



# Quantitative modelling predicts the impact of DNA methylation on RNA polymerase II traffic

Justyna Cholewa-Waclaw<sup>a,1</sup>, Ruth Shah<sup>a,1</sup>, Shaun Webb<sup>a</sup>, Kashyap Chhatbar<sup>a</sup>, Bernard Ramsahoye<sup>b</sup>, Oliver Pusch<sup>c</sup>, Miao Yu<sup>d</sup>, Philip Greulich<sup>e,f</sup>, Bartlomiej Waclaw<sup>g,2</sup>, and Adrian P. Bird<sup>a,2</sup>

<sup>a</sup>The Wellcome Centre for Cell Biology, University of Edinburgh, EH9 3BF Edinburgh, United Kingdom; <sup>b</sup>Institute of Genetics and Molecular Medicine, University of Edinburgh, Western General Hospital Campus, EH4 2XU Edinburgh, United Kingdom; <sup>c</sup>Center for Anatomy and Cell Biology, Medical University of Vienna, 1090 Vienna, Austria; <sup>d</sup>Ludwig Institute for Cancer Research, San Diego Branch, La Jolla, CA 92093; <sup>e</sup>Mathematical Sciences, University of Southampton, SO17 1BJ Southampton, United Kingdom; <sup>f</sup>Institute for Life Sciences, University of Southampton, SO17 1BJ Southampton, United Kingdom; and <sup>g</sup>School of Physics and Astronomy, University of Edinburgh, EH9 3FD Edinburgh, United Kingdom

Contributed by Adrian P. Bird, June 4, 2019 (sent for review March 1, 2019; reviewed by Martin Howard, Amos Tanay, and Pieter Rein ten Wolde)

Patterns of gene expression are primarily determined by proteins that locally enhance or repress transcription. While many transcription factors target a restricted number of genes, others appear to modulate transcription levels globally. An example is MeCP2, an abundant methylated-DNA binding protein that is mutated in the neurological disorder Rett syndrome. Despite much research, the molecular mechanism by which MeCP2 regulates gene expression is not fully resolved. Here, we integrate quantitative, multidimensional experimental analysis and mathematical modeling to indicate that MeCP2 is a global transcriptional regulator whose binding to DNA creates “slow sites” in gene bodies. We hypothesize that waves of slowed-down RNA polymerase II formed behind these sites travel backward and indirectly affect initiation, reminiscent of defect-induced shockwaves in nonequilibrium physics transport models. This mechanism differs from conventional gene-regulation mechanisms, which often involve direct modulation of transcription initiation. Our findings point to a genome-wide function of DNA methylation that may account for the reversibility of Rett syndrome in mice. Moreover, our combined theoretical and experimental approach provides a general method for understanding how global gene-expression patterns are choreographed.

MeCP2 | gene regulation | mathematical modelling | DNA methylation

Many eukaryotic chromatin-associated factors modulate transcription by binding to specific sites in gene promoters or enhancers (1, 2). Most transcription factors are thought to modulate the initiation rate of transcription by altering histone-DNA interactions (2, 3) or imposing promoter-proximal obstacles (4). However, transcription can also be affected by processes that occur in the bodies of genes. In particular, DNA methylation, which is widespread in gene bodies, appears to affect progression of RNA polymerase II (RNA Pol II) through densely methylated exons (5). The mechanism is unclear, but methyl-CpG binding proteins (6) may be involved. Since most gene bodies contain methylated CpGs, such proteins may have a global effect on transcription.

One putative global modulator is methyl-CpG binding protein 2 (MeCP2) (7, 8), which is highly expressed in neurons. *MECP2* mutations, including loss-of-function or gene duplication, lead to severe neurological disorders (9, 10). MeCP2 does not behave as a conventional transcription factor with discrete targets, as its binding site occurs on average every ~100 base pairs (bp). Evidence from in vitro systems (11, 12) and mouse models (13, 14) suggests that MeCP2 can mediate DNA-methylation-dependent transcriptional inhibition. Transcriptional changes in mouse brain when MeCP2 is absent or overexpressed are relatively subtle but widespread (15–17), and the molecular mechanisms underlying these changes are unknown.

Here, we set out to resolve the mechanism of MeCP2-dependent transcriptional regulation. Because MeCP2 binding sites occur in the vast majority of genes, we reasoned that most are likely to be influenced to some extent by its presence. To confront the technical and analytical challenges posed by modest changes in the expression of large numbers of genes, we adopted a quantitative approach

that combined deep, high-quality datasets obtained from a uniform population of Lund Human Mesencephalic (LUHMES)-derived human dopaminergic neurons (18) with computational modeling. We created a spectrum of LUHMES cell lines expressing distinct levels of MeCP2. Using an assay for transposase-accessible chromatin sequencing (ATAC-seq) and chromatin immunoprecipitation sequencing (ChIP-seq) together with mathematical modeling, we detected a robust footprint of MeCP2 binding to mCG in vivo and determined the amount of MeCP2 bound to DNA. Quantification of mRNA abundance by RNA-sequencing (RNA-seq) revealed a relationship between changes in transcription and the density of mCG on gene bodies. To explain this observation, we proposed and tested several distinct mechanistic models. The only model consistent with our experimental results was one in which MeCP2 leads to slowing down of RNA Pol II progression through a transcription unit. Importantly, mutant MeCP2 that is unable to bind the TBL1/TBLR1 subunits of the NCoR

## Significance

We introduce an interdisciplinary approach to understanding global modulation of gene expression in mammalian cells. Conventional transcription factors target a limited subset of genes, whereas global modulators bind the genome broadly. An example of the latter is MeCP2, which is mutated in neurological disorders. MeCP2 has millions of genomic binding sites, but its effects on gene expression are mostly small scale and incompletely understood at a mechanistic level. Using datasets from genetically modified human neurons, our mathematical approach rigorously distinguishes global effects from experimental noise. This allows us to integrate theory with experiments to discriminate competing mechanistic models. The results indicate that MeCP2 creates “roadblocks” in gene bodies that slow down elongating RNA polymerase II, leading to polymerase queuing.

Author contributions: J.C.-W., R.S., B.W., and A.P.B. designed research; J.C.-W., R.S., B.R., and M.Y. performed research; J.C.-W., R.S., and O.P. contributed new reagents/analytic tools; J.C.-W., R.S., S.W., K.C., P.G., B.W., and A.P.B. analyzed data; and J.C.-W., R.S., B.W., and A.P.B. wrote the paper.

Reviewers: M.H., John Innes Centre; A.T., Weizmann Institute; and P.R.t.W., University of Amsterdam.

The authors declare no conflict of interest.

This open access article is distributed under [Creative Commons Attribution License 4.0 \(CC BY\)](https://creativecommons.org/licenses/by/4.0/).

