2006; Mohn, Weber et al. 2009) using tiling arrays designed to encompass rhesus regions orthologous to human ESC H3K4me3-H3K27me3 bivalent domains (Bernstein, Mikkelsen et al. 2006) with additional extensive control regions. For each cell type, data from three biological replicates was averaged. The array provided us with comprehensive quantitative data on key genomic regions of diverse CpG contents and developmental relevance, something that is still difficult to achieve using alternative technologies (Cokus, Feng et al. 2008; Meissner, Mikkelsen et al. 2008).

We first assessed the overall degree of reprogramming in the CRES-2 line. Strikingly, although the fibroblasts and native ES cells showed considerable differences in methylation patterns (see below), these differences were almost completely reversed upon reprogramming, generating a CRES-2 DNA methylation pattern that is highly similar to that of the native stem cell (correlation between differential methylation,  $\rho=0.53$ ,  $P \ll 10^{-100}$ , **Fig 1B**, similar to correlation between biological replicates, **Fig S1**). As demonstrated in the HOXA and HOXD loci (**Fig 1C**), specific genomic regions are subject to different DNA methylation dynamics, including gain or loss of methylation in fibroblasts compared to ES cells (we denote these hyper- and hypomethylated regions, respectively). Interestingly, even though reprogramming is globally accurate, some relatively rare domains remain methylated in a fibroblast-like pattern in the CRES-2 reprogrammed line (marked as "Failed" in Fig 1C).

A set of differentially methylated regions (DMRs) was then extracted using a statistical procedure that searched blindly for significant methylation differences between any two of the three lines (native ESC, somatic, reprogrammed) (Methods). The median methylation values of each of the lines at each of the DMRs were clustered to provide an unbiased view of the global methylation dynamics in different

groups of loci. As shown in **Fig 2A**, we observed a similar number of hypo-DMRs (loss of methylation in fibroblasts compared to ESCs, 391 regions), and hyper-DMRs (gain of methylation in fibroblasts compared to ESC, 331 regions). Distributions of differential methylation values in hypo- and hyper-DMRs are shown in **Fig S2**. Importantly, in the vast majority of DMRs, the methylation levels in all three ES lines (including the reprogrammed CRES-2) were equivalent. A smaller group of 97 DMRs was characterized by high levels of methylation in the native ESCs and lower levels in the reprogrammed line. These DMRs either represent failure to reprogram the CRES-2 line or ESC culture hyper-methylation, but importantly they constitute only a minority of the detected DMRs, which we analyzed separately.

The natural grouping of DMRs into regions undergoing gain or loss of DNA methylation during differentiation is further supported by genomic properties of these loci. Regions gaining DNA methylation are typically larger than regions losing it ( $P < 6 \times 10^{-9}$  (KS), **Fig 2B**), the latter having a well defined distribution of lengths with mean around 2kb, suggesting association with more spatially-defined genomic elements. Furthermore, regions gaining DNA methylation have lower overall CpG content than regions losing DNA methylation ( $P < 4 \times 10^{-6}$  (KS), **Fig 2C**), and both groups are generally of much lower CpG content than classical CpG islands. Our analysis therefore suggests the existence of a group of large regions with low CpG content that are unmethylated in ESCs, gain methylation in fibroblasts, and are capable of losing it upon reprogramming. The data also suggests the existence of well-localized (1-3kb) regions with intermediate CpG content and high levels of ESC methylation, which lose methylation in fibroblasts and regain it following reprogramming. Additional differences between hyper- and hypo- DMRs are described in **Fig S3** and **Fig S4**.

How important are the DMRs we have characterized for programming and reprogramming? One cannot rule out the possibility that some of the hypo-DMRs represent culture effects that are accumulated deterministically and independently in the three ES lines we have analyzed, generating hypermethylation in regions that are normally never methylated. Similarly, it is possible that some of the hyper-DMRs represent accumulation of DNA methylation in the fibroblast culture. To further describe the universality and robustness of the monkey DMRs, we computed the differential DNA methylation between human ESCs and muscle tissues (Straussman, Nejman et al. 2009) in regions of the human genome that are orthologous to monkey DMRs and regions of high and low methylation. The data (**Fig 3A**) demonstrated very

good conservation of the monkey DMR methylation patterns, where hyper-DMRs have higher methylation in human muscle tissues ( $P < 2.2 \times 10^{-16}$ ) and hypo-DMRs have higher methylation in human ESCs ( $P < 0.0005$ ). Furthermore, analysis of the range of methylation of our DMRs across a panel of four human tissues (spleen, liver, colon and brain, (Irizarry, Ladd-Acosta et al. 2009)) shows that monkey DMRs, which are defined based on comparisons of fibroblasts and ESCs, are significantly more likely to have variable methylation in the human tissues ( $P < 10 \times 10^{-8}$  (hypo-DMRs),  $P < 6 \times 10^{-5}$  (hyper-DMRs), **Fig 3B**). These lines of evidence suggest that at least some of our DMRs are real targets of methylation changes during programming and reprogramming, and that culture effects cannot explain all of the epigenomic changes we observe between ESCs and fibroblasts.

Based on the observations on robustness and conservation of monkey DMRs in human, we next examined the extent to which hyper- and hypo-methylated DMRs are correlated with changes in gene expression or occupancy of epigenomic marks and regulatory factors. Using gene expression array data, we identified rhesus genes with induced or repressed expression in fibroblasts compared to ESCs. We then mapped array probes with gain or loss of DNA methylation to these regulated genes and computed the enrichment of hypo- and hyper-methylated probes around induced and repressed transcription start sites (TSSs). As shown in **Fig S5**, the well documented (Weber, Hellmann et al. 2007; Gal-Yam, Egger et al. 2008) anti-correlation between gene expression changes and DNA methylation changes is observed, where induced TSSs tend to lose DNA methylation and repressed TSSs gain them. Nevertheless, TSS-related changes in DNA methylation account for less than a quarter of the observed DMRs, even when assigning DMRs to TSSs as distant as 5kb and relaxing our definition of gene induction or repression. We note that the annotation of the monkey genome, which is based on mapping of known human genes, is clearly incomplete, but that such incompleteness is likely to affect only a small fraction of the DMRs. We must conclude that changes in DNA methylation during programming and reprogramming are only partially associated with changes in gene expression *in cis*, and that other factors also contribute to modify DNA methylation patterns from their ESC pattern to a somatic pattern and back, either affecting transcription through long range interactions or global epigenomic reorganization, or not affecting transcription at all.

We and others have previously proposed that Polycomb occupancy in ES cells predisposes genomic regions to retain high CpG content during evolution (Tanay,

O'Donnell et al. 2007) or to gain DNA methylation in cancer cells (Ohm, McGarvey et al. 2007; Schlesinger, Straussman et al. 2007; Widschwendter, Fiegl et al. 2007; Gal-Yam, Egger et al. 2008). We therefore computed the distribution of human ESC Suz12 occupancy (Lee, Jenner et al. 2006) in rhesus-mapped regions with low or high ES methylation and in DMRs. In accordance with previous reports, we observe a general deficit of Polycomb occupancy in regions of high methylation (across CpG contents classes) (Mikkelsen, Ku et al. 2007; Fouse, Shen et al. 2008). Moreover, hyper-DMRs in medium and high CpG contents are strongly enriched for high Suz12 occupancy in human ESCs (**Fig. 4A**). Surprisingly, Suz12 enrichment is also detected at hypo-DMRs, suggesting that a significant minority of the Polycomb targets in ES cells sustain significant levels of DNA methylation and are predisposed to lose these upon differentiation. Polycomb complexes are therefore correlated with DNA methylation programming and reprogramming as previously suggested, but their role may be rather heterogenic, either passive (by blocking de-novo or housekeeping methylation) or active (by promoting methylation or demethylation upon specific regulatory queues). Alternatively, other underlying uncharacterized epigenetic factors may facilitate DNA methylation changes in epigenetic hotspots, generating indirect correlation with Polycomb occupancy in these regions.

In search of additional factors involved in DNA methylation reprogramming we studied the genomic distribution of the CTCF genomic insulator protein using data from human fibroblasts (Kim, Abdullaev et al. 2007). CTCF is known to have DNA methylation-dependent activity at several key loci, including the H19 imprinting control DMR. CTCF binding is characterized by a highly specific and informative DNA binding motif, making it one of the mammalian DNA binding proteins with the highest in-vivo sequence specificity. Nevertheless, sequence based prediction of CTCF binding is still less than 50% accurate (Kim, Abdullaev et al. 2007). As shown in **Fig 4B**, DNA methylation can account for much of this limited specificity, since the distribution of CTCF binding in methylated regions with a CTCF binding site is essentially the same as that of regions without a CTCF binding site. On the other hand, CTCF binding in unmethylated CTCF binding sites is significantly higher ( $P < 4 \times 10^{-5}$  (KS)), representing a highly specific binding distribution. This supports the mutual exclusion of CTCF binding activity and DNA methylation, and suggests that some of the DMRs we detected may be related to changes in CTCF occupancy. Indeed, as shown in **Fig 4C**, many hypo-DMRs, but fewer hyper-DMRs, have an underlying CTCF binding site, suggesting a possible link between CTCF recruitment and DNA hypomethylation in fibroblasts, and between the re-methylation of these

sites upon reprogramming and CTCF loss. We did not detect significant correlation between DNA methylation changes and the pluripotency factors Oct4 and Nanog, as shown **Fig S6**.

We have considered several factors that may be correlated with DNA methylation changes during differentiation and de-differentiation. These include (i) the machinery that drives activity at TSSs, (ii) Polycomb complexes and (iii) CTCF factors. In all three cases, occupancy of some specialized protein complexes is generally excluded from regions of high DNA methylation. Changes in DNA methylation may therefore depend (directly or indirectly) on the recruitment or loss of these protein complexes. Using conservative thresholds (Methods), a systematic analysis suggests that 239 out of 391 (61%) hypo-DMRs and 179 out of 331 hyper-DMRs (54%) are correlated with at least one of these factors. Only a minority of the DMRs are associated with regulated TSSs (32/391 hypo-DMRs, 23/331 hyper-DMRs). CTCF sites are associated mostly with hypo-DMRs (106/391) (**Fig 4D**). We note that these numbers are based on comparison of monkey methylation data and human CTCF and Suz12 profiles, so we may be underestimating the overlaps in cases of evolutionary divergence.

A minority of the DMRs we detected show a consistent native ESC pattern (in two ESC lines) that differs from a pattern common to the fibroblasts and the reprogrammed ESC line (**Fig 5A**). These cases potentially represent failure to reprogram the DNA methylation patterns of the somatic cells to their ESC templates. To systematically quantify the efficiency of reprogramming we computed the ratio between the difference in the methylation median of the reprogrammed ESC and the fibroblast and the difference in the methylation median of the native ESC and the fibroblast (denoted as the *DMR reprogramming ratio*). A reprogramming ratio of one represents perfect reprogramming and smaller values represent imperfect reprogramming. As shown in **Fig 5B**, the distribution of reprogramming ratios for hyper-DMRs is centered near 1 (Median=0.84), with general tendency to values lower than one, but very few cases near zero. Reprogramming for hyper-DMRs is therefore close to perfect on average. In contrast, the distribution of reprogramming ratios for hypo-DMRs indicate poorer overall reprogramming (Median=0.70,  $P < 10^{-8}$  (KS), and several cases that partially or completely lack reprogramming (ratio<0.25, n=27). One possibility is that the distribution of reprogramming ratio reflects an ongoing process of ESC hypermethylation during reprogramming, which occurs at very different rates for different regions, making slowly reprogrammed regions appear



non-reprogrammed and rapidly reprogrammed regions appear perfectly reprogrammed. Another possibility is that reprogramming is terminated, or never occurs, in some of the low reprogramming-ratio DMRs. Both scenarios result in imperfections of the reprogrammed epigenome, but the implications on stem cell biology remain unclear. We did not detect significant systematic correlations between low reprogramming-ratio and other genomic features (**Fig S7**). Since our methylation profile covers selected parts of the genome, it is possible that additional DMRs are slowly or improperly reprogrammed during somatic cell nuclear transfer, and this may also be the case for the epigenomic state of stem cells derived by induction of pluripotency factors (Takahashi and Yamanaka 2006; Mikkelsen, Hanna et al. 2008).

