

Genome-wide distribution of 5-formylcytosine in embryonic stem cells is associated with transcription and depends on thymine DNA glycosylase

Eun-Ang Raiber^{1†}, Dario Beraldi^{1,2†}, Gabriella Ficzi³, Heather E Burgess^{3,4}, Miguel R Branco^{3,4}, Pierre Murat¹, David Oxley⁵, Michael J Booth¹, Wolf Reik^{3,4*†} and Shankar Balasubramanian^{1,2,6*}

Abstract

Background: Methylation of cytosine in DNA (5mC) is an important epigenetic mark that is involved in the regulation of genome function. During early embryonic development in mammals, the methylation landscape is dynamically reprogrammed in part through active demethylation. Recent advances have identified key players involved in active demethylation pathways, including oxidation of 5mC to 5-hydroxymethylcytosine (5hmC) and 5-formylcytosine (5fC) by the TET enzymes, and excision of 5fC by the base excision repair enzyme thymine DNA glycosylase (TDG). Here, we provide the first genome-wide map of 5fC in mouse embryonic stem (ES) cells and evaluate potential roles for 5fC in differentiation.

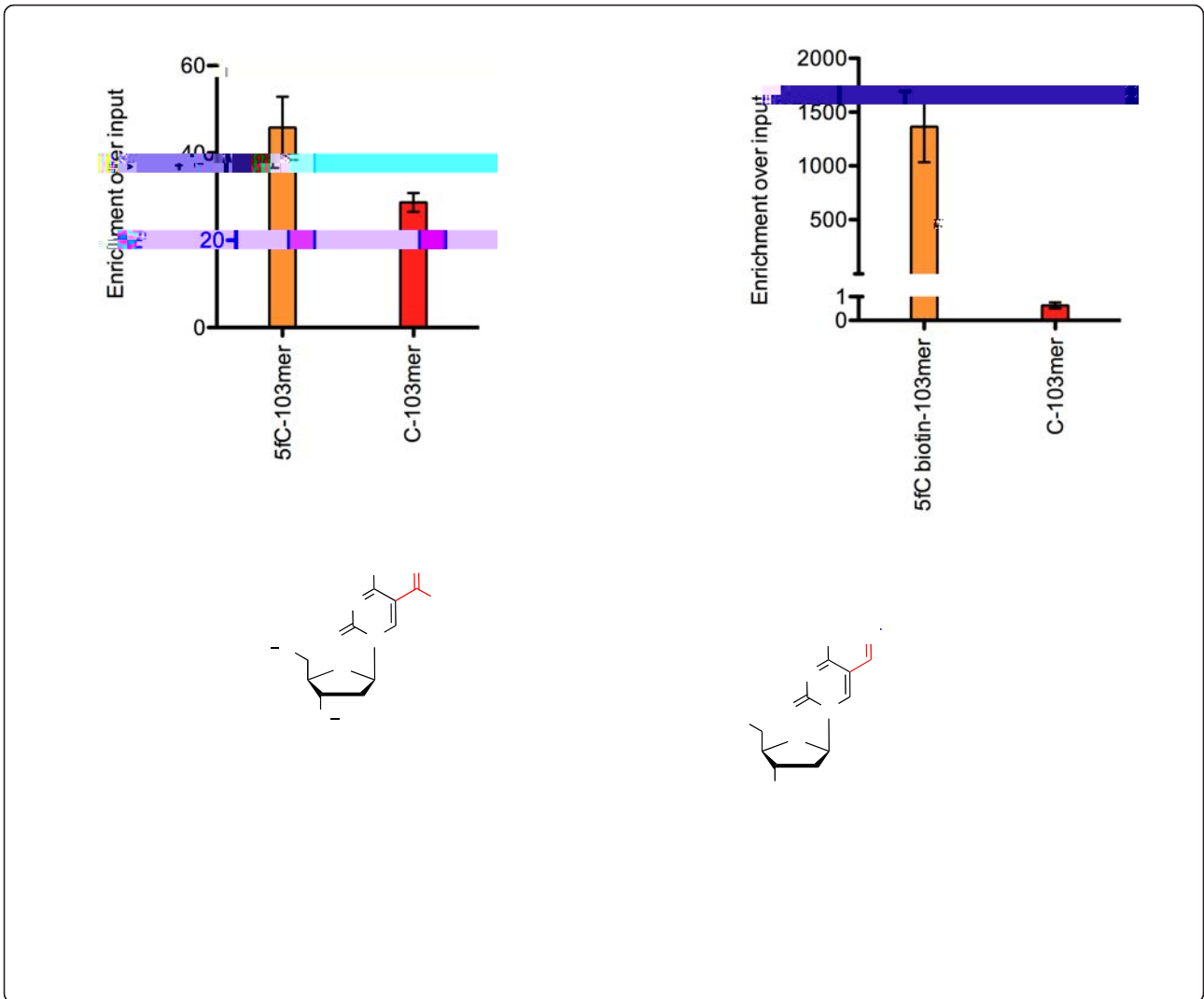
Results: Our method exploits the unique reactivity of 5fC for pulldown and high-throughput sequencing. Genome-wide mapping revealed 5fC enrichment in CpG islands (CGIs) of promoters and exons. CGI promoters in which 5fC was relatively more enriched than 5mC or 5hmC corresponded to transcriptionally active genes. Accordingly, 5fC-rich promoters had elevated H3K4me3 levels, associated with active transcription, and were frequently bound by RNA polymerase II. TDG down-regulation led to 5fC accumulation in CGIs in ES cells, which correlates with increased methylation in these genomic regions during differentiation of ES cells in wild-type and TDG knockout contexts.