Data deposition: The data reported in this paper have been deposited in the Gene Expression Omnibus (GEO) database, [www.ncbi.nlm.nih.gov/geo](http://www.ncbi.nlm.nih.gov/geo) (accession no. GSE125660). Computer programs, scripts, and data related to mathematical models have been deposited at the Edinburgh Data Share database, <https://doi.org/10.7488/ds/2568>.

<sup>1</sup>J.C.-W. and R.S. contributed equally to this work.

<sup>2</sup>To whom correspondence may be addressed. Email: [bwaclaw@ph.ed.ac.uk](mailto:bwaclaw@ph.ed.ac.uk) or [a.bird@ed.ac.uk](mailto:a.bird@ed.ac.uk).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1903549116/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1903549116/-DCSupplemental).

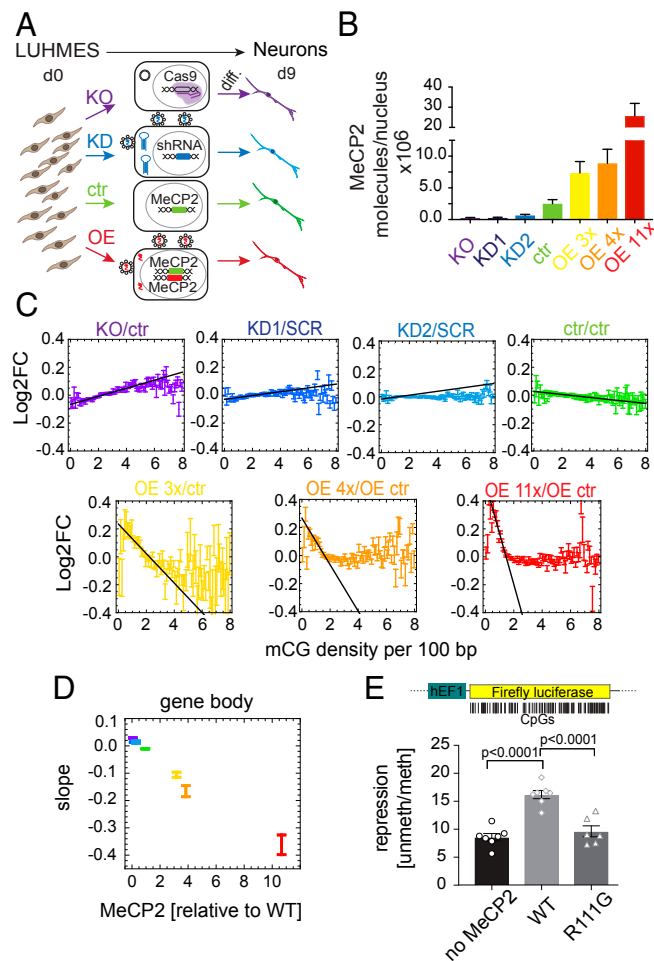
Published online July 9, 2019.

corepressor complex fails to repress efficiently, suggesting that repression depends upon this interaction.

## Results

### Global Changes in Transcription Correlate with MeCP2 Expression Level.

We created progenitor cell lines capable of differentiation to a uniform population of human neurons (*SI Appendix, Fig. S1 A–C*) that expressed seven widely different levels of MeCP2, including knockout (KO), wild-type (WT), and 11-fold overexpression (OE 11x) (Fig. 1 *A* and *B* and *SI Appendix, Fig. S1D* and *Table S1*). All lines differentiated into neurons with similar kinetics, expressed neuronal markers (*SI Appendix, Fig. S1E*), and had identical global levels of DNA methylation (~3.7% of all cytosines were methylated) (*SI Appendix, Fig. S2A*). Based on the known affinity of MeCP2 for methylated CG (mCG), we expected that the effect of MeCP2 on gene expression would depend on their mCG content.



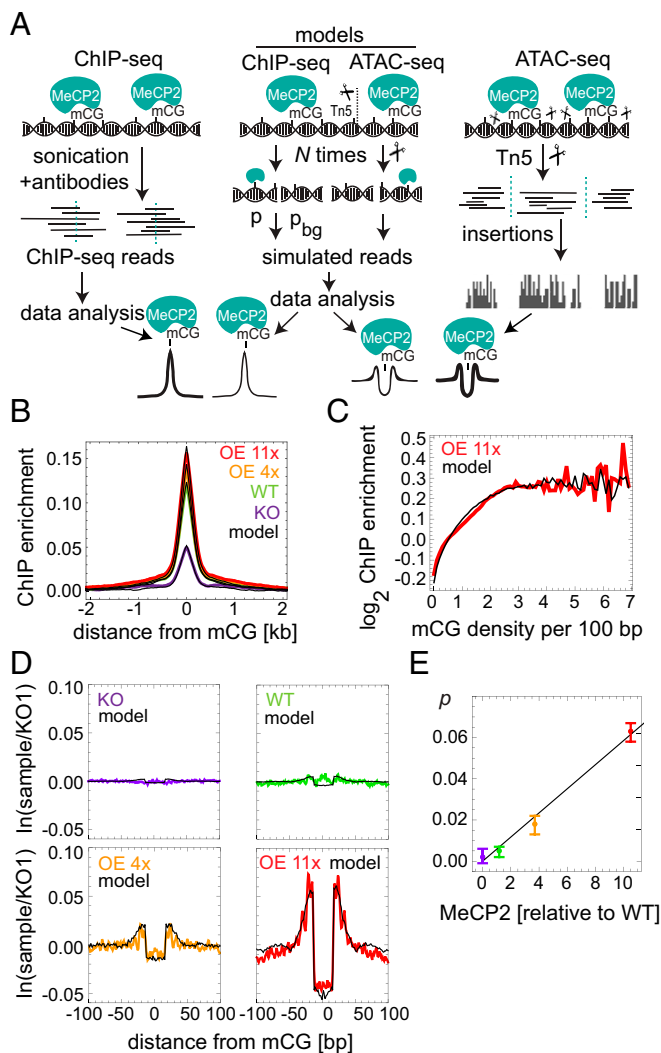
**Fig. 1.** Gene expression strongly correlates with gene body mCG density and MeCP2 abundance. (A) Experimental design (*Materials and Methods*). d0, day 0; d9, day 9. (B) Mean number of MeCP2 molecules per nucleus. (C) Log<sub>2</sub> fold change of gene expression (Log<sub>2</sub>FC) relative to appropriate controls (ctr, unmodified controls; SCR, scrambled shRNA control; OE ctr, over-expression control) for all seven levels of MeCP2, plotted against gene body mCG density. All Log<sub>2</sub>FC values have been shifted so that Log<sub>2</sub>FC averaged over all genes is zero. Black line indicates the maximum slope. (D) The maximum slope for gene bodies varies proportionally to MeCP2 abundance. (E) Ratio between luciferase expressions from an unmethylated and gene-body methylated constructs, for three cases: no MeCP2, WT MeCP2, and a methyl-CpG binding domain mutant R111G that is unable to bind mCG. Points show individual replicates. In all panels, error bars represent  $\pm$ SEM.