Our experiments and analysis, together with other recent measurements of DNA methylation in mouse and human ES cells (Farthing, Ficz et al. 2008; Fouse, Shen et al. 2008; Meissner, Mikkelsen et al. 2008; Deng, Shoemaker et al. 2009), differentiated cell lines (Meissner, Mikkelsen et al. 2008; Mohn, Weber et al. 2008), and somatic tissues (Rakyan, Down et al. 2008; Irizarry, Ladd-Acosta et al. 2009; Straussman, Nejman et al. 2009), outline a rather dynamic picture of the DNA methylation landscape. Comparing fibroblasts to ES cells, a large number of regions are either hyper- or hypo- methylated, most of which have medium CpG content. Such changes in DNA methylation are correlated with changes in TSS activity, Polycomb occupancy or CTCF occupancy for at least 58% of the cases we have profiled. We hypothesize that for the remaining regions, other protein complexes, or more accurate information on the current protein complexes, may account for the observed methylation dynamics. Interestingly, very little dynamic is observed in the methylation of high CpG content CpG islands, which are generally devoid of methylation in ES cells and differentiated cells. These CpG islands are very frequently located next to developmental regulators and are occupied by Polycomb complexes in ES cells. Many of these CpG islands are aberrantly methylated in cancer, but we have not found significant data suggesting their normal hypermethylation in the present study or any of the other recent high throughput studies. We therefore believe that future experiments quantifying DNA methylation programming/reprogramming should carefully distinguish between different classes of CpG rich regions and avoid focusing on promoters or CpG islands alone. Finally, the epigenome of reprogrammed ES cells is shown here to follow closely that of native ES cells, matching the striking similarity of gene expression in native and SCNT stem cells (Byrne, Pedersen et al. 2007) or iPS cells (Takahashi and Yamanaka 2006). However, we detect several exceptions to this general trend.

Further analysis of these exceptions should clarify whether the failure to reprogram specific genomic domains has functional consequences. Slow- or limited-reprogramming DMRs may also serve as key examples to contrast those many regions that are reprogrammed efficiently, leading to better understanding of the epigenomics of stem cells and the dynamics of DNA methylation and demethylation in general.

**Accession numbers.** DNA methylation profiles were deposited in GEO, accession GSEXXXX (submission in process).

**Acknowledgments.** We thank E. Segal, E. Meshorer and I. Simon for critical reading of the manuscript and to members of the Tanay lab for discussion. Research in the lab of AT was supported by ISF grant 1372/08 and by the ISF converging technologies program. AT is an Alon fellow. Funding in SM lab was provided by start-up funds from Oregon National Primate Research Center and Oregon Stem Cell Center, grants from Stem Cell Research Foundation and NIH NS044330, HD18185 and RR00163. The AP lab is supported by the University of Iceland research council, the Icelandic Centre for Research (Project: 070260022) and Marie Curie International reintegration grant (MIRG-CT-2007-46510). We thank Steingrimur Einarsson and his colleagues at Nimblegen for technical advice.

## METHODS

*Stem cells culture.* Methods for isolation and culture of monkey ES cells from IVF, SCNT and parthenogenetic embryos used in this study were reported previously (Mitalipov, Kuo et al. 2006; Byrne, Pedersen et al. 2007; Dighe, Clepper et al. 2008). Briefly, ES cells were grown on feeder layers (mouse embryonic fibroblasts, mEFs) in DMEM/F12 medium with glucose and without sodium pyruvate, supplemented with 1% nonessential amino acids, 2 mM L-glutamine, 0.1 mM  $\beta$ -mercaptoethanol and 15% FBS at 37°C, 3% CO<sub>2</sub>, 5% O<sub>2</sub> and 92% N<sub>2</sub>. Culture medium was changed daily and the ES cell colonies were typically split every 5-7 days by manual dissociation and the collected clumps were replated onto fresh mEFs.

*MeDIP.* MeDIP was performed as previously described (Gal-Yam, Egger et al. 2008) with the following alterations: 10  $\mu$ g of sonicated genomic DNA (300–1000 bp in length) was denatured, incubated O/N at 4°C with 10  $\mu$ g of anti-methyl cytosine antibody (Diagenode, Belgium), and subsequently with 40  $\mu$ l of Dynabeads (M-280 Sheep anti-Mouse IgG –  $6.7 \times 10^8$  beads/ml) (Invitrogen) for 2 h at 4°C. The beads were washed and incubated with digestion buffer and proteinase K for 3 h at 50°C and the DNA was extracted by phenol chloroform and EtOH precipitation. For array experiments the output from 3 MeDIP reactions was combined (total of 30  $\mu$ g starting DNA) to constitute one replicate. The sonicated DNA served as input. MeDIP arrays were performed in biological triplicates.

*Array Design.* We collected a set of human ESC bivalent domains (Bernstein, Mikkelsen et al. 2006), and combined them with additional methylation related domains and control regions. We mapped these regions from the human genome to the rhesus genome using the UCSC liftOver program and tiled them with probes at 100 bp resolution (Roche-Nimblegen).

*Sample preparation and array hybridization.* The MeDIP DNA was amplified with a Sigma GenomePlex Complete Whole Genome Amplification (WGA) Kit using protocol developed in the Farnham lab (O'Geen et al. 2006). The amplified samples were column cleaned with the GenElute PCR clean-Up Kit. DNA quality and quantity were assessed with a Nanodrop device, and the size distribution estimated on 1% Agarose gel. The IP samples were labelled with Cy5 dye-labelled 9mers(blue), the reference samples with Cy3 dye-labelled 9mers(pink) from Trilink biotech, and after EtOH washing, drying and rehydrating, each sample was requantified with

NanoDrop. Adhering to the Nimblegen protocol, we pooled 6 ug of each sample and 6 ug of appropriate reference into the same tube prior to hybridization. The hybridizations were conducted with the Nimblegen Hybridization kit and X1 mixers, and placed on a 4 bay station for 18 hours. The slides were then washed and spot intensity in the two channels was recorded with Pix 4000B Scanner. The data were synchronized with the NimbleScan software, and exported for analyses. MeDIP data was normalized as described (Gal-Yam, Egger et al. 2008)(**Fig S6**), with the exception of subtraction of M.Sssl data, which was omitted since the main application of the data was the analysis of DMRs.

*Detection of DMRs.* To detect DMRs, we screened genomic windows of size 500bp to 20kb. We computed for each cell type (ORMES-22, Fibroblast, CRES-2, ORMES-9) the distribution of methylation values for all probes in the window (using triplicates as independent observations). We then tested the difference between any two distributions using Kolmogorov-Smirnov statistics and scored the window using the lowest p-value thus derived. Given the p-values for all windows, we selected the lowest p-value windows while excluding window overlap. This resulted in a set of non-overlapping DMRs with locally optimal p-value. We used a p-value threshold of  $10^{-5}$  to generate the set of DMRs analyzed, and a set of non overlapping 2kb regions with p-value larger than  $10^{-3}$  and median methylation higher than 0.3 or lower than 0 for background regions with high and low methylation respectively. We note that the p-values we computed are not corrected for the correlation between adjacent probes, but that such correction would affect all genomic loci uniformly, effectively only changing the thresholds we used (**Fig 3 and Fig 4**).

*Comparison to human data.* We renormalized the MeDIP data of Straussman et al. as described above, and computed the mean MeDIP signal for regions in the human genome that were orthologous to DMRs or background monkey regions. We disregarded regions that were not covered on the human array (which focused only on CpG islands). We used the data from Irizarry et al. as provided by the authors, computing the minus average CHARM levels for genomic regions that were orthologous to the monkey DMR and background regions.

*Genomic and epigenomic analysis.* To define rhesus TSS we used human genes mapped onto the monkey genome (UCSC). Gene expression data (GSE7748) was mapped onto these genes based on an overlap between the gene expression probe and the mapped gene body. Induced and repressed genes were defined as having at

least one mapped gene expression probe with  $\log_2(\text{fibroblast/ESC}) > 1.5$ . We considered DMRs as associated with a regulated TSS (**Fig 4**) if they overlapped the region 1kb around the TSS. Comparison of rhesus methylation and human ESC Suz12 or human fibroblast CTCF was carried out by mapping monkey data onto the human genome. Suz12 ChIP (Lee, Jenner et al. 2006) and CTCF ChIP (Kim, Abdullaev et al. 2007) data were renormalized as previously described. To generate the counts in **Fig 4** we associated a DMR with CTCF or Suz12 if at least one of the probes in the DMR had a normalized CTCF or Suz12 ChIP value larger than 1.5.

## REFERENCES:

- Allegrucci, C., Y. Z. Wu, et al. (2007). "Restriction landmark genome scanning identifies culture-induced DNA methylation instability in the human embryonic stem cell epigenome." Hum Mol Genet **16**(10): 1253-68.
- Bernstein, B. E., T. S. Mikkelsen, et al. (2006). "A bivalent chromatin structure marks key developmental genes in embryonic stem cells." Cell **125**(2): 315-26.
- Brunner, A. L., D. S. Johnson, et al. (2009). "Distinct DNA methylation patterns characterize differentiated human embryonic stem cells and developing human fetal liver." Genome Res.
- Byrne, J. A., D. A. Pedersen, et al. (2007). "Producing primate embryonic stem cells by somatic cell nuclear transfer." Nature **450**(7169): 497-502.
- Cokus, S. J., S. Feng, et al. (2008). "Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA methylation patterning." Nature **452**(7184): 215-9.
- Deng, J., R. Shoemaker, et al. (2009). "Targeted bisulfite sequencing reveals changes in DNA methylation associated with nuclear reprogramming." Nat Biotechnol **27**(4): 353-60.
- Dighe, V., L. Clepper, et al. (2008). "Heterozygous embryonic stem cell lines derived from nonhuman primate parthenotes." Stem Cells **26**(3): 756-66.
- Farthing, C. R., G. Ficuz, et al. (2008). "Global mapping of DNA methylation in mouse promoters reveals epigenetic reprogramming of pluripotency genes." PLoS Genet **4**(6): e1000116.
- Fouse, S. D., Y. Shen, et al. (2008). "Promoter CpG Methylation Contributes to ES Cell Gene Regulation in Parallel with Oct4/Nanog, PcG Complex, and Histone H3 K4/K27 Trimethylation." Cell Stem Cell **2**(2): 160-169.
- Gal-Yam, E. N., G. Egger, et al. (2008). "Frequent switching of Polycomb repressive marks and DNA hypermethylation in the PC3 prostate cancer cell line." Proc Natl Acad Sci U S A **105**(35): 12979-84.
- Irizarry, R. A., C. Ladd-Acosta, et al. (2008). "Comprehensive high-throughput arrays for relative methylation (CHARM)." Genome Res **18**(5): 780-90.
- Irizarry, R. A., C. Ladd-Acosta, et al. (2009). "The human colon cancer methylome shows similar hypo- and hypermethylation at conserved tissue-specific CpG island shores." Nat Genet **41**(2): 178-86.
- Keshet, I., Y. Schlesinger, et al. (2006). "Evidence for an instructive mechanism of de novo methylation in cancer cells." Nat Genet **38**(2): 149-53.
- Kim, T. H., Z. K. Abdullaev, et al. (2007). "Analysis of the vertebrate insulator protein CTCF-binding sites in the human genome." Cell **128**(6): 1231-45.
- Lee, T. I., R. G. Jenner, et al. (2006). "Control of developmental regulators by Polycomb in human embryonic stem cells." Cell **125**(2): 301-13.
- Meissner, A., T. S. Mikkelsen, et al. (2008). "Genome-scale DNA methylation maps of pluripotent and differentiated cells." Nature **454**(7205): 766-70.
- Mikkelsen, T. S., J. Hanna, et al. (2008). "Dissecting direct reprogramming through integrative genomic analysis." Nature **454**(7200): 49-55.
- Mikkelsen, T. S., M. Ku, et al. (2007). "Genome-wide maps of chromatin state in pluripotent and lineage-committed cells." Nature **448**(7153): 553-60.
- Mitalipov, S., H. C. Kuo, et al. (2006). "Isolation and characterization of novel rhesus monkey embryonic stem cell lines." Stem Cells **24**(10): 2177-86.

