

# 1 **normR: Regime enrichment calling for ChIP-seq data**

2 Johannes Helmuth<sup>1</sup>, Na Li<sup>1</sup>, Laura Arrigoni<sup>2</sup>, Kathrin Gianmoena<sup>3</sup>, Cristina Cadenas<sup>3</sup>, Gilles  
3 Gasparoni<sup>4</sup>, Anupam Sinha<sup>5</sup>, Philip Rosenstiel<sup>5</sup>, Jörn Walter<sup>4</sup>, Jan G. Hengstler<sup>3</sup>, Thomas  
4 Manke<sup>2</sup> and Ho-Ryun Chung<sup>1,\*</sup>

5 <sup>1</sup> Otto-Warburg-Laboratory: Epigenomics at Max Planck Institute for Molecular Genetics,  
6 Ihnestrasse 63-73, 14195 Berlin, Germany.

7 <sup>2</sup> Max Planck Institute of Immunobiology and Epigenetics, Stübeweg 51, 79108 Freiburg,  
8 Germany

9 <sup>3</sup> Leibniz Research Centre for Working Environment and Human Factors at the TU  
10 Dortmund, Ardeystrasse 67, 44139 Dortmund

11 <sup>4</sup> The Department of Genetics and Epigenetics, University of Saarland, Campus A2.4 66123  
12 Saarbrücken, Germany.

13 <sup>5</sup> Institute of Clinical Molecular Biology, Christian-Albrechts-University of Kiel,  
14 University Hospital Schleswig Holstein - Campus Kiel, Schittenhelmstr. 12, 24105 Kiel,  
15 Germany

16 \* Corresponding author: Ho-Ryun Chung, Max Planck Institute for Molecular Genetics,  
17 Ihnestrasse 63-73, 14195 Berlin, Germany, Tel.: + 49 - 30 - 84 13 1122, Fax: + 49 - 30 - 84  
18 13 1960 Email: [ho-ryun.chung@molgen.mpg.de](mailto:ho-ryun.chung@molgen.mpg.de)

19 **Running title:** normR regime enrichment calling

20 **Keywords:** Peak Calling, ChIP-seq data Normalization, ChIP-seq Regime Identification,  
21 ChIP-seq Heterogeneity, ChIP-seq Difference Calling

22

## 23 **Abstract**

24 ChIP-seq probes genome-wide localization of DNA-associated proteins. To mitigate technical  
25 biases ChIP-seq read densities are normalized to read densities obtained by a control. Our  
26 statistical framework “normR” achieves a sensitive normalization by accounting for the effect  
27 of putative protein-bound regions on the overall read statistics. Here, we demonstrate  
28 normR’s suitability in three studies: (i) calling enrichment for high (H3K4me3) and low  
29 (H3K36me3) signal-to-ratio data; (ii) identifying two previously undescribed H3K27me3 and  
30 H3K9me3 heterochromatic regimes of broad and peak enrichment; and (iii) calling  
31 differential H3K4me3 or H3K27me3-enrichment between HepG2 hepatocarcinoma cells and  
32 primary human Hepatocytes. normR is readily available on  
33 <http://bioconductor.org/packages/normr>

34

## 35 **Introduction**

36 Chromatin Immunoprecipitation followed by high-throughput sequencing (ChIP-seq; [1]) is a  
37 widely used method for the genome-wide localization of DNA-associated proteins, such as  
38 transcription factors or histone modifications. In brief, after crosslinking with formaldehyde  
39 the chromatin is sheared and the resulting chromatin fragments are enriched by  
40 immunoprecipitation for the protein of interest. The precipitate is reverse-crosslinked to  
41 obtain DNA fragments, which are amplified and then sequenced. The reads generated in this  
42 way are then aligned to a reference genome and genomic loci bound by the protein are  
43 inferred by an accumulation of sequencing reads. Due to the genome-wide scalability and  
44 cost-efficiency of ChIP-seq, hundreds of distinct proteins and their modifications have been  
45 assayed to study underlying mechanisms of molecular function in different cell types [2,3].  
46 Consequently, a huge resource of protein location information is available to be readily  
47 integrated into studies at hand.

48 ChIP-seq data are used to characterize transcription factor binding sites [4], chromatin  
49 landscapes [5,6] or functional elements, like enhancers [7,8]. Specifically, most ChIP-seq  
50 experiments aim to study protein binding sites in the context of gene regulation. For example,  
51 the lineage-specific binding of transcription factors orchestrates differentiation pathways [9].  
52 Furthermore, ChIP-seq signals of histone modifications are predictive for promoter activity  
53 [10] and enhancer competence [11].