DNA methylation was therefore quantified for all genes in WT neurons by using whole-genome bisulfite sequencing [Tet-assisted bisulfite sequencing (TAB-seq)] (*SI Appendix, Fig. S2 B and C*). We calculated total methylation (total mCG,  $N_{\text{mCG}}$ ) as the number of mCG dinucleotides, mCG density ( $\rho_{\text{mCG}}$ ) as the number of mCGs per 100 bp, and mCG mean as the percentage of mCG in all CG dinucleotides. To determine the effects of MeCP2 on transcription, we performed RNA-seq on all seven cell lines. We included all expressed protein-coding genes (~17,000 genes) in our analysis. Most genes responded to MeCP2, but changes were small, precluding definition of a subset of affected genes (*SI Appendix, Fig. S3A*). To enhance a possible relationship between expression changes and DNA methylation that otherwise might be obscured by other regulatory mechanisms and statistical noise, genes were binned according to methylation density, considering gene bodies and promoters separately.

The average change in expression vs. appropriate controls [ $\log_2$  fold change of gene expression (Log<sub>2</sub>FC)] showed a strong relationship to mCG density ( $\rho_{\text{mCG}}$ ) in gene bodies (Fig. 1C). The effect was the strongest for  $\rho_{\text{mCG}} = 0.8\text{--}4.0$  mCG per 100 bp, which includes the vast majority of genes (*SI Appendix, Fig. S3B*). The apparent stimulation of expression at very low mCG densities in OE neurons is discussed in *SI Appendix*. Moreover, the maximum slope of the Log<sub>2</sub>FC vs.  $\rho_{\text{mCG}}$  in gene bodies (Fig. 1C, black lines) was strikingly proportional to MeCP2 levels (Fig. 1D). In contrast, plots of Log<sub>2</sub>FC vs.  $\rho_{\text{mCG}}$  in promoter regions showed a slope close to zero, indicating minimal dependence on promoter methylation (*SI Appendix, Fig. S3C*). No clear dependence on MeCP2 level was observed for Log<sub>2</sub>FC vs. total gene body mCG or mCG mean (*SI Appendix, Fig. S3 D and E*). These results indicated that the gene-body mCG density is the strongest predictor of MeCP2-dependent transcriptional changes. This relationship was not affected when data were filtered by significance, gene length, or promoter methylation (*SI Appendix, Fig. S4 A–D*). Moreover, the relationship was maintained even when intronic reads were analyzed, suggesting that pre-mRNA is affected in the same way as processed RNA (*SI Appendix, Fig. S4E*). To test for a causal relationship, we transfected cells with two versions (methylated or unmethylated gene body) of a luciferase reporter gene with a methylation-free promoter in the presence of WT or the DNA binding mutant MeCP2 [R111G] (*SI Appendix, Fig. S5 A and B*). We observed a twofold repression of methylated vs. unmethylated luciferase gene body in the presence of WT MeCP2, compared with either no MeCP2 or mutant MeCP2 (Fig. 1E).

**MeCP2 Binds Predominantly mCG Genome-Wide.** To map the binding of MeCP2 in human neurons, we performed MeCP2 ChIP-seq for KO, WT, OE 4x, and OE 11x and developed a computer model that simulated the ChIP-seq procedure and MeCP2 binding in vivo (Fig. 2A). As expected, ChIP enrichment was proportional to the level of MeCP2 in each cell line (*SI Appendix, Fig. S6 A–C*) and showed a strong peak centered at the mCGs in MeCP2-positive lines (Fig. 2B), as well as a correlation between MeCP2 enrichment and mCG density (Fig. 2C). Conversely, enrichment was absent at nonmethylated CGs (*SI Appendix, Fig. S6E*).

To derive an independent measure of absolute MeCP2 density on the DNA and to detect its molecular footprint with high resolution, we performed ATAC-seq (19), in which transposase Tn5 cuts exposed DNA to reveal DNA accessibility within chromatin (Fig. 2A). In agreement with the ChIP-seq data, ATAC-seq Tn5 insertion profiles (Fig. 2D) showed a graded depletion of insertion sites centered around mCG in WT, OE 4x, and OE 11x neurons, whose amplitude was proportional to MeCP2 concentration (Fig. 2E) and therefore represented a “molecular footprint” of MeCP2 binding in vivo. The size and amplitude of the footprint agreed well with a computer model of ATAC-seq and MeCP2 binding (Fig. 2D, black lines) and previous in vitro data (20, 21), confirming that



**Fig. 2.** MeCP2 occupancy on the DNA is proportional to mCG density and MeCP2 level. (A) MeCP2 ChIP- and ATAC-seq experimental procedures and their in silico counterparts.  $p_{bg}$  and  $p$  are probabilities of background and MeCP2-bound reads, respectively. Tn5 insertion sites (scissors) occur in exposed DNA regions. (B) ChIP-seq enrichment profiles centered at mCG dinucleotides for different cell lines. Black lines represent in silico profiles fitted to the experimental data. (C) MeCP2 ChIP-seq enrichment data in OE 11x/KO (red) as a function of mCG density. (D) Average depletion profiles (logarithm of the ratio between the number of Tn5 insertions in a given cell line and KO1; two to four biological replicates) in the  $\pm 100$ -bp regions surrounding mCG dinucleotides. Black lines represent computer simulations of the model fitted to the data. (E) Predicted fraction of mCGs occupied by MeCP2 vs. MeCP2 level obtained from depletion profiles in D. Error bars represent  $\pm 5$ EM.

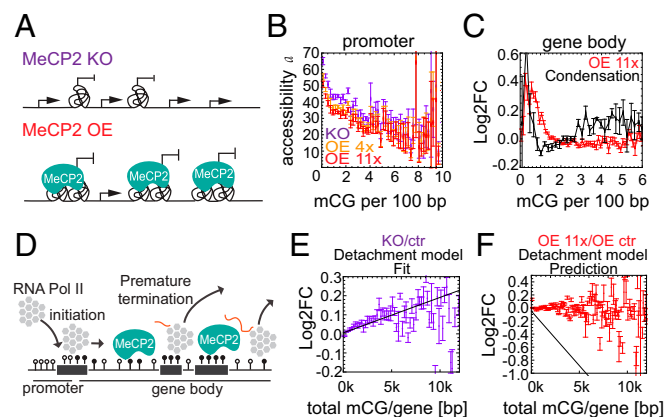
MeCP2 occupies 11 bp of DNA in living cells. No depletion of insertion sites was observed over unmethylated CG (SI Appendix, Fig. S6F). The model revealed that only 6.3% of mCG sites are actually occupied by MeCP2 in OE 11x neurons, falling to <1% occupancy in WT (Fig. 2E), perhaps due in part to occlusion by nucleosomes. Excellent agreement between the models and ATAC-seq and ChIP-seq data allowed us to predict MeCP2 occupancy from mCG density and MeCP2 level in each cell line (Fig. 2E and SI Appendix, Fig. S6D).