- Mohn, F., M. Weber, et al. (2008). "Lineage-specific polycomb targets and de novo DNA methylation define restriction and potential of neuronal progenitors." Mol Cell **30**(6): 755-66.
- Mohn, F., M. Weber, et al. (2009). "Methylated DNA immunoprecipitation (MeDIP)." Methods Mol Biol **507**: 55-64.
- Ohm, J. E., K. M. McGarvey, et al. (2007). "A stem cell-like chromatin pattern may predispose tumor suppressor genes to DNA hypermethylation and heritable silencing." Nat Genet **39**(2): 237-42.
- Rakyan, V. K., T. A. Down, et al. (2008). "An integrated resource for genome-wide identification and analysis of human tissue-specific differentially methylated regions (tDMRs)." Genome Res **18**(9): 1518-29.
- Reik, W. (2007). "Stability and flexibility of epigenetic gene regulation in mammalian development." Nature **447**(7143): 425-32.
- Rollins, R. A., F. Haghghi, et al. (2006). "Large-scale structure of genomic methylation patterns." Genome Res **16**(2): 157-63.
- Schlesinger, Y., R. Straussman, et al. (2007). "Polycomb-mediated methylation on Lys27 of histone H3 pre-marks genes for de novo methylation in cancer." Nat Genet **39**(2): 232-6.
- Straussman, R., D. Nejman, et al. (2009). "Developmental programming of CpG island methylation profiles in the human genome." Nat Struct Mol Biol **16**(5): 564-71.
- Takahashi, K. and S. Yamanaka (2006). "Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors." Cell **126**(4): 663-76.
- Tanay, A., A. H. O'Donnell, et al. (2007). "Hyperconserved CpG domains underlie Polycomb-binding sites." Proc Natl Acad Sci U S A **104**(13): 5521-6.
- Weber, M., J. J. Davies, et al. (2005). "Chromosome-wide and promoter-specific analyses identify sites of differential DNA methylation in normal and transformed human cells." Nat Genet **37**(8): 853-62.
- Weber, M., I. Hellmann, et al. (2007). "Distribution, silencing potential and evolutionary impact of promoter DNA methylation in the human genome." Nat Genet **39**(4): 457-66.
- Widschwendter, M., H. Fiegler, et al. (2007). "Epigenetic stem cell signature in cancer." Nat Genet **39**(2): 157-8.

**Fig 1: Profiling Monkey ES methylation. A) Experimental design.** We studied DNA methylation in native ES cells (ORMES-22), fibroblasts and ES cells generated by somatic cell nuclear transfer (CRES-2). We also assayed a distinctly different native ES line (the homozygous parthenote ORMES-9) to control for ES line specific effects. DNA methylation profiles in these four cell types were assayed using MeDIP and tiling arrays. Values of zero correspond to average genomic methylation. **B) Near perfect DNA methylation reprogramming in CRES-2.** Shown are the differential methylation values for 380,000 array probes covering orthologous human K4-K27 bivalent domains and selected DNA methylation hotspots. The differences between Fibroblasts and the two stem cells lines are highly correlated, showing that at the global level, reprogramming of the fibroblast epigenome during nuclear transfer is near perfect. **C) Conserved and differential methylation in HOX clusters.** Shown are the methylation profiles at the HOXA (upper) and HOXD clusters (lower), which were tiled completely on our array and reflect an excellent overall correlation between the native and reprogrammed ES cells. Regions undergoing fibroblasts hypermethylation (hyper-DMRs - red) or hypomethylation (hypo-DMR - green) are highlighted. In contrast to the good overall correspondence between native and reprogrammed ES cells methylation, a small region in the HOXA cluster (marked blue), shows a CRES-2 methylation pattern that is similar to the fibroblast profile, suggesting incomplete reprogramming or independent hypermethylation in OMRES-22 and ORMES-9.

**Fig 2: Differentially methylated regions (DMRs). A) Global patterns of methylation reprogramming.** DMRs were statistically extracted from the data by comparing methylation in all pairs of cell types, thereby not pre-assuming any type of organization. Median methylation values for each DMR over all cell types were then clustered (k-means). Shown are the color coded methylation values of each DMR (rows), organized into clusters showing higher methylation in fibroblasts than in native stem cells (hyper-DMRs) and clusters showing lower methylation in fibroblasts than in native stem cells (hypo-DMRs). Overall, the clusters reflect different basal levels of methylation across the genome, but good correspondence between methylation in the different ES lines. An important exception is a cluster including DMRs with significantly higher methylation in ORMES-22 than in CRES-2. Some of the DMRs in this cluster may reflect ORMES-22 line-specific effects and were excluded from further analysis. Other DMRs in this cluster are also hypermethylated in the ORMES-9 line, were classified as "failed reprogramming" DMRs and were analyzed separately. **B) Distribution of DMR sizes.** Shown is the distribution of



sizes of genomic intervals determined to be hyper- (red) and hypo- (green) DMRs. Hypo-DMRs have a more specific length distribution, peaking around 2KB. **C) DMR CpG content.** The average number of CpGs in 500bp windows was computed for each DMR (each CpG was counted twice) and the distribution of CpG contents for hyper- and hypo- DMRs was plotted. Hyper-DMRs have a lower overall CpG content. Importantly, both types of DMRs generally occupy regions of low to medium CpG content, and are not observed in classical CpG islands (CpG content > 50).

**Fig 3: Monkey DMRs are conserved in human ESC and tissues. A) Muscle-ES differential methylation.** Shown are boxplots of the DNA methylation differences between human muscle tissues and human ESCs (Straussman, Nejman et al. 2009), computed for regions of the human genome that are orthologous to monkey hyper- and hypo- DMRs (red and green), or to regions with low or high monkey ES methylation (blue and yellow). Since the human data span only CpG islands, the statistics only cover regions with intermediate or high CpG content. P values indicate the significance (using KS test) of difference between hyper-DMRs and low ES methylation regions, and between hypo-DMRs and high ES methylation regions. **B) Range of methylation across a panel of human tissues.** Shown are boxplots for the differences between the minimum and maximum DNA methylation in human brain, colon, spleen and liver (Irizarry, Ladd-Acosta et al. 2009), for regions that are orthologous to monkey DMRs or regions of high and low monkey ES methylation (same color scheme as in (A)). P values indicate the significance of difference between hypo-DMRs and high ES methylation (for CpG range 0-15) and between hyper-DMRs and low ES methylation (for CpG range 15-40).

**Fig 4: CTCF and Polycomb are correlated with differential methylation. A) Suz12 occupancy.** Shown are box plots for average human Suz12 occupancy on mapped monkey DMRs and background regions. We separately plot groups of regions with different levels of CpG content, dissected into hyper-DMRs (red), hypo-DMRs (green), regions with low ES methylation (blue) and regions with high ES methylation (yellow). In general, regions with high methylation have low Suz12 levels (e.g. lower than regions with low methylation, see CpG content 15-40). Moreover, regions with higher CpG content (>40) that are hyper and hypo methylated are enriched in Suz12 targets. **B) CTCF occupancy at CTCF motifs.** Shown are distributions of CTCF binding levels in three groups of genomic loci: 1) background regions lacking CTCF motifs and having low DNA methylation (gray), 2) regions of high methylation featuring CTCF binding motifs (red) 3) regions with low DNA

methylation lacking CTCF binding motifs (green). **C) CTCF binding capacity at DMRs.** Shown are cumulative probability distributions for the predicted binding energy of the CTCF motif in hyper-DMRs (red) and hypo-DMRs (green). About 15% of the hypo-DMRs have a strong CTCF binding site, much higher than the percentage for hyper-DMRs. **D) Combinatorial analysis.** Shown are counts of DMRs associated with combinations of regulated TSS, Suz12 hotspot or CTCF binding site. More than half of the DMRs have at least one factor associated with them.

**Fig. 5: Partial and failed reprogramming. A) Failure to reprogram DMRs.** Shown are examples of DMRs in which the reprogrammed ES DNA methylation pattern follows the fibroblast pattern. These stand in marked contrast to the overall genomic trend (e.g., **Fig. 1**) and may represent complete lack of reprogramming, partial reprogramming that could not complete, or ongoing reprogramming with much slower kinetics than the genomic trend. **B) Reprogramming ratios.** Reprogramming ratios were computed as the ratio of the difference between the reprogrammed ES and fibroblast methylation medians and the difference between the native ES and fibroblast methylation medians. Ratio of 1 indicates perfect reprogramming, and a ratio of 0 represents no reprogramming. Plotted is the distribution of reprogramming ratios of hypo-DMRs and hyper-DMRs. Data is only shown for DMRs that had similar methylation levels in the two native ES lines (ORMES-22 and ORMES-9).

**Table S1: DMR properties.**

**Fig S1: Correlation of biological replicates.** Scatter plots depicting the correlation between normalized methylation values of two biological replicates for each cell line.

**Fig S2: Differential methylation distributions. A)** The distribution of the difference in methylation between ORMES-22 and fibroblasts in hyper-DMRs **B)** The distribution of the difference in methylation between ORMES-22 and fibroblasts in hypo-DMRs.

**Fig S3: Sequence analysis of hypo- and hyper- DMRs at low CpG levels. A) dinucleotide distribution.** The frequency of dinucleotides was computed for hypo-DMRs (green), hyper-DMRs (red) and background regions (grey) with average CpG content less than 15. Background regions were chosen as regions in which the difference between ORMES-22, fibroblasts and CRES-2 was insignificant (KS p-value > 0.001). **B) CpG sequence context.** Shown are the frequency profiles for dinucleotides following a CpG in hypo-DMRs (green), hyper-DMRs (red) and background regions (grey) in which the average CpG content is less than 15.

**Fig S4: Sequence analysis of hypo- and hyper- DMRs at intermediate CpG levels. A) dinucleotide distribution. B) CpG sequence context.** As Fig S3, for average CpG content between 15 and 40.

**Fig S5: DMRs at regulated promoters.** Induced and repressed genes (here at least two fold) were extracted using gene expression array profiles on fibroblasts, ORMES-22 and CRES-2. Individual tiling array probes were classified as hyper-methylated and hypo-methylated according to their fibroblast and ORMES-22 array values (normalized array difference > 0.6), and were further partitioned according to their CpG content (low: CpG < 15, intermediate: 15 < CpG < 50, high CpG > 50). We computed the number of probes at each distance relative to induced (upper) and repressed (lower) TSSs. The  $\log_2$  of the ratio between this number and the number expected by chance (assuming hyper- and hypo- methylated probes are randomly distributed) is plotted. As expected from the general anti-correlation between gene expression and DNA methylation, induced genes are enriched with hypo-methylated probes and anti-enriched for hyper-methylated probes, while repressed genes are enriched with hyper-methylated probes.

**Fig S6: Occupancy of pluripotency factors. A) Oct4 occupancy.** Shown are box plots for average human Oct4 occupancy on mapped monkey DMRs and background regions. We plot separately groups of regions with different levels of CpG content, dissected into groups of hyper-DMRs (red), hypo-DMRs (green), regions with low ES methylation (blue) and regions with high ES methylation (yellow). **B) Nanog occupancy.** As part A, for Nanog.