54 The identification of regions bound by a protein of interest requires the discrimination of  
55 enrichment against background. Intuitively, a high number of ChIP sequencing reads should  
56 map to protein-bound regions, where the average number of reads in these regions depends on  
57 the “binding mode” of the protein of interest. For example, transcription factors and certain  
58 histone modifications, such as H3K4me3, are characterized by a localized read accumulation  
59 with a high signal-to-noise ratio [12]. Some histone modifications, such as H3K9me3,  
60 H3K27me3, and H3K36me3 are characterized by a more delocalized read accumulation with

61 a substantially lower signal-to-noise ratio.

62 Technical biases introduced during the ChIP-seq procedure lead to accumulation of reads in  
63 regions that are devoid of the protein [13-15]. These biases arise by copy number variations,  
64 sequencing biases, mapping ambiguities, and the chromatin structure [13,16]. These biases  
65 are also discernable in control experiments, *i.e.* they can be accounted for by comparing the  
66 ChIP read coverage to a control experiment without specific enrichment, such as the input  
67 chromatin to the ChIP.

68 The comparison of the read counts in the ChIP to those in the control requires normalization  
69 to account for, both, the differences in the sequencing depth, and the effects of enrichment by  
70 the ChIP. Ideally, such a normalization should yield a normalization factor that corrects the  
71 average ratio between ChIP- and control read counts in background regions [17-19]. Thus, a  
72 proper normalization requires the identity of background regions. On the other hand, the  
73 discrimination of enriched and background regions requires normalization itself –  
74 normalization and discrimination of enrichment against background are two faces of the same  
75 coin.

76 Earlier approaches estimate the normalization factor either by the ratio of sequencing depths  
77 (*e.g.* MACS [20] and DFilter [21]), by the ratio of ChIP- and control read counts summed  
78 over *ad hoc* chosen background regions with fixed width (*e.g.* CisGenome [22], SPP [23] and  
79 MUSIC [24]), or by identifying background regions and their width using a data-driven  
80 approach (*e.g.* NCIS [17] or SES [19]). After normalization these approaches identify  
81 enriched regions and equate them to protein binding sites or modifications. All these  
82 approaches discriminate a single signal regime from the background. However, a qualitative  
83 separation of this signal regime, *e.g.* into moderately and highly enriched regimes, could  
84 distinguish genomic loci that are bound by the protein in only a subpopulation of cells in the  
85 sample from those that are bound in the majority of cells in the sample. Those analyses of  
86 ChIP-seq sample heterogeneity cannot be performed using existing methods.

87 The discrimination of signal against background is not only required to determine protein

88 binding sites it is also required for identifying regions that are differentially bound in two  
89 conditions, *e.g.* control and disease. Most approaches (*e.g.* [20]) aimed at identifying  
90 differentially bound regions concentrate on the modeling of condition-specific exclusive  
91 enrichment. In addition, other methods [25-27] employ a three-state Hidden Markov Model to  
92 additionally identify condition-specific changes of signal within regions of concurrent ChIP  
93 enrichment. Therein, a computationally intensive training is done to learn a hidden state  
94 representation of the data. Consequently, the regional ChIP read coverage is  
95 “interpolated”/”smoothed” based on the read coverage in adjacent genomic loci. This data  
96 abstraction sacrifices a statistically sound null hypothesis.

97 Here, we describe a data-driven robust and broadly applicable approach for simultaneous  
98 normalization and difference calling in ChIP-seq data called normR (recursive acronym:  
99 “normR obeys regime mixture rules”). normR models ChIP- and control read counts by a  
100 binomial mixture model. One component models the background, while one or more other  
101 components model the signal. As a proof of principle, normR is applied in three scenarios:  
102 Firstly, we show that normR achieves robust enrichment calling for both high (H3K4me3)  
103 and low (H3K36me3) signal-to-noise ratio ChIP-seq data. High specificity and sensitivity of  
104 normR is confirmed by functional outputs like gene expression and DNA methylation state.  
105 Secondly, we use normR to characterize two previously undetectable enrichment regimes for  
106 H3K27me3 and H3K9me3 in hepatocarcinoma HepG2 cells. Finally, the translational normR  
107 approach is shown to confidently call differences between primary human hepatocytes and  
108 HepG2 cells for both high (H3K4me3) and low (H3K27me3) signal-to-noise ratio histone  
109 modification ChIP-seq data. Here, we uncover potential epigenetic alterations introduced by  
110 the cancer-associated immortalization of primary liver cells. Thus, normR is a versatile tool  
111 that can identify enriched regions, distinct enrichment regimes and differences between  
112 conditions using a simple binomial mixture model and robust statistics.