**MeCP2 Does Not Regulate Transcription via Condensation of Chromatin or Premature Termination.** To interpret these results mechanistically, we considered mathematical models based on a commonly accepted paradigm for gene expression (SI Appendix,

Fig. S7A) (22). In the first class of models, named condensation models (Fig. 3A), MeCP2 affects the rate of transcription initiation via changes in chromatin structure. The possibility that MeCP2 affects the initiation rate  $\alpha$  by binding to promoters was rejected because it would imply a stronger correlation between gene expression and  $\rho_{mCG}$  in promoters than in gene bodies, contrary to our observations (SI Appendix, Fig. S3C). MeCP2 could hypothetically affect the fraction  $f$  of cells with specific genes in the ON state via some long-distance mechanism involving binding to gene bodies and leading to changes in the degree of chromatin openness near promoters. However, mapping chromatin accessibility by using ATAC-seq showed that, while there is a weak correlation between MeCP2 and accessibility (Fig. 3B), it cannot account for the observed Log2FC in gene expression (Fig. 3C).

We next considered potential effects of MeCP2 on the elongation phase of transcription. The detachment model posits that MeCP2 causes transcription to prematurely abort (Fig. 3D). Since the probability of termination increases with each blocking site, under this model, the Log2FC is a function of the total number of mCGs ( $N_{mCG}$ ) in the gene:  $\text{Log2FC} = -\gamma (M - 1)N_{mCG}$ , where  $M$  is MeCP2 concentration relative to WT, and the parameter  $\gamma$  is proportional to the probability that Pol II aborts transcription when it encounters MeCP2 or an MeCP2-induced chromatin modification. The unknown parameter  $\gamma$  can be obtained by fitting the model to the Log2FC (KO/WT) data (Fig. 3E). We found that the model failed to reproduce the Log2FC vs.  $N_{mCG}$  relationship for the OE 11x cell line (Fig. 3F). The model also failed to correctly predict the observed relationship between Log2FC and mCG density in gene bodies (SI Appendix, Fig. S7B and C). Therefore, it is unlikely that MeCP2 affects transcription via premature termination.

**MeCP2 Creates “Dynamical Obstacles” That Impede Transcriptional Elongation.** Finally, we considered a “congestion model,” whereby Pol II pauses when it encounters MeCP2 itself or an induced, transient structural modification of chromatin (Fig. 4A). The parameters were: the fraction  $p$  of mCGs bound by MeCP2, MeCP2 turnover (unbinding) rate  $k_u$ , and (specific to each gene) the length  $L$  of the gene, the density  $\rho_{mCG}$  of mCGs, and the initiation rate  $\alpha$ . Fig. 4B shows the transcription rate for OE 11x predicted by the



**Fig. 3.** MeCP2 does not regulate transcription via condensation of chromatin or premature termination. (A) A cartoon of the condensation model. Tangles represent regions of condensed chromatin that are inaccessible to RNA Pol II. (B) Chromatin accessibility (measured by ATAC-seq) at promoters rapidly decreases with increasing promoter methylation. In contrast, MeCP2 has a minor effect on accessibility (curves for OEs 4x and 11x are slightly lower than for KO). (C) The condensation model disagrees with Log2FC(OE 11x/KO) obtained from RNA-seq. (D) Schematic representation of the detachment model. (E) Log2FC (gene expression) for KO/ctr (purple) vs. the total number of mCGs per gene. Black lines represent predictions of the detachment model. Error bars represent  $\pm 5$ EM. (F) As in E for OE 11x/OE ctr (red).

model as a function of  $\alpha$ , for different mean MeCP2 densities ( $\rho\rho_{mCG}$ ). The assumed value of  $k_u = 0.04 \text{ s}^{-1}$  is compatible with the reported in vivo residence time of MeCP2 on chromatin [25–40 s (23)]. Inspired by nonequilibrium statistical mechanics approaches that have been utilized to model one-dimensional transport (24, 25), we expected a nonequilibrium phase transition from a low-density to a maximal-current (congested) phase as the initiation rate or the density of obstacles increase beyond a critical point. Indeed, all curves in Fig. 4B have a characteristic shape: a linear relationship  $J \sim \alpha$  for small  $\alpha$ , followed by saturation at high initiation rates. Saturation occurs due to congestion as polymerases queue upstream of obstacles (Movies S1 and S2). However, even in the nonsaturated regime of intermediate  $\alpha$ , excluded-volume interactions between polymerases that have been slowed down by an obstacle cause a density shockwave that propagates backward (Fig. 4C). A small increase in the density of polymerases near the promoter decreased the rate of Pol II binding to the transcription start site (TSS). Thus, even though MeCP2 does not directly affect Pol II initiation, it does so indirectly by shockwaves that form behind MeCP2-induced obstacles in gene bodies (Fig. 4D). To test the model against RNA-seq data, we estimated average initiation rates for genes with similar mCG densities by fitting the model to Log2FC data from one of the cell lines [OE 11x/OE control (ctr); Fig. 4E, Left and SI Appendix, Fig. S8F]. We then used the model to predict Log2FC for the remaining six cell lines. The model strikingly reproduced the data (Fig. 4E for OE 4x and KO) as well

as the slopes of the Log2FC plots for all seven cell lines (Fig. 4F). A similar behavior occurred in a modified model in which Pol II slowed down (rather than completely stopped) on permanent or long-lasting structural modifications of chromatin (SI Appendix, Fig. S8A–E and Movie S3). We conclude that both congestion models are compatible with the experimental data presented in Fig. 1C and D. The models also predict that Log2FC should decrease with increasing expression (measured as transcripts per million reads), in agreement with the data (SI Appendix, Fig. S8G).