**Fig S7: Properties of low reprogramming ratio DMRs. A) Reprogrammed DMR size distribution.** Shown is the log size distribution of genomic intervals determined to be hypo-DMRs with reprogramming ratio less than 0.25 (limited, blue), and hypo-DMRs with reprogramming ratio greater than 0.5 (extensive, black). **B) Reprogrammed DMR CpG content.** The average number of CpGs in 500bp windows was computed for each DMR (each CpG was counted twice) and the distribution of CpG content for limited-reprogramming hypo-DMRs and extensive-reprogramming hypo-DMRs was plotted. **C) Reprogrammed DMRs Suz12 occupancy.** Shown is the distribution of Suz12 occupancy in human ESCs on mapped monkey hypo-DMRs with limited reprogramming (blue) and in hypo-DMRs with extensive reprogramming (black). **D) Reprogrammed DMR CTCF occupancy.** Shown is the distribution of CTCF occupancy in human fibroblasts on mapped monkey hypo-DMRs with limited reprogramming (blue) and in hypo-DMRs with extensive reprogramming (black). **E) Reprogrammed DMR telomere distance distribution.** The distance from telomere was computed for each DMR and the distributions of telomere distance of limited reprogrammed hypo-DMRs and extensively reprogrammed hypo-DMRs were plotted.

**Fig S8: CpG content vs. array binding ratio for all MeDIP samples.** Shown are the means of MeDIP binding ratios for every biological replicate and cell line, computed for bins of weighted probes' CpG contents, (hyb\_e curves, see Methods, Gal-Yam, Egger et al. 2008) and different G and C contents.

**Fig S9: qPCR validation.** Q-PCR amplification of genomic DNA fragments was purified by pull-down with anti-5 methylcytosine antibody. Paternally methylated *H19* and maternally methylated *SNRPN* promoters were used as control genes. Enrichment of *H19* (black column) and *SNRPN* (blue column) was observed in IVF-derived ES cells, SCNT-derived ES cells and donor fibroblasts, consistent with CpG methylation on both the paternal and maternal alleles in these XY cells. As expected, enrichment of *H19* was observed in sperm and enrichment of *SNRPN* was observed

in homozygous parthenogenetic ES cells. Enrichment values were plotted after normalization with 10% input DNA. Data represents the mean  $\pm$  S.E.M. (n = 4).

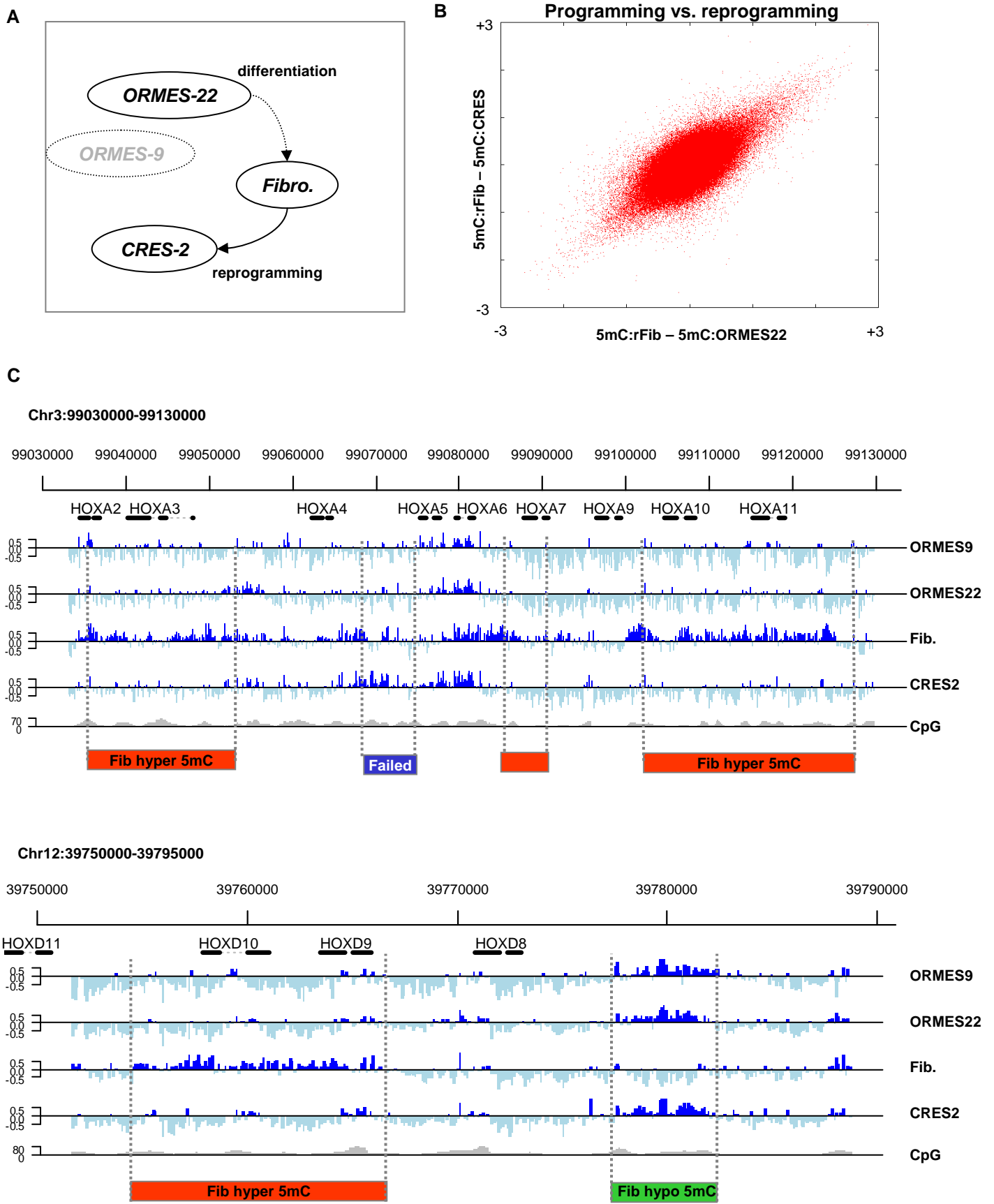
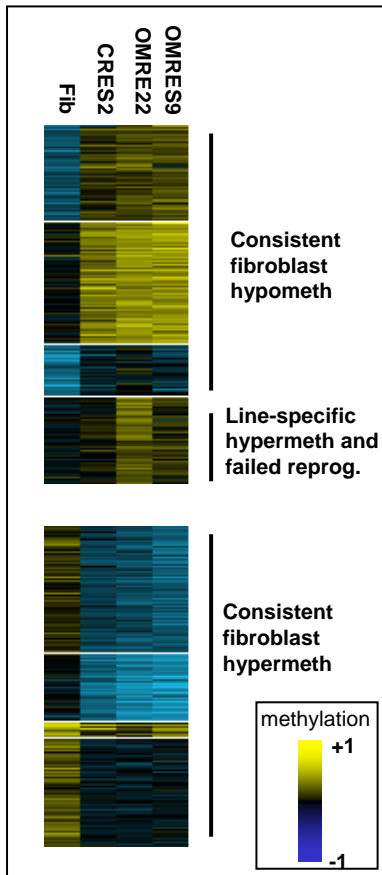
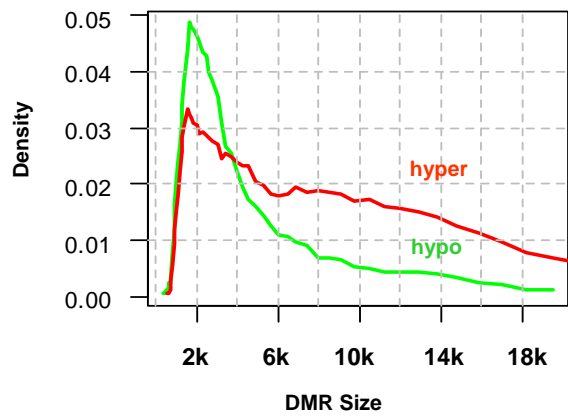


Figure 1

A



B



C

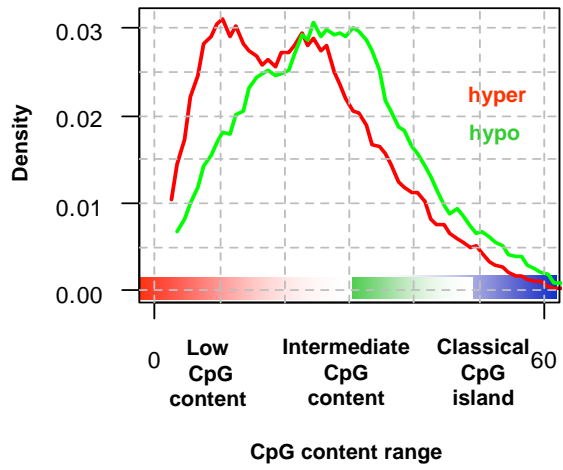
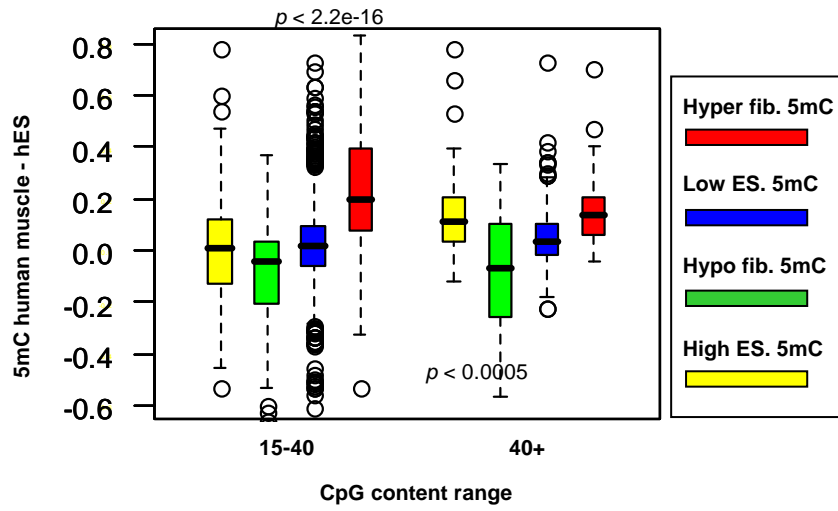