## 113 **Results**

### 114 **The normR Framework**

115 During a ChIP experiment antibodies are used to enrich chromatin fragments carrying the  
116 protein of interest from a population of fragments obtained by sonication of chromatin. These  
117 antibodies bind preferentially but not exclusively to protein-DNA complexes. Hence, ChIP  
118 only enriches rather than selects protein containing chromatin fragments. Bearing this in mind,  
119 ChIP can be envisioned as a sampling process where the probability to draw a fragment  
120 depends on the presence or absence of the protein. If present, the probability is high, if absent,  
121 the probability is lower but not zero. The spatial distribution of the fragments sampled in this  
122 way is then estimated by mapping the sequenced ends (reads) of these fragments to the  
123 reference genome.

124 To infer regions bound by a protein of interest the read densities obtained by ChIP-seq  
125 experiment are compared to the corresponding counts obtained by a control experiment *e.g.*  
126 by sequencing the sonicated chromatin (input). A region should be called “enriched by the  
127 ChIP” only if the number of reads from the ChIP is sufficiently greater than that expected  
128 relative to the control. Such an approach addresses a number of systematic biases, like copy  
129 number variations, sequencing biases, mapping ambiguities and chromatin structure  
130 [13,15,16]. To this end, a proper normalization of the read count densities is essential: For  
131 example, if we sequence twice as many reads in the ChIP than in the control, the read counts  
132 per region in the ChIP should be greater than in the control. In the absence of enrichment by  
133 the ChIP, we expect twice as many reads per region in the ChIP than in the control. In the  
134 presence of enrichment by the ChIP, the read counts in the region associated with the protein  
135 should be much higher than in the control, but what happens to the read counts in the  
136 remaining regions?

137 Sequencing the ChIP and control libraries is a multinomial sampling process, which induces  
138 dependencies between the regions. As the total number of reads obtained from one

139 sequencing run is fixed and finite, the increase of reads in some regions due to ChIP  
140 enrichment leads to a decrease in remaining regions, *i.e.* background regions  $B$ . Returning to  
141 our example, this implies that the number of reads in non-enriched regions in the ChIP should  
142 be less than twice the number from the control. In particular, the normalization factor  $c_B$  is  
143 less than two which relates the number of reads in ChIP-seq  $s_i$  to the ones in control  $r_i$  by  
144  $s_i \approx c_B \times r_i$  for background regions  $i \in B$ .  $c_B$  depends on the average enrichment achieved  
145 by the ChIP and the number of enriched regions — it shrinks as, both, the number of  
146 enriched regions and the level of enrichment in these regions increases. Critically,  $c_B$  is  
147 required to define a statistically sound Null hypothesis for testing whether the observed ChIP  
148 read counts are sufficiently greater than expected given the control. Moreover, the more  
149 regions are enriched, the lower the signal-to-noise ratio becomes at a fixed sequencing depth  
150 [12]. The estimation of  $c_B$  requires the identity of background regions, albeit the identification  
151 of the background requires normalization itself. Thus, ChIP-seq normalization and the  
152 identification of enriched regions are two sides of the same problem.