Conclusions: Collectively, our data suggest that 5fC plays a role in epigenetic reprogramming within specific genomic regions, which is controlled in part by TDG-mediated excision. Notably, 5fC excision in ES cells is necessary for the correct establishment of CGI methylation patterns during differentiation and hence for appropriate patterns of gene expression during development.

Background

labeled chromatograph and tandem liquid chromatograph - mass spectrometry. Quantification of 5fC in genomic ES cell DNA showed this modified base to be present at around a level of 0.02 to 0.002% of all cytosine species, which is roughly 10- to 100-fold lower than those of 5hmC [4,5]. In ES cells, TET1 and TET2 are highly expressed and considered to play roles in reprogramming 5mC and control of the differentiation potential [6,7]. 5fC levels dramatically decrease with ongoing differentiation, suggesting its potential involvement during epigenetic reprogramming [5]. Indeed, immunostaining of zygotes that undergo global demethylation has shown that the appearance of 5fC and 5caC in the male pronucleus is associated with Tet3-mediated loss of 5mC [8].

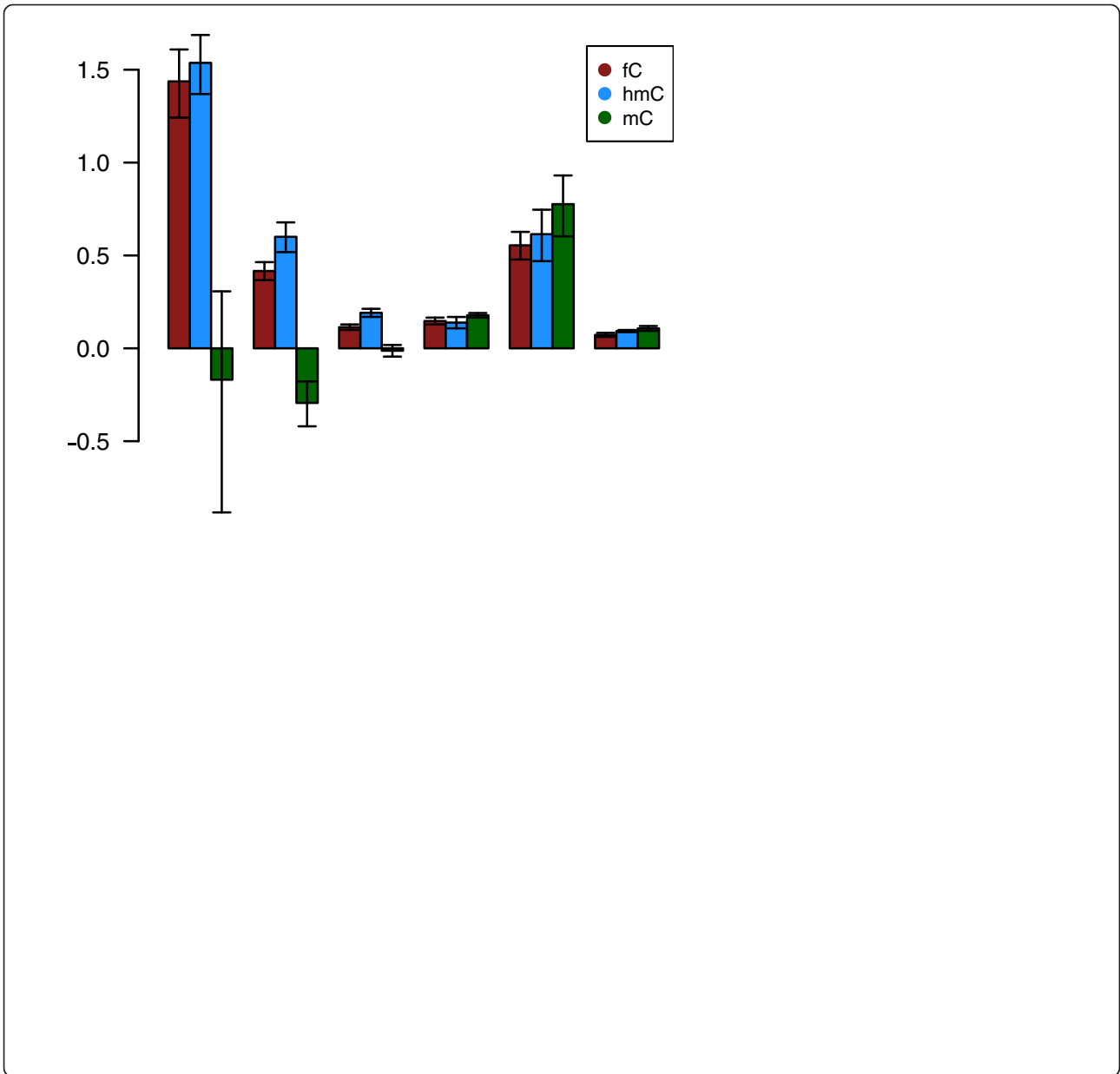
Bisulfite treatment and subsequent high throughput sequencing (BS-Seq) has been the gold standard for the detection of cytosine methylation. This method, however, does not distinguish 5mC from 5hmC or cytosine from 5fC and 5caC. Specific antibodies have been used to enrich and map 5mC (methylated DNA immunoprecipitation (MeDIP)-Seq) and 5hmC (hydroxymethylated DNA