Figure 2

**A**



**B**

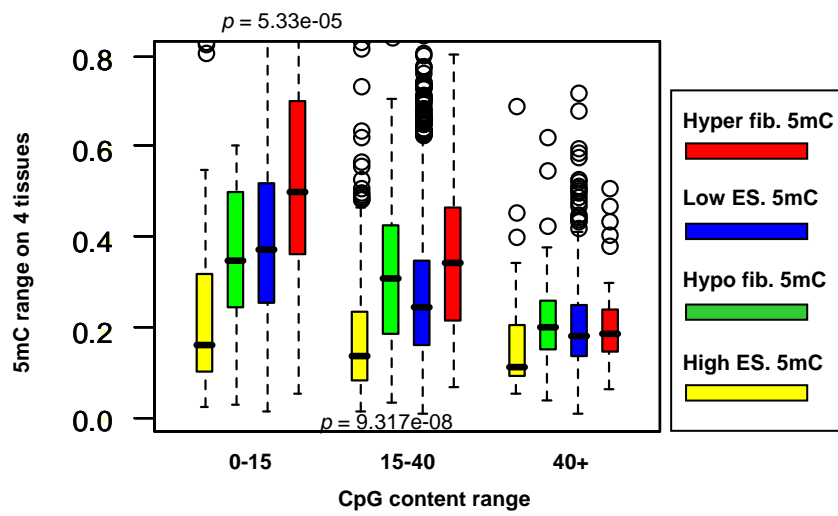
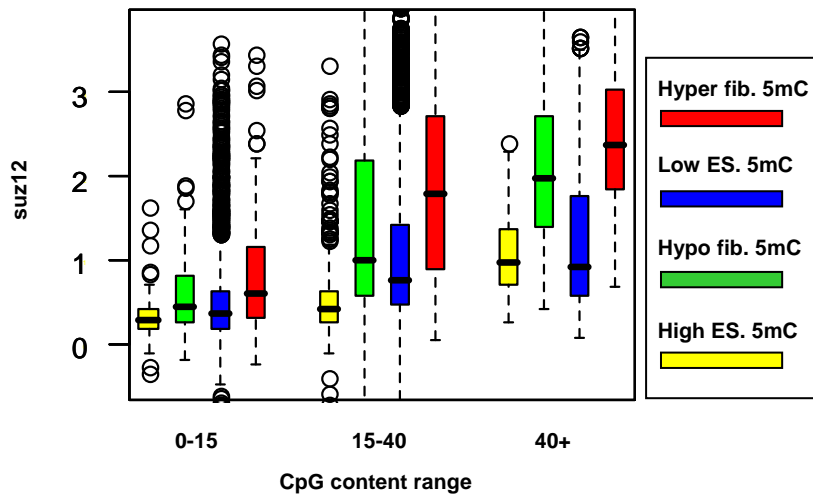


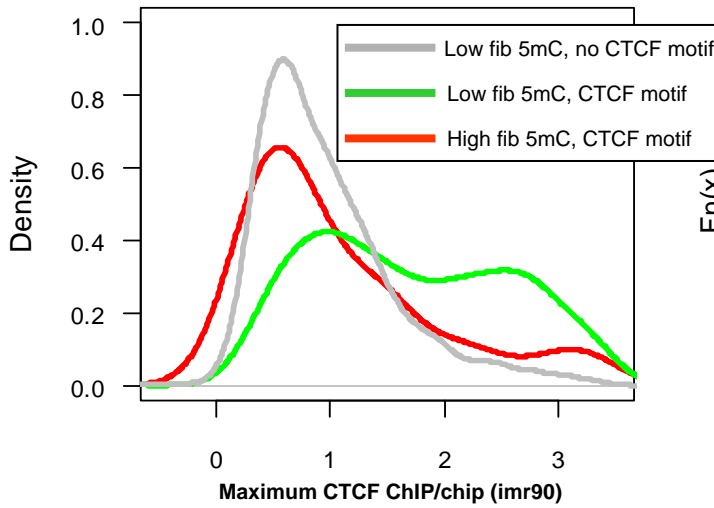
Figure 3



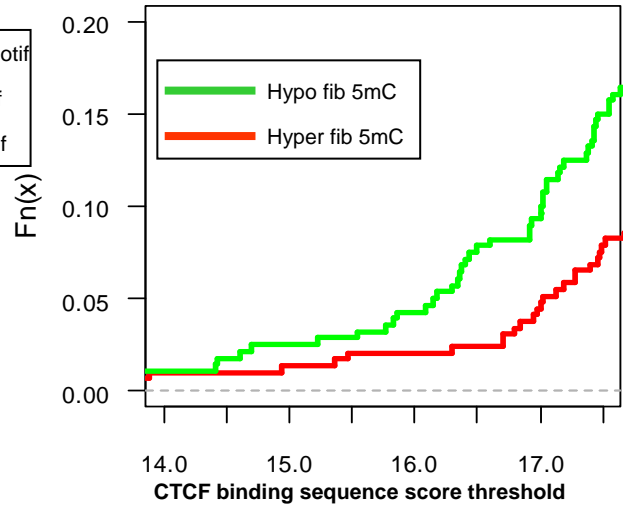
A



B



C

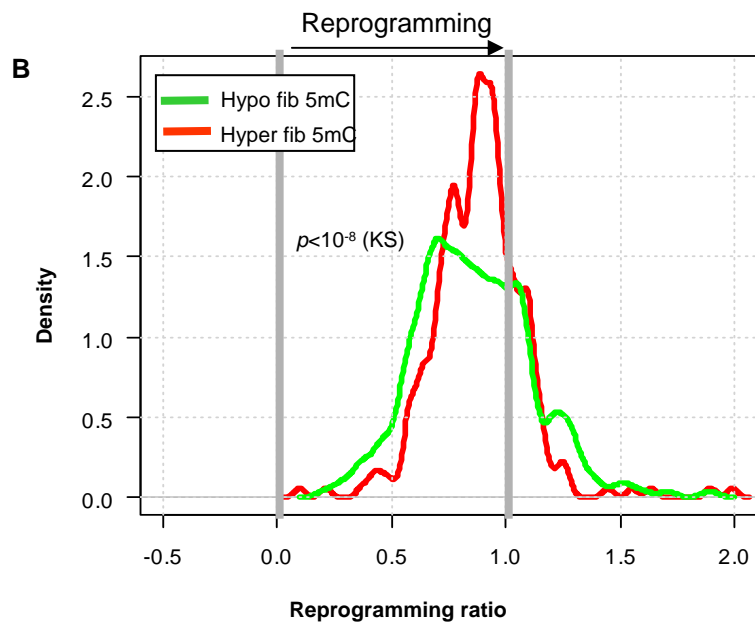
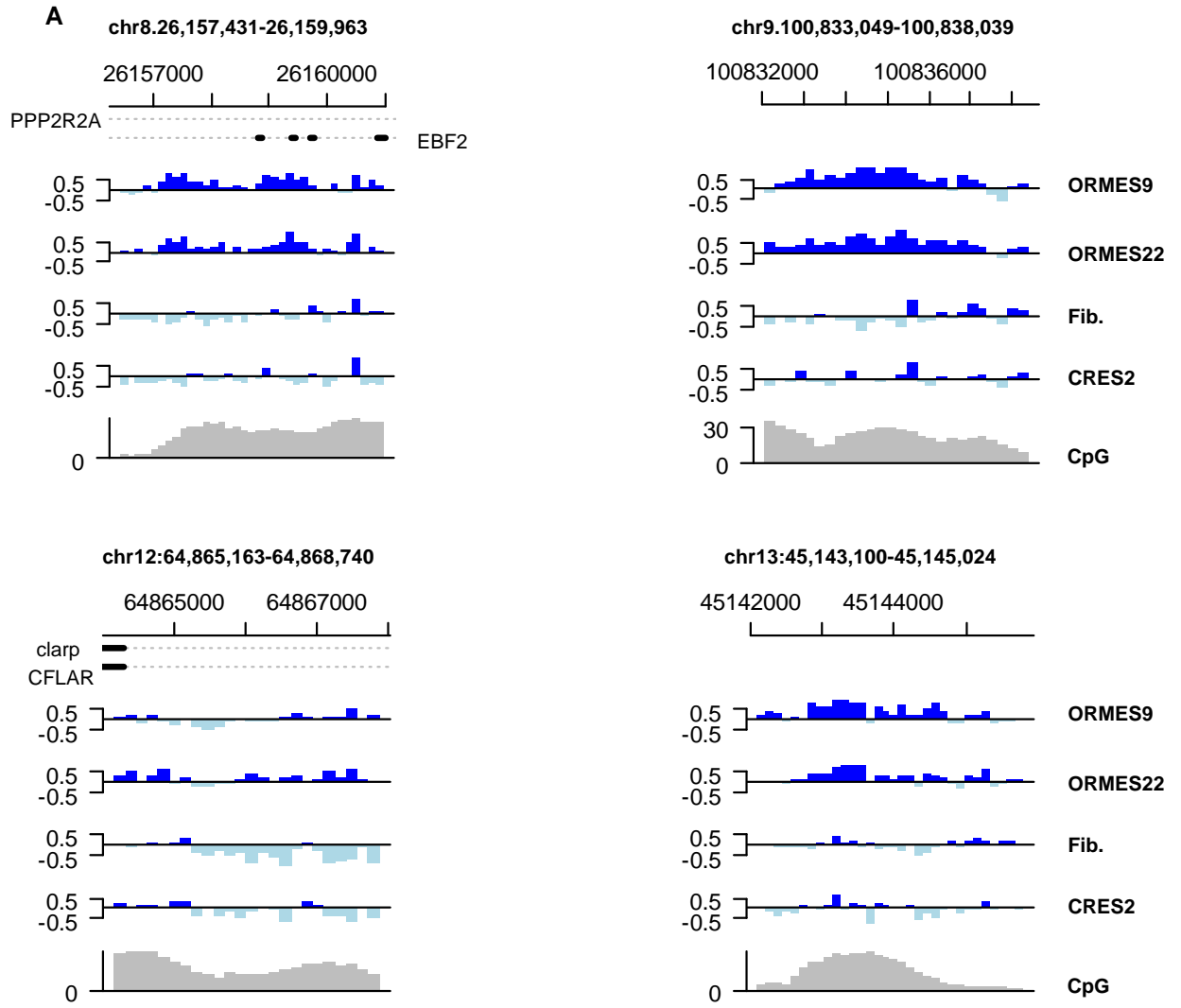


D

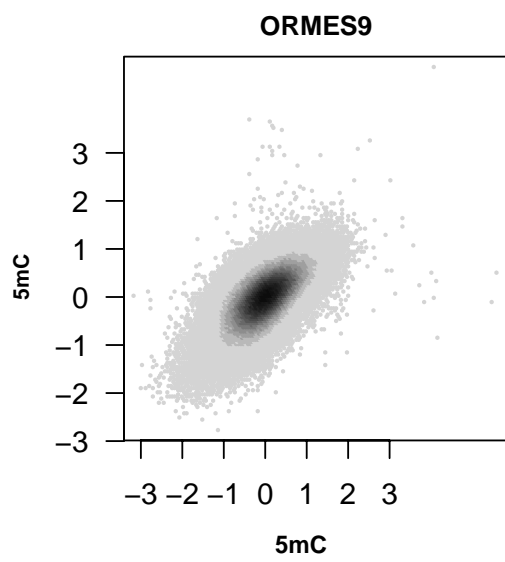
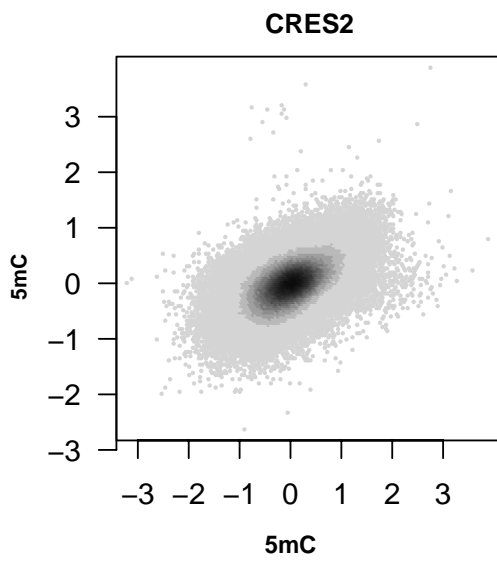
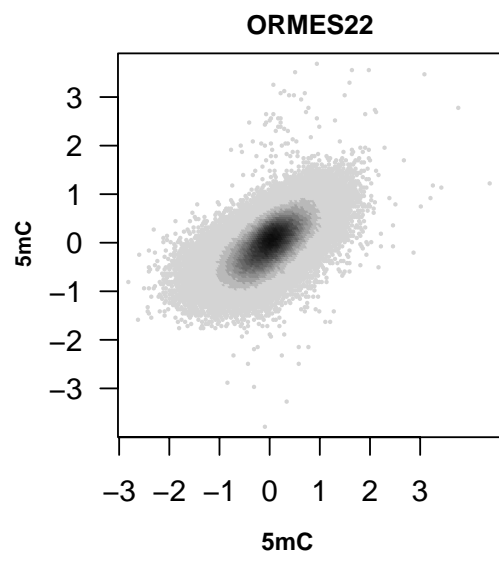
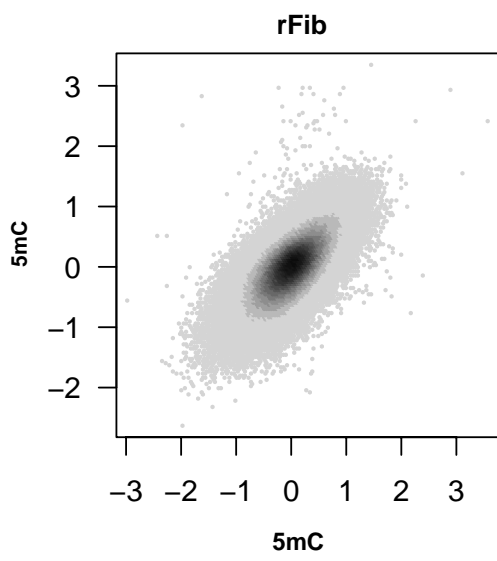
**Combinatorial analysis**