### MeCP2 Binding to Both DNA and NCoR Are Essential to Slow Down RNA Pol II.

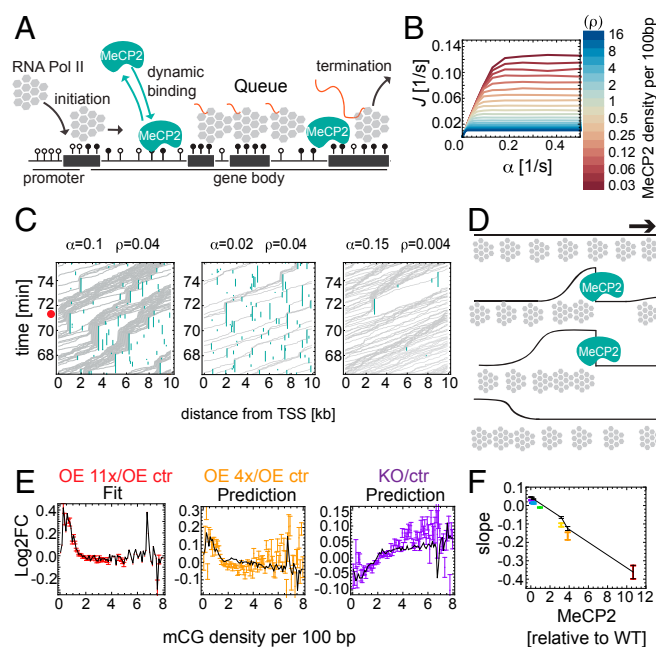
To address the question of whether MeCP2 impedes Pol II progression directly by steric interference or indirectly by altering chromatin structure [e.g., by histone deacetylation (26)], we overexpressed mutated forms of MeCP2 in the presence of WT MeCP2. The mutants were either unable to bind methylated DNA (R111G) (27) or unable to recruit the histone deacetylase complex NCoR (R306C) (14, 28) (Fig. 5A and SI Appendix, Fig. S9A). As expected, sevenfold overexpression of MeCP2–R111G caused no mCG-density-dependent transcriptional changes (Fig. 5B and C and SI Appendix, Fig. S9B and C). The R306C mutant, on the other hand, was predicted to repress transcription if inhibition is directly due to MeCP2 binding to DNA, but not if inhibition is mediated via the corepressor. In fact, 11-fold overexpression of MeCP2–R306C relative to WT MeCP2 caused only a small perturbation of gene expression, indicating a significant loss of DNA methylation-dependent repression (Fig. 5B and D and SI Appendix, Fig. S9B and C). The weak slope may represent minor direct interference of DNA-bound MeCP2–R306C with transcription. As neither mutant falls on the line defining the linear relationship between gene repression and MeCP2 concentration (Fig. 5E), our findings favor a predominantly indirect mechanism of repression, whereby corepressor recruitment alters the chromatin state to impede transcription.

### Concluding Remarks

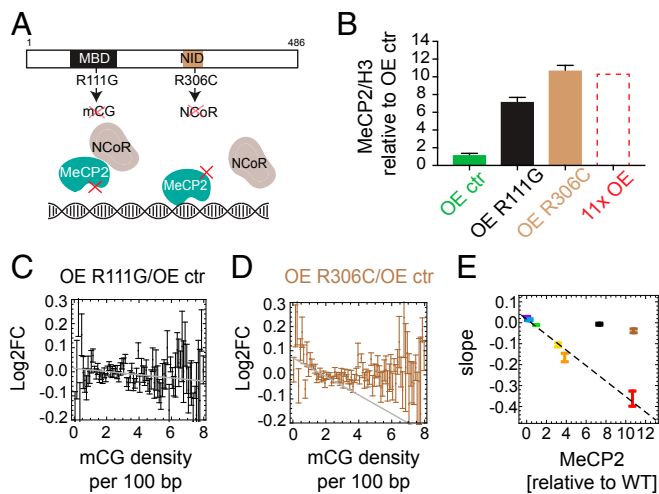
In summary, a close alliance between mathematical modeling and molecular biology has allowed us to discriminate molecular mechanisms underlying the relatively subtle global effects of MeCP2 on global gene expression. The proposed mechanism relies on MeCP2–NCoR interaction that slows down the progression of Pol II during transcription elongation. A candidate mediator of this effect is histone modification, in particular histone deacetylation, as cell-transfection assays using methylated reporters demonstrate that repression depends upon histone deacetylase activity (11, 12). According to this scenario, MeCP2 recruitment of the histone deacetylase corepressor NCoR would restrain transcription, perhaps by causing tighter binding of nucleosomes to DNA (26). To explain the dramatic reversibility of Rett syndrome in animal models (29), we propose that, in the absence of MeCP2, DNA methylation patterns are unaffected, allowing the reexpressed WT protein to bind within gene bodies and commence normal modulation of transcriptional elongation. We suggest that the congestion model may apply to proteins other than MeCP2. For example, other chromatin-binding factors that bind short (and thus abundant) motifs, including other methyl-binding proteins, may modulate gene expression by a similar mechanism.

### Materials and Methods

**Cell Lines.** The procedure for culture and differentiation of the LUHMES cell line was described (18). To create two independent *MECP2* KO lines, we used CRISPR-mediated gene disruption (30). To generate MeCP2 knockdowns, several shRNAs against MeCP2 were designed by using Sigma-Aldrich Mission shRNA online software. Two shRNAs were chosen and cloned into a pLKO.1 vector including scrambled shRNA as a control, and lentiviruses were created (SI Appendix, Table S2). To increase the level of MeCP2, we created lentiviruses expressing MeCP2 from two alternative promoters in the pLKO.1 vector: synapsin and cytomegalovirus. Calculation of SD, SEM, and *t* tests for



**Fig. 4.** Mathematical modeling indicates that MeCP2 slows down transcriptional elongation. (A) Schematic representation of the dynamical obstacles model. (B) Transcription rate  $J$  predicted by the model, plotted as a function of the initiation rate  $\alpha$ , for different mean MeCP2 densities in gene bodies. (C) Space–time plots (kymographs) representing Pol II moving along the gene. Queues of Pol II induced by MeCP2 can reach the TSS (red dot) and block initiation if both the initiation rate ( $\alpha$ ) and the density of MeCP2 ( $\rho$ ) are sufficiently high (C, Left). (D) Schematic representation of Pol II (gray) density shockwaves forming behind MeCP2 (blue). Black line is the local density of Pol II. (E) Log2FC (gene expression) vs. mCG density in gene bodies obtained in computer simulations of the dynamical obstacles model (black solid lines) fitted to the OE 11x/OE ctr RNA-seq dataset (red) agrees well with experimental data for OE 4x/OE ctr (orange) and KO/ctr (purple) datasets. Error bars represent  $\pm$ SEM. (F) The maximum slope of Log2FC (gene expression) vs. mCG density in gene bodies, predicted by the dynamical obstacles model (black line). Points are experimental slopes from Fig. 1C.