immunoprecipitation (hMeDIP)-Seq) [9]. The use of chemical labeling is an alternative method to enrich and sequence 5hmC in the genome [10,11]. The most recent breakthrough in this field came with two new methods allowing the measurement of 5hmC at single base resolution [12,13]. While various techniques for genome-wide analysis of 5mC and 5hmC are available, there is currently no method that allows the positional mapping of 5fC in



relative enrichment to the input library rather than the absolute 5fC levels. TET1 binding sites (data taken from [15]) were enriched in 5hmC and 5fC, but not 5mC, which is in accordance with the fact that TET1 is the catalyst for the generation of 5hmC and 5fC (Figure s4 in Additional file 1). The genome-wide distribution of 5fC followed a similar pattern to 5hmC with enrichments in euchromatic regions, including CpG islands (CGIs), exons and promoters (Figure 2a; Additional file 2). We also looked at the 5' UTR of LINE1 and the intracisternal A particle long terminal repeat (IAP LTR), all of which showed enrichments of 5fC in contrast to the depletion in the gene body of LINE and also other retrotransposon elements (Figure 2b; Additional file 1, Figure s5). The 5' UTR of LINE1 displayed high levels of 5hmC, medium levels of 5fC and low levels of 5mC. In contrast, IAP LTR had low levels of 5hmC, medium levels of 5fC and high levels of 5mC, demonstrating that the kinetics at each oxidation

stage depends on the genomic context. It remains to be seen if these patterns are associated with active demethylation.