Regulated TSS	SUZ12	CTCF	hyper DMRs	hypo DMRs
1	0	0	23	32
1	1	0	11	11
1	0	1	5	11
1	1	1	8	4
0	1	0	93	90
0	1	1	26	32
0	0	1	13	59
0	0	0	152	152
<b>Total</b>			<b>331</b>	<b>391</b>

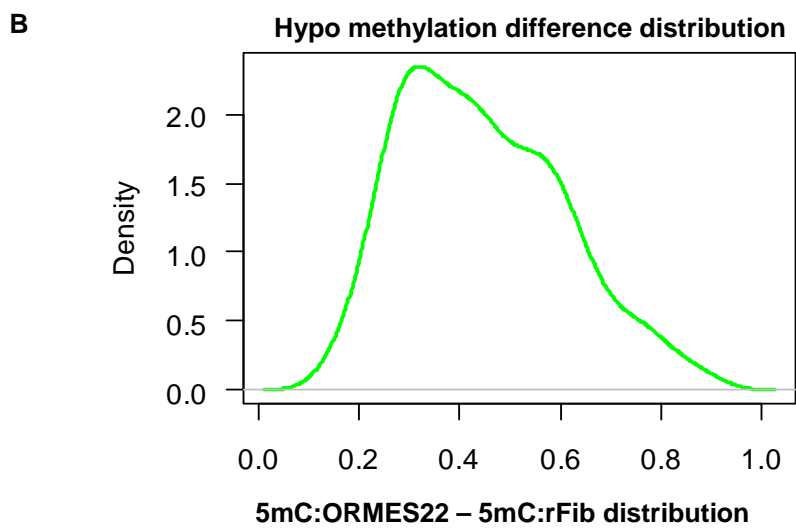
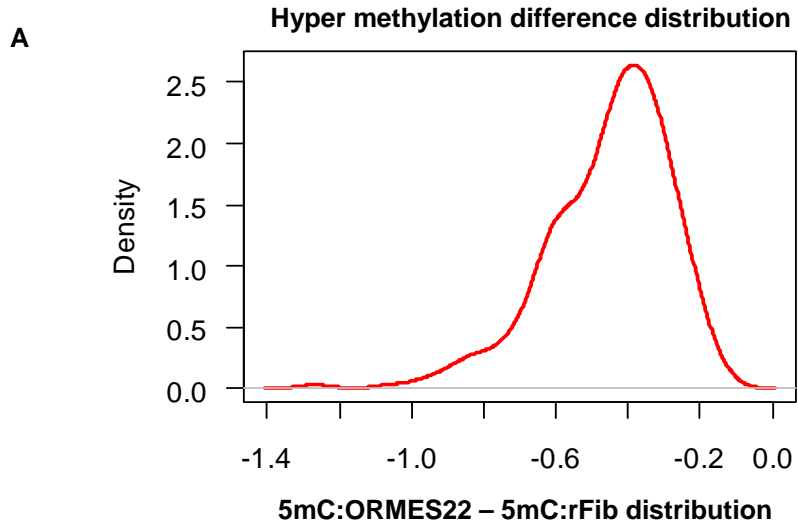
Figure 4



**Figure 5**



**Figure S1**



**Figure S2**

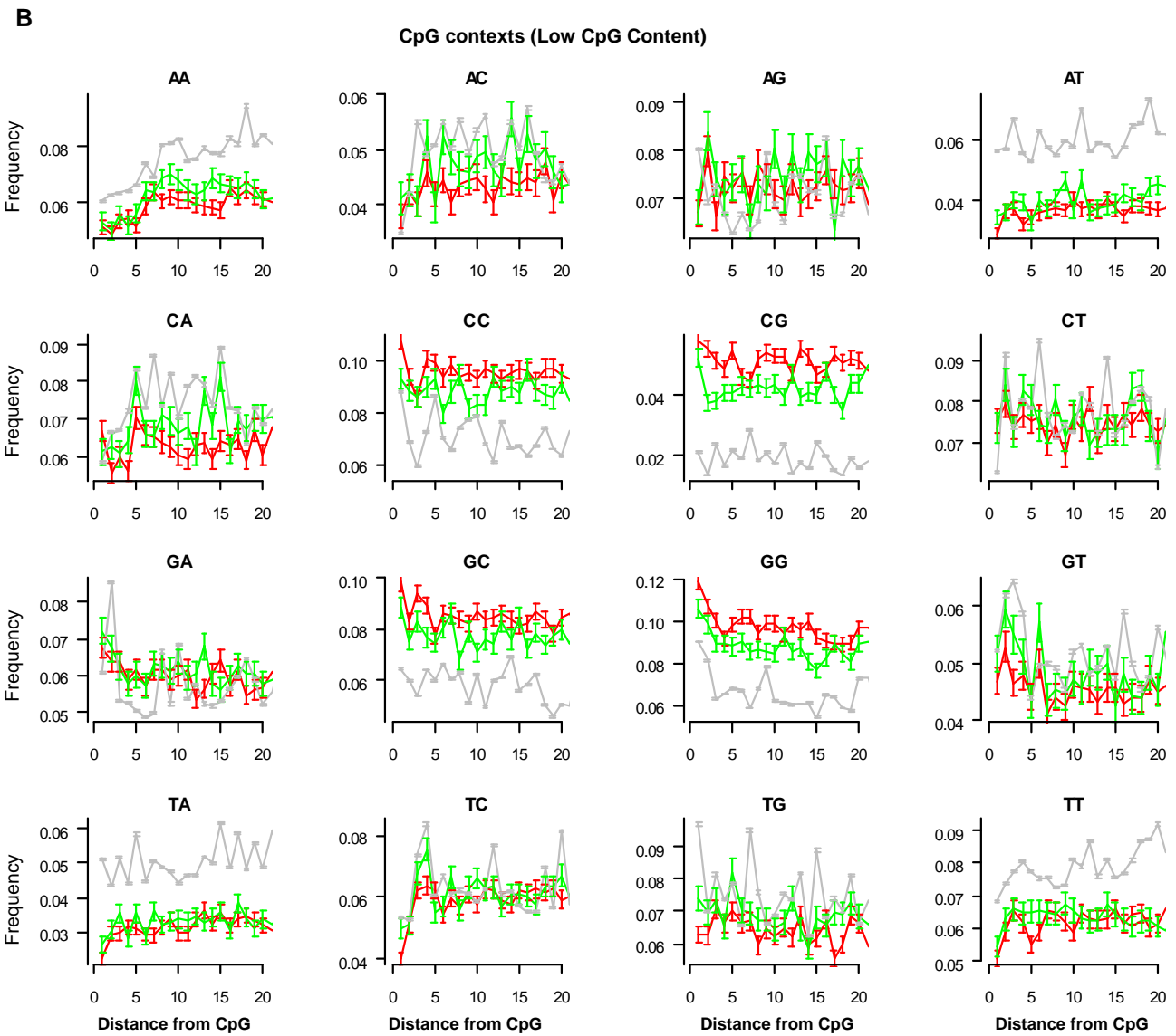
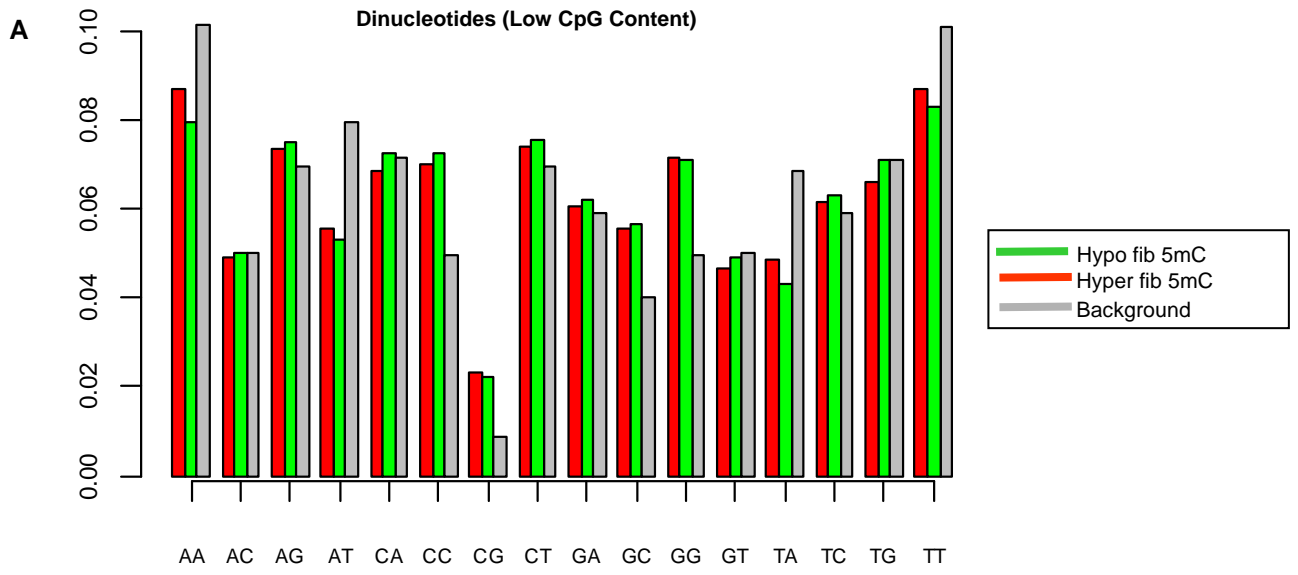


