

Inferring gene regulation from 6-base sequencing data

William Stark¹, Mark Hill¹, Jack Monahan¹, Jean Teyssandier¹, Chenfu Shi¹, Nicola Wong¹, Hugo Sepulveda³, Isaac Lopez-Moyado², Nicholas Harding¹, Anjana Rao², Páidí Creed¹

¹ biomodal Ltd, The Trinity Building, Chesterford Research Park, Cambridge, UK
² Division of Signaling and Gene Expression, La Jolla Institute for Immunology, La Jolla, CA, USA
³ Universidad Andrés Bello, Santiago, Chile

biomodal

1. Introduction

DNA methylation, an epigenetic modification, plays a key role in the regulation of gene expression. Specifically, the addition of a methyl group to the 5th carbon of cytosine (5mC) is broadly associated with transcriptional repression, with patterns of methylation established during cell fate determination constraining the transcriptional programs in cells. Demethylation via oxidation of 5-methylcytosine to 5-hydroxymethylcytosine (5hmC) is performed by TET enzymes and is reflected by higher levels of 5hmC in the gene bodies of actively transcribed genes and at active or lineage-specific enhancers.

Disentangling the roles of these two distinct modifications in gene regulation has been constrained by technological limitations, with most sequencing approaches conflating the two modifications into a single measure representing 5mC or 5hmC. Recent developments in sequencing technologies have enabled base-resolved simultaneous measurement of 5mC and 5hmC. Utilising these developments, we have generated data-sets with whole-genome 5mC and 5hmC measurements paired with RNA-sequencing data, across several different cell and tissue types. We show that we can see distinct patterns of 5mC and 5hmC across different cell types, that patterns of 5mC and 5hmC reflect tissue-specific gene expression, and that machine learning models can be trained to predict gene expression from 5mC and 5hmC.

2. duet evoC 6-base sequencing [A,C,T,G,5mC,5hmC]

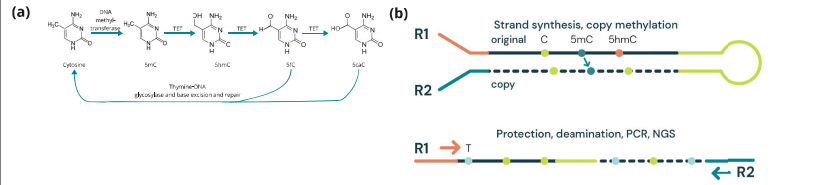
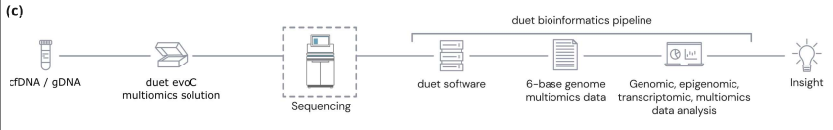


Figure 1. (a) TET-mediated demethylation pathway (b) duet multiomics solution evoC - a 6-base sequencing technology that reads all four canonical bases plus 5mC and 5hmC¹ via strand copy, 5mC copy and 5mC + 5hmC protection enzymatic steps. (c) The duet multiomics solution evoC works as an end-to-end solution comprising reagents & bioinformatics pipeline