**Fig. 5.** MeCP2 slows down transcription via a mechanism involving NCoR. (A) Location of two binding domains in MeCP2 that are relevant for the proposed mechanism: methyl-CpG binding domain (MBD) and NCoR-interaction domain (NID). The mutation R111G causes MeCP2 to lose the ability to bind specifically to mCG. The mutation R306C prevents MeCP2 from binding the NCoR complex. (B) Level of MeCP2 (Western blot) in two overexpressed mutant cell lines (R111G and R306C) and the overexpression control cell line (OE ctr). OE 11x is shown for comparison. Values are averaged over three biological replicates and normalized by the level of histone H3. (C) Log<sub>2</sub>FC (expression) of OE R111G/OE ctr shows almost no dependence on mCG density in gene bodies (black). The gray line shows the maximum slope. (D) Log<sub>2</sub>FC (expression) of OE R306C/OE ctr shows a small negative correlation with gene body mCG density (brown). The gray line shows the maximum slope. (E) Maximum slopes for all cell lines including OE R111G (black) and OE R306C (brown) from C and D vs. MeCP2 level (Western blot). In all plots, error bars represent  $\pm$ SEM.

qPCR, Western blots, methylation, and total RNA quantification using high-performance liquid chromatography (HPLC) were performed by using GraphPad Prism (version 7).

**Repression Assay.** CpG-free vector containing Firefly Luciferase with CpGs was methylated by M.Sss1 methyltransferase in the presence or absence of S-adenosyl-L-methionine. Mouse embryonic fibroblasts were transfected by using Lipofectamine 2000 with three plasmids containing Firefly Luciferase, Renilla Luciferase, and MeCP2. Luciferase activity measurements were performed by using the Dual Luciferase assay kit (Promega), according to manufacturer's protocol.

**Library Preparation for Illumina Sequencing.** All libraries were sequenced as 75- or 100-nucleotide-long paired-end reads on HiSeq 2000 and HiSeq 2500 Illumina platforms. Methyome of WT LUHMES-derived neurons at day 9 was obtained by TAB-seq according to the published protocol (31). The RNA-seq library was prepared according to manufacturer's protocol for the ScriptSeq Complete Gold kit (human/mouse/rat). Total RNA was isolated from all generated cell lines (SI Appendix, Table S1) at day 9 of differentiation by using either the RNeasy Mini kit or the AllPrep DNA/RNA Mini kit (Qiagen). ATAC-seq in four cell lines (KO, WT, OE 4x, and OE 11x; SI Appendix, Table S1) was performed as in ref. 32.

MeCP2 ChIP-seq was performed by using LUHMES-derived neurons at day 9 of differentiation with four levels of MeCP2: KO, WT, OE 4x, and OE 11x (SI Appendix, Table S1). Libraries were prepared by using the NEBNext Ultra II DNA library Prep kit (NEB) for both immunoprecipitations and corresponding inputs.

**Data Processing of Raw Reads from Illumina Sequencing.** All reads were quality-controlled, trimmed to remove adapters (Trimomatic) (33) and duplicated reads, and mapped to the human hg19 reference genome. Bismark (34) was used to extract cytosine methylation from TAB-seq. All raw data were deposited in the Gene Expression Omnibus database (accession no. GSE125660) (35).

**RNA-Seq Data Analysis.** We used a subset of protein-coding genes with sufficient methylation coverage (bisulfite sequencing;  $\geq 80\%$  C with coverage  $\geq 20$ ),

and gene bodies 1 kb or longer. This resulted in 15,382 genes of the initial 17,764 protein-coding genes (86%). In all plots of Log<sub>2</sub>FC of differential gene expression, we shifted the Log<sub>2</sub>FC values so that the average Log<sub>2</sub>FC in the range of mCG density  $\rho_{mCG} \in [1,6]/100$  bp was zero for all samples. This was motivated by the difficulty in determining the absolute levels of expression, since we did not quantify total mRNA.

**ChIP-Seq Enrichment Profiles.** We first obtained accumulated counts (the number of reads)  $c_i^x$  that overlapped with  $i$ -th base pair to the right ( $i > 0$ ) or left ( $i < 0$ ) from feature  $x$  ( $x = \text{mCG, non-mCG, ...}$ ). We then calculated enrichment profiles as

$$f_i = \frac{\text{Norm}_1(c_i^{\text{ChIP},x}[i])}{\text{Norm}_1(c_i^{\text{input},x}[i])} - 1,$$

where  $c_i^{\text{ChIP},x}$  and  $c_i^{\text{input},x}$  are accumulated counts from ChIP and input (genomic) DNA sequencing, respectively, and  $\text{Norm}_1(c)[i]$  normalizes the counts profiles such that their flanks have values close to one:

$$\text{Norm}_1(c)[i] = \frac{c_i}{\left(\sum_{j=-301}^{-500} c_j + \sum_{j=301}^{500} c_j\right) / 400}.$$

We considered a particular C to be methylated if it was methylated in 100% of the reads, and the coverage was at least 5. We considered a C to be unmethylated if it did not show up in any of the ChIP-seq reads as methylated.

**Computer Model of ChIP-Seq.** We assumed that MeCP2 occupies methylated cytosines with probability  $p$  times the probability of binding to a particular motif. Binding probabilities for different motifs are based on known binding affinities (36) and relative binding strengths (15). To create simulated ChIP fragments, we assumed that if a DNA fragment contained at least one MeCP2 bound to it, it would be present in the simulated ChIP-seq. Fragments that do not contain any MeCP2 may still be present in the ChIP-seq data with probability  $p_{bg}$ , which accounts for "background" reads in ChIP-seq, even in the absence of MeCP2. This is similar to previous models of ChIP-seq (37); even the best ChIP-seq libraries can have a significant level of background reads ( $p_{bg}$  close to 1) (38). We also added CG and length bias and processed simulated reads in the same way as the experimental ChIP data.

For each ChIP-seq dataset, we fitted the simulated profile (parametrized by  $p, p_{bg}$ ) to the experimental profile. Any  $p \leq 0.1$  gives a good fit (SI Appendix, Fig. S6D), indicating that  $p \sim 0.1$  is the upper bound on mCG occupancy in 11x OE. We used best-fit parameters to predict profiles on features other than mCG (SI Appendix, Fig. S6E).

**ATAC-Seq Footprints.** ATAC-seq was analyzed in a similar way to ChIP-seq, except that we used fragments' endpoints (Tn5 insertion sites) to generate accumulated counts  $n_i$ . We calculated the insertion profiles as

$$f_i = \ln \left[ \frac{\text{Norm}_2(n_i^{\text{cell line}})}{\text{Norm}_2(n_i^{\text{KO}})} \right],$$

where  $n_i^{\text{cell line}}$  and  $n_i^{\text{KO}}$  are the insertion counts profiles for a given cell line and KO1, respectively, and  $\text{Norm}_2$  normalizes the counts profiles such that their flanks have values close to one:

$$\text{Norm}_2(n_i) = \frac{n_i}{\left(\sum_{j=-41}^{-50} n_j + \sum_{j=41}^{50} n_j\right) / 20}.$$

**Computer Model of ATAC-Seq.** We used the same binding model as in the ChIP-seq simulations. We assumed that MeCP2 occupies 11 bp (20) and that the protein is centered on an mC. We simulated the action of the Tn5 transposase by splitting the sequence into fragments in areas free of MeCP2, and we included Tn5 sequence bias and CG and length bias. The model has three parameters: the density  $p$  of MeCP2 on mCxx, the average density of insertion (cut) sites  $t$ , and the GC bias  $b$ . We processed simulated DNA fragments in the same way as described above for the experimental data. We examined the role of the parameters on the shape and depth of the simulated footprint of MeCP2 and concluded that the footprint is not affected as long as the test and control samples have been processed in a similar way. To extract MeCP2 occupancy  $p$  from ATAC-seq data, we fitted the model (free parameters  $p, t$ , and a fixed  $b = 6.0$ ) to experimental footprints for all four cell lines. The relationship is linear (Fig. 2E), with the best-fit  $p = 0.0058 \times M_{\text{cell line}}/M_{\text{WT}}$ .