The profiles shown in Figure 2c represent the enrichment levels of cytosine modifications for all genes separated into CGI- or non CGI-containing genes. In CGI-containing promoters there is a sharp enrichment peak of 5fC at the transcription start site and a slightly less localized enrichment of 5hmC with a depletion of 5mC at the transcription start site. In contrast, the profile of non-CGI promoter regions of the reference genes showed a much less pronounced increase in the levels of both 5mC and 5fC upstream of the transcription start site; these then remain at a constant level throughout the gene bodies. Overall, our analyses show that, depending on genomic regions, we observed different distributions of 5fC, 5hmC and 5mC, which suggests that the kinetics of processing 5mC are distinct between genomic regions. That 5fC is



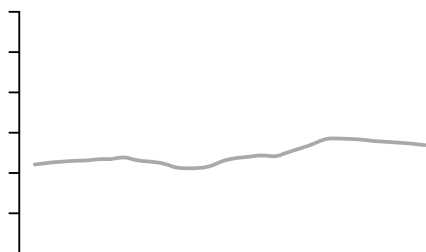
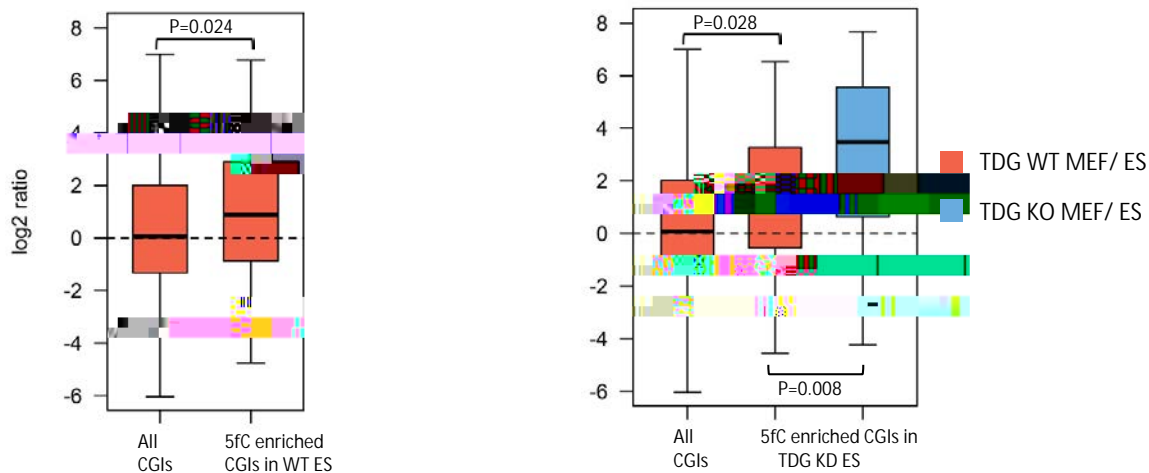
especially enriched in CGIs also supports the role of 5fC in the maintenance of histone methylation in these regions in ES cells.