3. Tissue-specific 5mC and 5hmC and gene expression

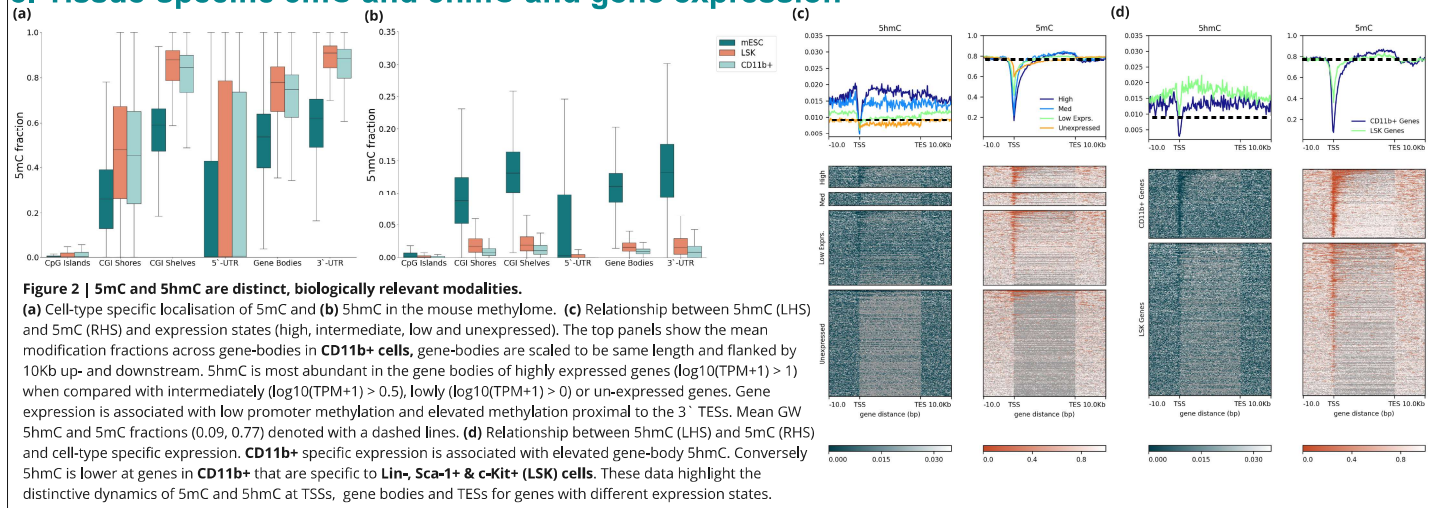
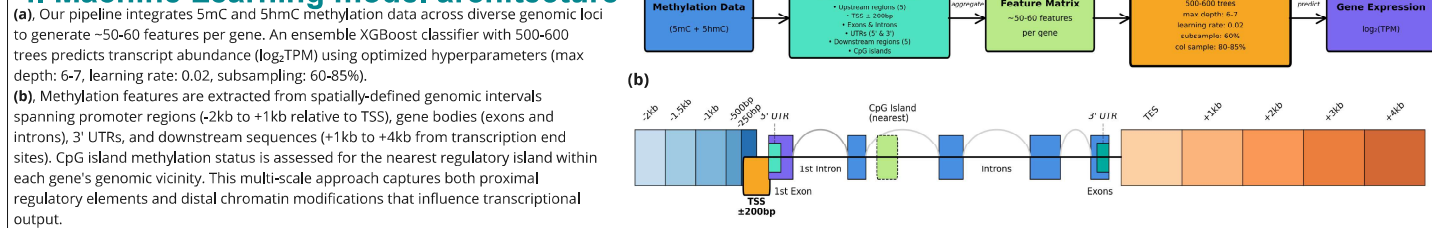
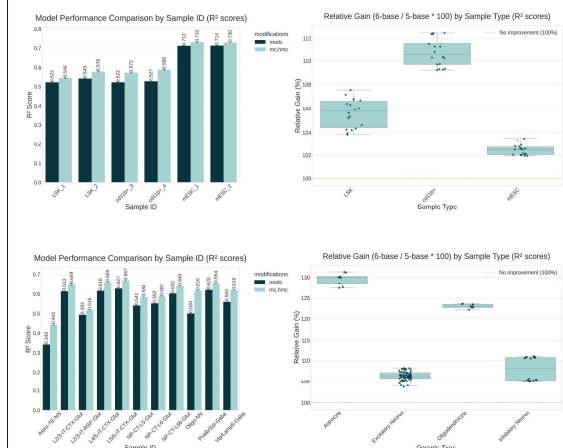


Figure 2 | 5mC and 5hmC are distinct, biologically relevant modalities. (a) Cell-type specific localisation of 5mC and (b) 5hmC in the mouse methylome. (c) Relationship between 5hmC (LHS) and 5mC (RHS) and expression states (high, intermediate, low and unexpressed). The top panels show the mean modification fractions across gene-bodies in **CD11b+** cells, gene-bodies are scaled to be same length and flanked by 10kb up- and downstream. 5hmC is most abundant in the gene bodies of highly expressed genes ($\log_{10}(\text{TPM}+1) > 1$) when compared with intermediately ($\log_{10}(\text{TPM}+1) > 0.5$), lowly ($\log_{10}(\text{TPM}+1) > 0$) or un-expressed genes. Gene expression is associated with low promoter methylation and elevated methylation proximal to the 3' TSSs. Mean GW 5hmC and 5mC fractions (0.09, 0.77) denoted with a dashed lines. (d) Relationship between 5hmC (LHS) and 5mC (RHS) and cell-type specific expression. **CD11b+** specific expression is associated with elevated gene-body 5hmC. Conversely 5hmC is lower at genes in **CD11b+** that are specific to **Lin-, Sca-1+ & c-Kit+ (LSK) cells**. These data highlight the distinctive dynamics of 5mC and 5hmC at TSSs, gene bodies and TESs for genes with different expression states.

4. Machine Learning model architecture



5. Predict RNA-seq from 6-base data



Cross-sample validation demonstrates robust predictive performance with R^2 values ranging from 0.52-0.73 across diverse cell types and tissues. The 5mC+5hmC methylation feature set consistently outperforms 5mC-only models (modC), with mean improvements of 4-12% across stem cell lines (LSK, mESC), differentiated cells (cd11b+), and primary tissues. Box plots stratified by sample type reveal that astrocytes and oligodendrocytes exhibit the highest relative performance gains (25-30%), while excitatory neurons show more modest improvements (5-7%). The consistent performance enhancement across cell types validates the biological relevance of hydroxymethylation in transcriptional regulation prediction.

6. Conclusion

We have shown that patterns of 5hmC and 5mC provide a view on tissue-specific gene expression. In particular, in **CD11b+** myeloid cells from mouse, levels of 5hmC in the gene body were shown to track with the level of gene expression and to distinguish genes which were uniquely expressed in **CD11b+** cells relative to LSK cells (a type of hematopoietic stem cell from which **CD11b+** cells differentiate). Building on this we showed that we can build models to predict gene expression using features derived from 5mC and 5hmC, with models based on this feature set consistently outperforming models with features derived from modified Cytosine data (which does not discriminate between 5mC and 5hmC). This illustrates the potential for using 6-base sequencing, which provides accurate measurements of both 5mC and 5hmC, to characterize the transcriptomic programs of cells directly from DNA. Notably, this could enable transcriptional profiling in cases where RNA is not accessible or difficult to extract, for example in cell-free DNA or FFPE-preserved tissue samples

7. References

1. Fullgrabe J. et al. Simultaneous sequencing of genetic and epigenetic bases in DNA. Nat Biotechnol. 2023 Oct;41(10):1457-1464.



biomodal