**Chromatin Accessibility from ATAC-Seq.** For each gene, we calculated its mean insertion count  $\bar{n}$  and selected regions ("insertion peaks") in which  $n_i > 4\bar{n}$ . Accessibility was defined as the sum of all insertions in the peaks divided by the "background"  $\bar{n}$ :

$$a = \frac{\sum_i n_i}{\bar{n}}$$

**Mathematical Models of Gene Expression.** The condensation model assumes that the fraction  $f_i$  of cells in which gene  $i$  is actively transcribed depends on promoter openness  $a_i$  (measured by ATAC-seq), which in turn depends on the level  $M$  of MeCP2 and gene methylation  $\rho_i$ :  $f_i = f_i(M, \rho_i) \propto a_i = a(M, \rho_i)$ . The model predicts that  $\text{Log2FC}_{X/Y}$  of the ratio of gene expression of cell line X vs. cell line Y should yield the same curve (plus a constant) as the logarithm of the ratio of accessibilities of X vs. Y when plotted as a function of  $\rho_{\text{mCG}}$ . Data did not support this model (Fig. 3C). The detachment model poses that the probability that RNA Pol II successfully terminates is  $P = (1 - \lambda)^n \cong e^{-\lambda n}$ , where  $n$  is the number of "abort sites" on the gene, proportional to the number of MeCP2 molecules on the gene, and  $\lambda$  is the abortion probability. We show that

$$\text{Log2FC}_{X/Y} = \text{const} - \gamma \left( \frac{M_X}{M_Y} - 1 \right) n,$$

where  $\gamma \propto \lambda$  is an unknown parameter identical for all cell lines, and  $M_X, M_Y$  are MeCP2 levels in cell lines X and Y. The model was rejected (Fig. 3F).

We considered two mechanisms by which MeCP2 could affect elongation. To implement the slow-sites model, we used the totally asymmetric simple exclusion process (TASEP) with open boundaries (24). A gene is represented as a chain of  $L$  sites. Each site (equivalent to 60 bp of the DNA) is either occupied by a particle (RNA Pol II) or is empty. Particles enter the chain at site  $i = 1$  with rate  $\alpha$  (transcription initiation rate), move along the chain, and exit at site  $i = L$  with rate  $\beta = 1 \text{ s}^{-1}$ . Sites can be "fast" or "slow." Slow sites represent mCGs affected by the interaction with MeCP2, whereas fast sites are all other sites (methylated or not). Particles jump with rate  $v = 1 \text{ s}^{-1}$  (equivalent to Pol II speed  $\sim 60 \text{ bp/s}$ ) on fast sites and  $v_s = 0.05 \text{ s}^{-1}$  on slow sites. Slow sites are randomly and uniformly distributed with density  $\rho_s =$

$\rho_{\text{mCG}}$ , where  $p$  is the probability that an mCG is occupied by MeCP2. To relate this model to the mRNA-seq differential expression data, we calculated  $\text{Log2FC}$  as

$$\text{Log2FC}_{X/Y} = \log_2 \frac{J(\alpha, \rho_{s,X})}{J(\alpha, \rho_{s,Y})},$$

where  $\rho_{s,X} = \rho_{\text{mCG}} p_X$ ,  $\rho_{s,Y} = \rho_{\text{mCG}} p_Y$ , in which  $p_X, p_Y$  are MeCP2 occupation probabilities for cell lines X, Y. In the above expression, we know all quantities except the initiation rate  $\alpha$ , which we fit to the OE 11x data.

The dynamical obstacles model is very similar, with two exceptions: (i) Pol II always moves with the same speed  $v$  (no slow sites) as long as it is not blocked by other polymerases and obstacles; and (ii) obstacles bind and unbind dynamically from the methylated sites. We assumed that unbinding occurs with rate  $k_u$  per obstacle, whereas binding occurs with rate  $k_b \rho$  per unoccupied mCG. Obstacles do not bind if an mCG is already occupied by an obstacle or a polymerase. We assumed that obstacles are not restricted to accessible mCGs and that their density on actively transcribed genes may be higher than  $p$  obtained from ATAC-seq, but still proportional to MeCP2 level. We found that  $p = M/M_{\text{OE11x}}$  reproduces  $\text{Log2FC}$  data for all cell lines. Computer programs, scripts, and data related to mathematical models have been deposited at the Edinburgh Data Share database (39).

Additional details for materials and methods are provided in *SI Appendix*.

**ACKNOWLEDGMENTS.** We thank Tanja Waldmann for introducing us to the LUHMES cell line, Beatrice Alexander-Howden for technical support, David Kelly for microscopy assistance, Martin Waterfall for help with fluorescence-activated cell sorting, Jim Selfridge for help with preparing samples for HPLC, and Sabine Lager and John Connelly for critical reading of the manuscript. The work has made use of resources provided by the Edinburgh Compute and Data Facility (<https://www.ed.ac.uk/information-services/research-support/research-computing/ecdf>) and was supported by a Wellcome Trust Programme Grant and Investigator Award (to A.P.B.). A.P.B. is a member of the Simons Initiative for the Developing Brain. J.C.-W. was supported by a grant from the Rett Syndrome Research Trust. R.S. was supported by a Wellcome Trust 4-y PhD studentship. B.W. was supported by a Royal Society of Edinburgh Personal Research Fellowship.