5fC is associated with active gene expression in ES cells. We identified CGIs that showed a significant difference in 5fC enrichment compared to 5mC and 5hmC, and further characterized them using gene ontology categories. Therefore, we associated each island with the nearest gene within 5 kb and searched for overrepresented categories in this set. Gene ontology analysis of the 5fC-enriched genes identified pathways that were associated with transcription regulation (Table s4 in Additional file 1). We also examined

the correlation between 5fC at CGI gene promoters and their transcription levels using published gene expression data [9]. Specifically, we compared gene expression levels for cases where one cytosine modification in the CGI promoter region was relatively more enriched than other cytosine derivatives. We found that genes whose CGI promoters were 5fC-rich (relative to 5mC or 5hmC) showed higher expression than the overall group of CGI assigned genes (Figure 3a; Table s5 in Additional file 1). This suggests that the shift in equilibrium between the different cytosine modifications at promoter sites may be linked to mechanisms that control gene activity. Consistent with this observation, when genes were categorized as low, medium

levels by six-fold, consistent with its role in erasing 5fC, whereas methylation levels stayed constant (Figure s8 in Additional file 1). In general, we found that more than 98% of 5fC-enriched regions from TDG KD overlapped with those found in the siRNA control. Genome-wide 5fC mapping of the TDG KD showed 5fC-enriched sites were distributed with a reduced overall coverage of the genome (5fC sites distributed over 138 Mb in contrast to 415 Mb in the control). Thus, 5fC must be present at higher levels and/or higher density in the enriched sites for the TDG KD. This also indicates that the formation of 5fC marks at those remaining 277 Mb must be via a distinct pathway that is TDG-dependent, perhaps involving TET recruitment by TDG. It can also mean that the loss of 5fC in these particular regions is TDG-independent via an alternative pathway.

We then compared the enriched regions between TDG KD and siRNA control and found that 5% (out of 138 Mb) were significantly more highly enriched than in the



from the absence of TDG in the pluripotent stage of the early embryo may promote an even higher gain of methylation during development. In addition, TDG may also be acting in complementary pathways at these target CGIs to remove excess DNA methylation - for example, repairing mismatches resulting from 5mC deamination.

We also analyzed the 5fC distribution following siRNA-mediated down-regulation of TET1, which led to a 50% decrease in genomic 5fC as measured by mass spectrometry (data not shown). Due to the presence of TET2 in ES cells, which presumably overlaps with TET1 in binding to chromatin in many genomic regions, we concluded that

In order to control the specificity of the reaction, the same reaction was carried out on ODN1; in the absence of the oxidation step, only the starting material was recovered (Figure s1B in Additional file 1).

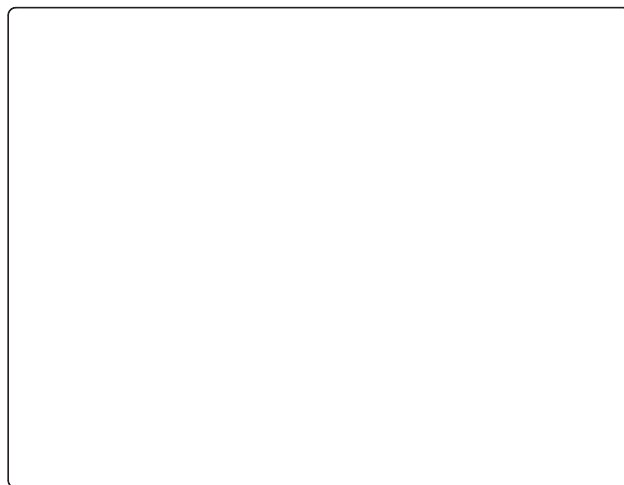
Biotin-labelling of fC in genomic DNA samples

Genomic DNA was prepared by sonicating genomic DNA extracted from mouse embryonic stem cells J1. Genomic

siGENOME non-targeting siRNA#2 (catalogue number D-001210-02; sequence not available). Cells were harvested after three rounds of transfection for DNA/RNA isolation.

Bioinformatics and data analysis

Reads in fastq format obtained from the Illumina sequencing pipeline have been aligned against the mouse genome (NCBI version mm9) using bwa [21]



the manuscript; HB carried out sample preparation, assisted with data analysis and provided feedback on the manuscript; MB assisted with data analysis and provided feedback on the manuscript; PM assisted with sample analysis and provided feedback on the manuscript; DO assisted with sample analysis; MJB