Figure S3

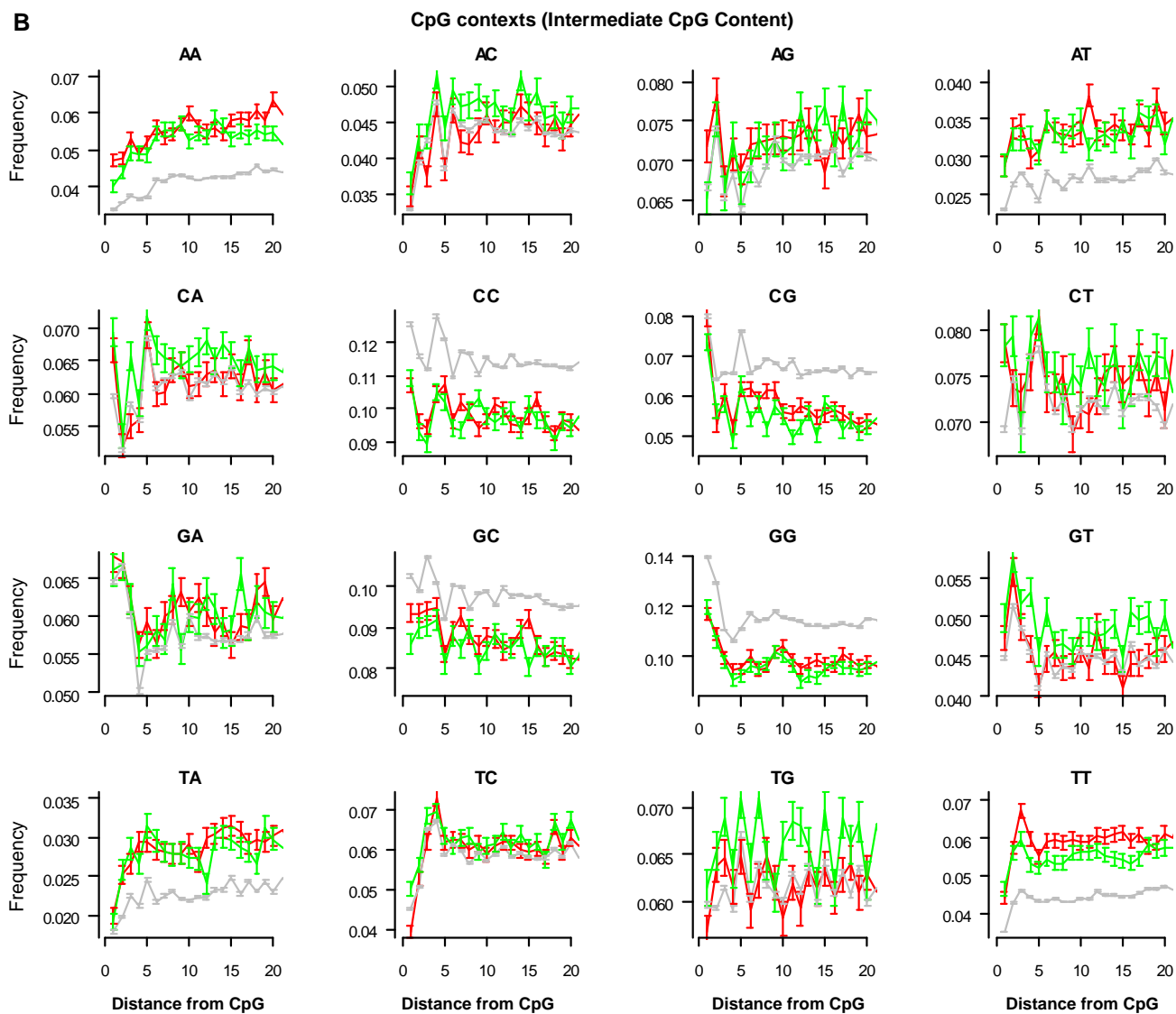
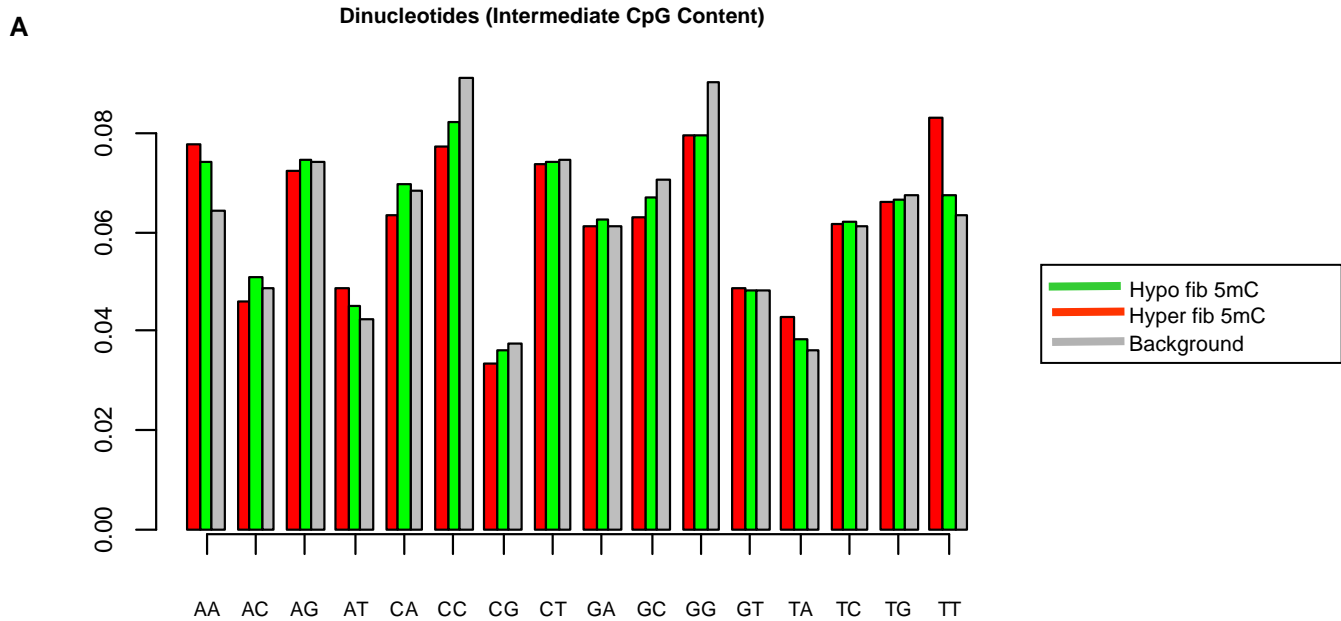


Figure S4

A

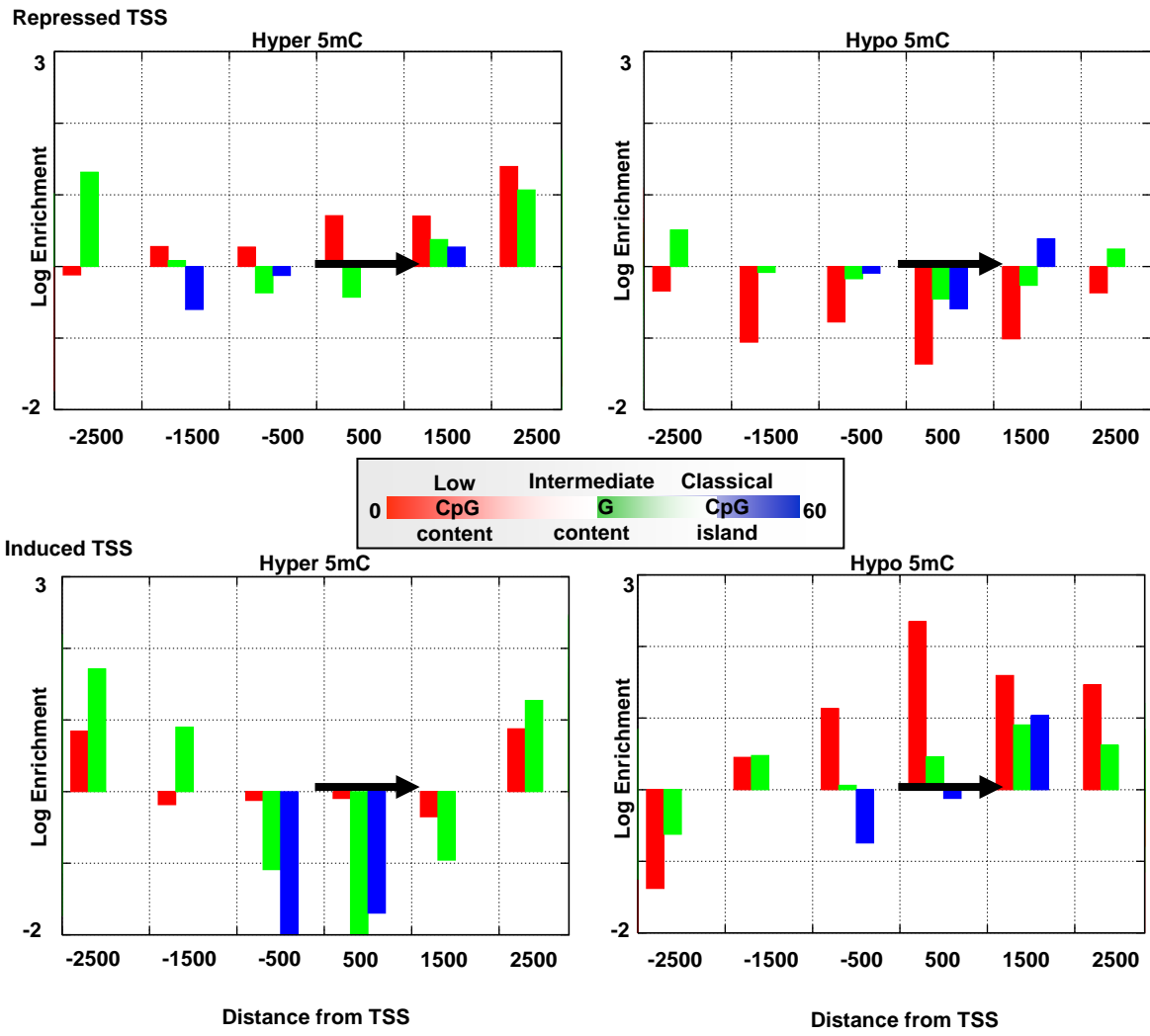
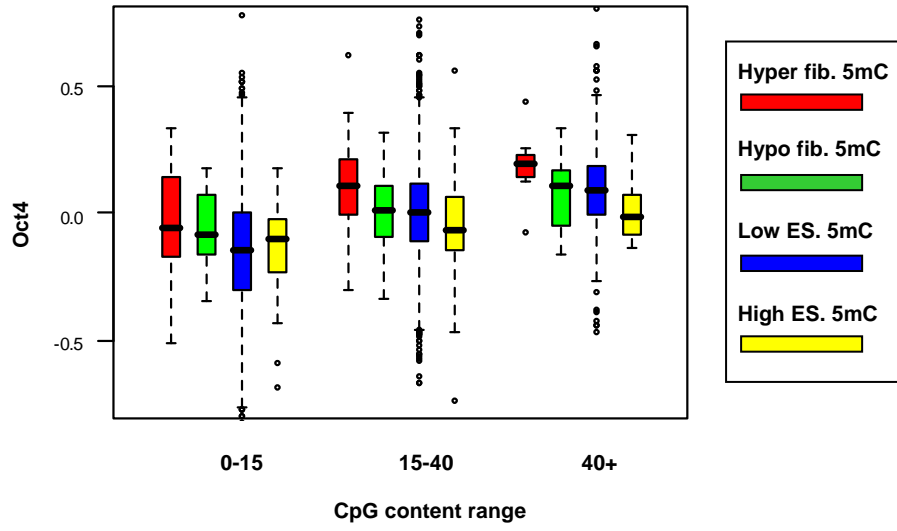