1. C. Y. Lin *et al.*, Transcriptional amplification in tumor cells with elevated c-Myc. *Cell* **151**, 56–67 (2012).
2. S. Berry, C. Dean, M. Howard, Slow chromatin dynamics allow polycomb target genes to filter fluctuations in transcription factor activity. *Cell Syst.* **4**, 445–457.e8 (2017).
3. K. Ouararhni *et al.*, The histone variant mH2A1.1 interferes with transcription by down-regulating PARP-1 enzymatic activity. *Genes Dev.* **20**, 3324–3336 (2006).
4. N. Hao, A. C. Palmer, I. B. Dodd, K. E. Shearwin, Directing traffic on DNA—How transcription factors relieve or induce transcriptional interference. *Transcription* **8**, 120–125 (2017).
5. I. Jonkers, J. T. Lis, Getting up to speed with transcription elongation by RNA polymerase II. *Nat. Rev. Mol. Cell Biol.* **16**, 167–177 (2015).
6. B. Hendrich, A. Bird, Identification and characterization of a family of mammalian methyl-CpG binding proteins. *Mol. Cell Biol.* **18**, 6538–6547 (1998).
7. J. D. Lewis *et al.*, Purification, sequence, and cellular localization of a novel chromosomal protein that binds to methylated DNA. *Cell* **69**, 905–914 (1992).
8. P. J. Skene *et al.*, Neuronal MeCP2 is expressed at near histone-octamer levels and globally alters the chromatin state. *Mol. Cell* **37**, 457–468 (2010).
9. R. E. Amir *et al.*, Rett syndrome is caused by mutations in X-linked MECP2, encoding methyl-CpG-binding protein 2. *Nat. Genet.* **23**, 185–188 (1999).
10. M. B. Ramocki, Y. J. Tavyev, S. U. Peters, The MECP2 duplication syndrome. *Am. J. Med. Genet. A.* **152A**, 1079–1088 (2010).
11. X. Nan, F. J. Campoy, A. Bird, MeCP2 is a transcriptional repressor with abundant binding sites in genomic chromatin. *Cell* **88**, 471–481 (1997).
12. X. Nan *et al.*, Transcriptional repression by the methyl-CpG-binding protein MeCP2 involves a histone deacetylase complex. *Nature* **393**, 386–389 (1998).
13. J. Guy, B. Hendrich, M. Holmes, J. E. Martin, A. Bird, A mouse MeCP2-null mutation causes neurological symptoms that mimic Rett syndrome. *Nat. Genet.* **27**, 322–326 (2001).
14. M. J. Lyst *et al.*, Rett syndrome mutations abolish the interaction of MeCP2 with the NCoR/SMRT co-repressor. *Nat. Neurosci.* **16**, 898–902 (2013).
15. S. Lagger *et al.*, MeCP2 recognizes cytosine methylated tri-nucleotide and di-nucleotide sequences to tune transcription in the mammalian brain. *PLoS Genet.* **13**, e1006793 (2017).
16. B. Kinde, D. Y. Wu, M. E. Greenberg, H. W. Gabel, DNA methylation in the gene body influences MeCP2-mediated gene repression. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 15114–15119 (2016).
17. H. W. Gabel *et al.*, Disruption of DNA-methylation-dependent long gene repression in Rett syndrome. *Nature* **522**, 89–93 (2015).
18. D. Scholz *et al.*, Rapid, complete and large-scale generation of post-mitotic neurons from the human LUHMES cell line. *J. Neurochem.* **119**, 957–971 (2011).
19. J. D. Buenrostro, P. G. Giresi, L. C. Zaba, H. Y. Chang, W. J. Greenleaf, Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome position. *Nat. Methods* **10**, 1213–1218 (2013).
20. X. Nan, R. R. Meehan, A. Bird, Dissection of the methyl-CpG binding domain from the chromosomal protein MeCP2. *Nucleic Acids Res.* **21**, 4886–4892 (1993).
21. T. Nikitina *et al.*, MeCP2-chromatin interactions include the formation of chromatosome-like structures and are altered in mutations causing Rett syndrome. *J. Biol. Chem.* **282**, 28237–28245 (2007).
22. V. Shahrezaei, J. F. Ollivier, P. S. Swain, Colored extrinsic fluctuations and stochastic gene expression. *Mol. Syst. Biol.* **4**, 196 (2008).
23. R. J. Klose *et al.*, DNA binding selectivity of MeCP2 due to a requirement for AT sequences adjacent to methyl-CpG. *Mol. Cell* **19**, 667–678 (2005).
24. R. A. Blythe, M. R. Evans, Nonequilibrium steady states of matrix-product form: A solver's guide. *J. Phys. A Math. Theor.* **40**, R333–R441 (2007).
25. B. Derrida, An exactly soluble non-equilibrium system: The asymmetric simple exclusion process. *Phys. Rep.* **301**, 65–83 (1998).
26. G. E. Zentner, S. Henikoff, Regulation of nucleosome dynamics by histone modifications. *Nat. Struct. Mol. Biol.* **20**, 259–266 (2013).
27. S. Kudo *et al.*, Heterogeneity in residual function of MeCP2 carrying missense mutations in the methyl CpG binding domain. *J. Med. Genet.* **40**, 487–493 (2003).
28. V. Krusvee *et al.*, Structure of the MeCP2-TBLR1 complex reveals a molecular basis for Rett syndrome and related disorders. *Proc. Natl. Acad. Sci. U.S.A.* **114**, E3243–E3250 (2017).
29. J. Guy, J. Gan, J. Selfridge, S. Cobb, A. Bird, Reversal of neurological defects in a mouse model of Rett syndrome. *Science* **315**, 1143–1147 (2007).
30. R. R. Shah *et al.*, Efficient and versatile CRISPR engineering of human neurons in culture to model neurological disorders. *Wellcome Open Res.* **1**, 13 (2016).
31. M. Yu *et al.*, Base-resolution analysis of 5-hydroxymethylcytosine in the mammalian genome. *Cell* **149**, 1368–1380 (2012).
32. J. D. Buenrostro, B. Wu, H. Y. Chang, W. J. Greenleaf, *ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide* (John Wiley & Sons, Inc., Hoboken, NJ, 2015).
33. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
34. F. Krueger, S. R. Andrews, Bismark: A flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* **27**, 1571–1572 (2011).
35. J. Cholewa-Waclaw, *et al.*, Study of MeCP2 in LUHMES-derived neurons. Gene Expression Omnibus. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE125660>. Deposited on 25 January 2019.
36. M. J. Sperlazza, S. M. Bilinovich, L. M. Sinanop, F. R. Javier, D. C. Williams, Jr, Structural basis of MeCP2 distribution on non-CpG methylated and hydroxymethylated DNA. *J. Mol. Biol.* **429**, 1581–1594 (2017).
37. H. Xu *et al.*, A signal-noise model for significance analysis of ChIP-seq with negative control. *Bioinformatics* **26**, 1199–1204 (2010).
38. K. Liang, S. Keleş, Normalization of ChIP-seq data with control. *BMC Bioinformatics* **13**, 199 (2012).
39. B. Waclaw, *et al.*, Supplementary data for the manuscript "Quantitative modelling predicts the impact of DNA methylation on RNA polymerase II traffic." Edinburgh DataShare, University of Edinburgh. <https://doi.org/10.7488/ds2568>. Deposited 7 June 2019.