Figure S5

**A**



**B**

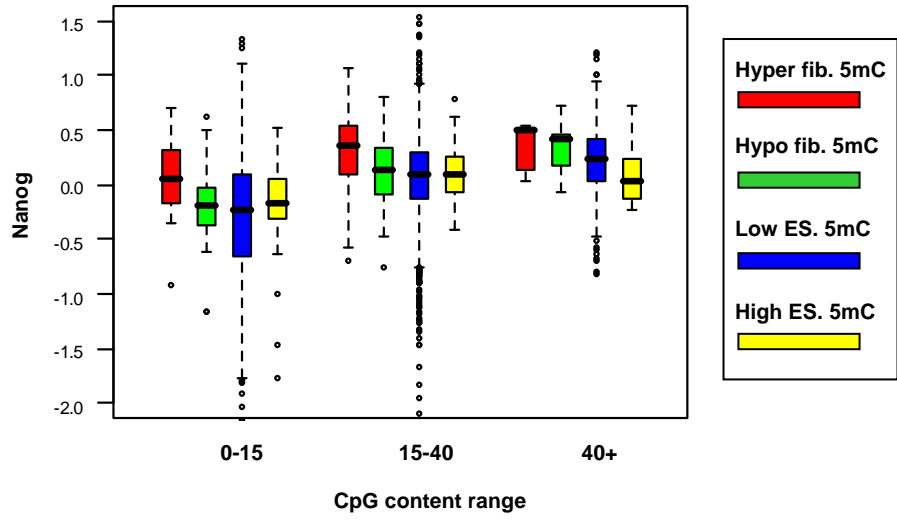


Figure S6



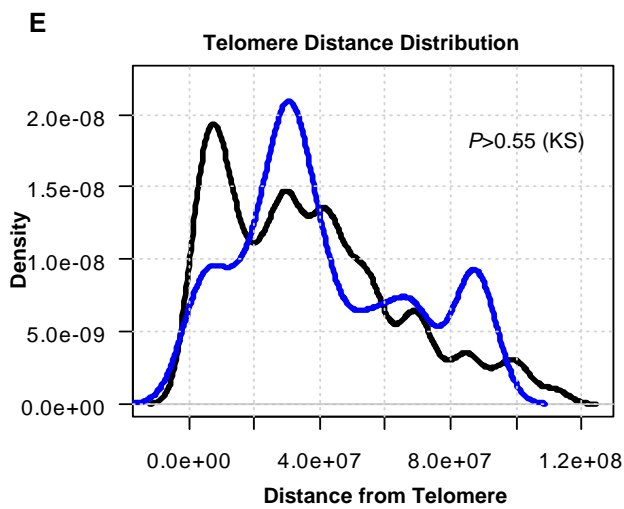
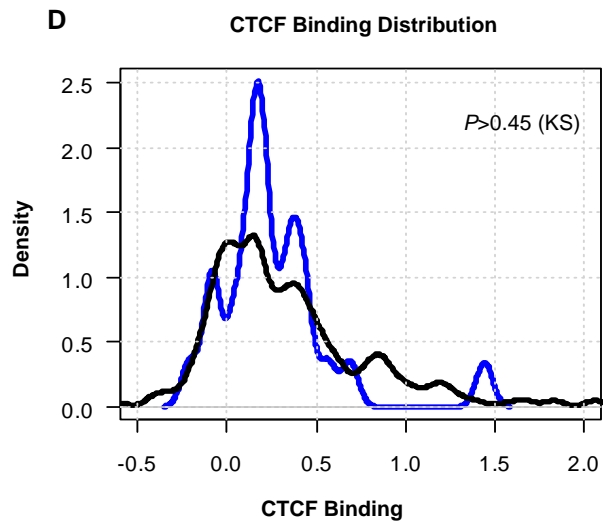
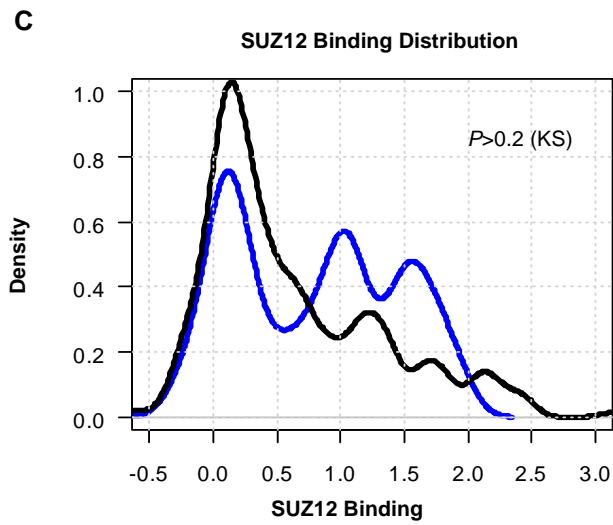
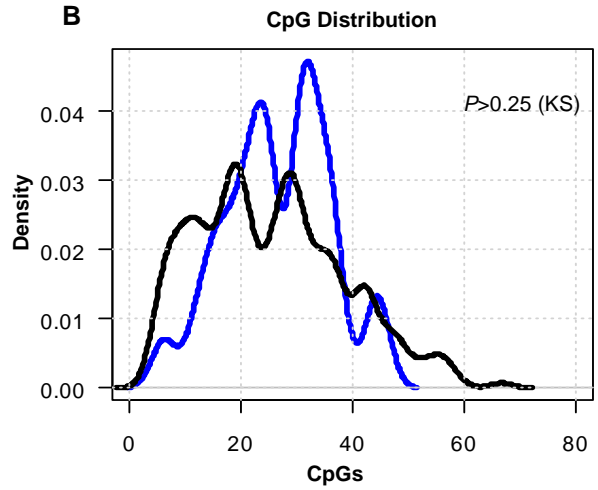
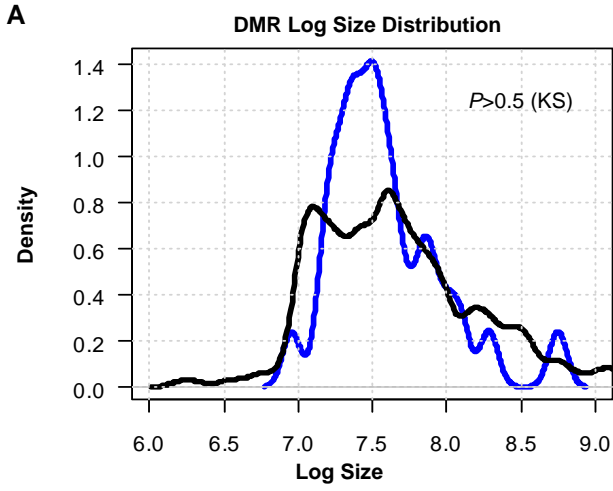
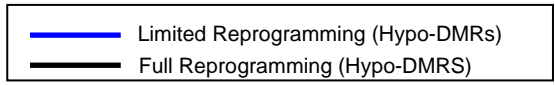


Figure S7

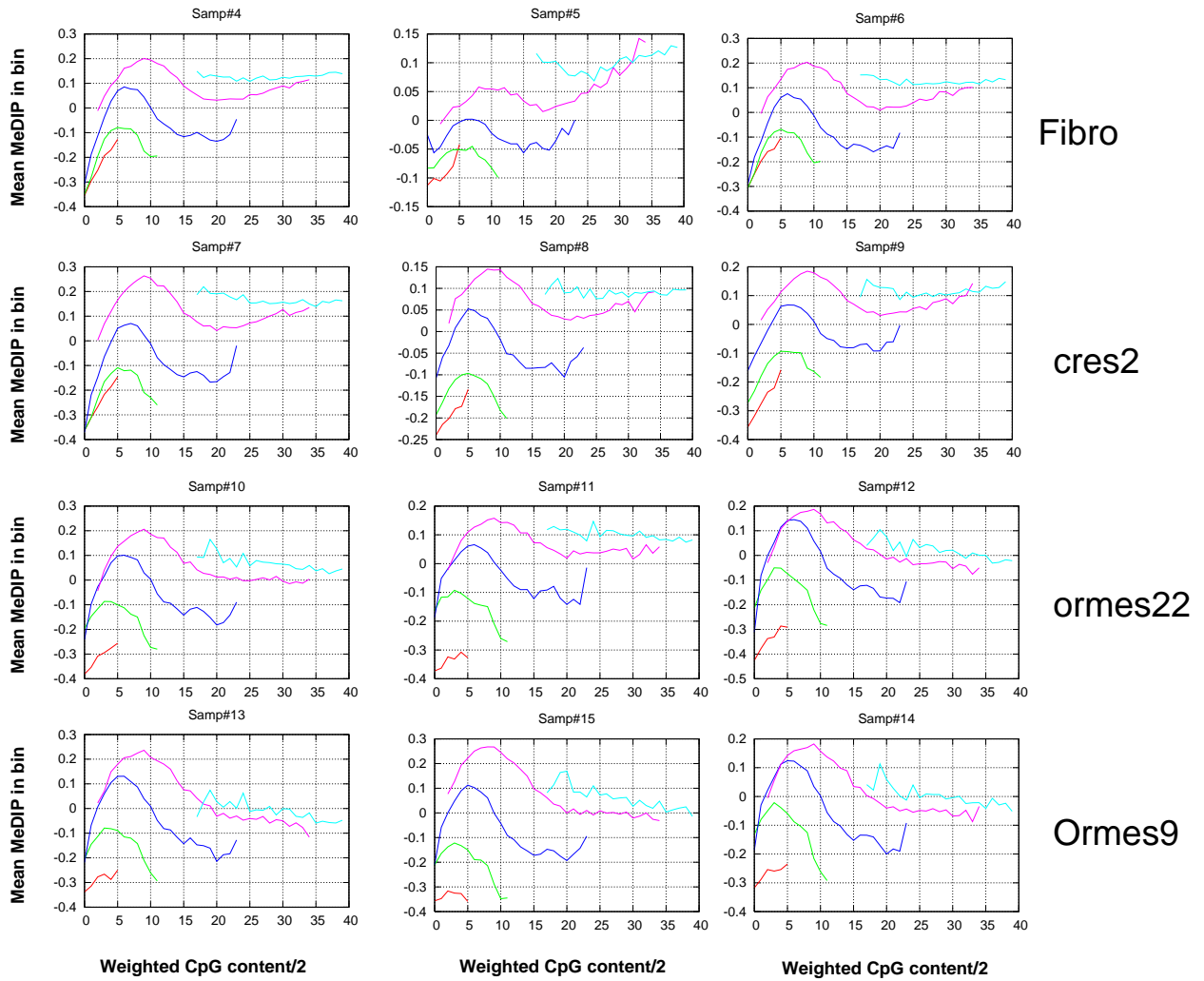


Figure S8

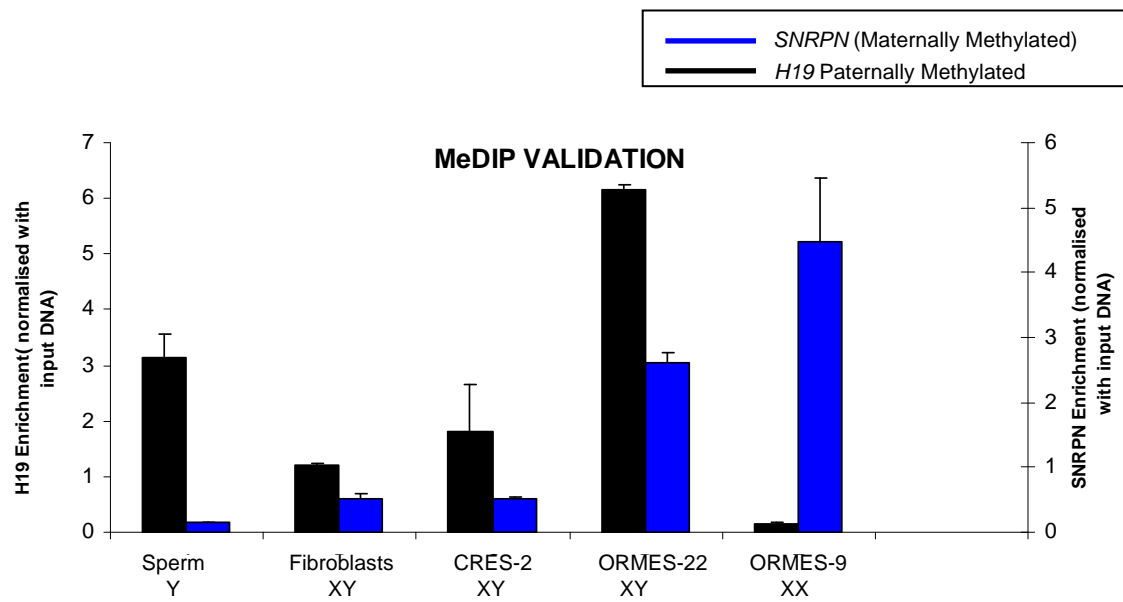


Figure S9