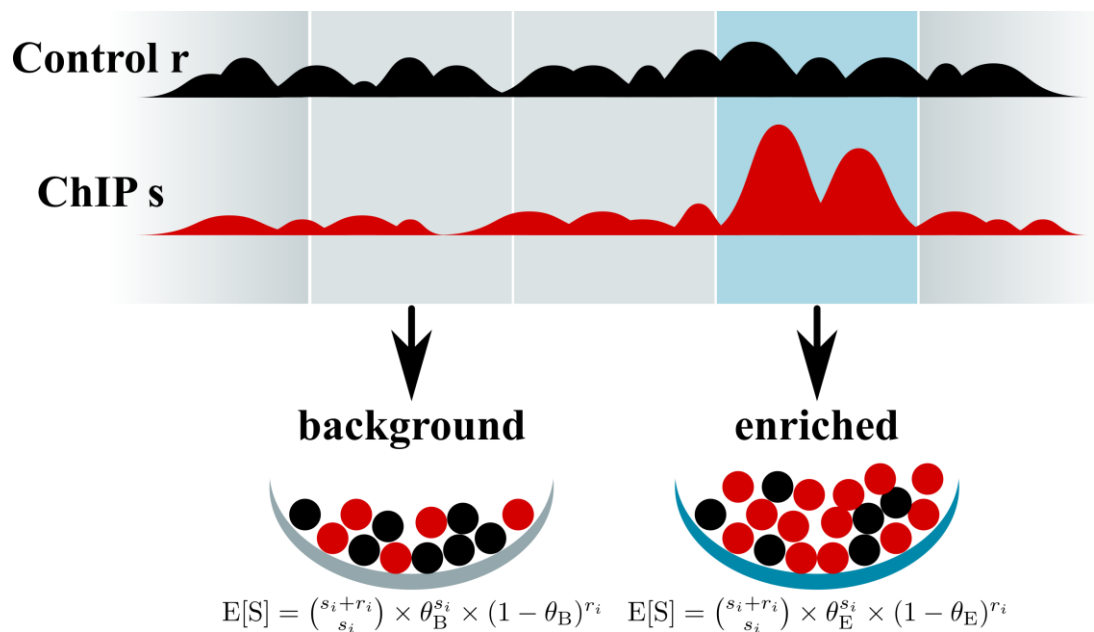
153 To tackle this problem we model the read counts from the ChIP and control by a binomial  
154 mixture model (Methods; Figure 1). In its simplest incarnation we use two components, *i.e.*  
155 background and enriched, to normalize and call enrichment over the control (referred to as  
156 “enrichR”). The model has in total three free parameters, *i.e.*  $\theta_B, \theta_E$  and  $\pi_B$ .  $\theta_B$  and  
157  $\theta_E$  represent the expected fraction of reads in the ChIP over the sum of reads from ChIP and  
158 control per region for the background and the enriched regions, respectively.  $\pi_B$  is the  
159 proportion of regions that belong to the background  $\pi_B$  (the proportion of regions that are  
160 enriched is simply  $\pi_E = 1 - \pi_B$ ). Given this model we derive the following likelihood  
161 function:

$$\mathcal{L} = \prod_i \binom{s_i + r_i}{s_i} (\pi_B \times \theta_B^{s_i} \times (1 - \theta_B)^{r_i} + (1 - \pi_B) \times \theta_E^{s_i} \times (1 - \theta_E)^{r_i}),$$

162 where  $s_i$  ( $r_i$ ) corresponds to the number of reads in the ChIP (control) for regions  $i = 1, \dots, n$ .

163 We fit these parameters using the expectation-maximization algorithm [28] on the closed

164 form solution (Methods). From the discussion above we expect that  $\theta_B \leq \frac{N_{ChIP}}{N_{ChIP} + N_{control}} = \theta^*$ ,  
 165 where  $N_{ChIP}$  ( $N_{control}$ ) is the total number of reads in the ChIP (control) and  $\theta^*$  denotes the  
 166 expected fraction of reads from ChIP-seq taking into account only sequencing depth. Equality  
 167 holds only in case of no enrichment, or  $\pi_B = 1$ . The last implicit “parameter” is the definition  
 168 of regions. We use non-overlapping fixed width regions because it is robust and appropriate  
 169 for most downstream analyses [5,29-31].



171 **Figure 1. The normR Framework.** Reads in control  $r$  and ChIP  $s$  are modeled as a binomial  
 172 mixture model with multiple components. Here, two components model the expected fraction  
 173 of reads in the ChIP over the sum of reads from ChIP and control per region for background  
 174  $\theta_B$  and the enriched  $\theta_E$ . By accounting for the effect of ChIP enrichment on the background  
 175 read statistics a statistical sound Null hypothesis is formed.

176 The identification of enriched regions across the genome is based on the fitted model: Given  
 177 the control read count, the ChIP read count in each region is compared to the expected read  
 178 count under the fitted background model. Using a binomial test statistically significant  
 179 deviations from the background model are recovered. The null distribution of p-values from a  
 180 binomial test is discrete and impedes the correction for multiple testing. By filtering out low  
 181 power tests (*i.e.* low count regions) with the T method [32], the p-value distribution becomes  
 182 more uniform and the p-values can be adjusted for multiple testing. Filtered p-values are then



183 transformed to q-values [33]. Enriched regions are reported if they fall below a user-specified  
184 threshold.

185 In addition to enrichR, we provide two augmented realizations of normR (Methods): (i)  
186 “regimeR” models multiple enrichment components defined by  $\theta_{E_j}$  with  $j = 1, \dots, m$  to  
187 identify ChIP enrichment regimes; and (ii) “diffR” models the expected fraction of reads in a  
188 depleted (control-enriched) component defined by  $\theta_D$  in addition to  $\theta_B$  and  $\theta_E$  yielding a  
189 direct comparison of two ChIP experiments. After assessing significance against  $\theta_B$  every  
190 region is assigned to a component by *Maximum a posteriori* assignment.

191 Based on the fitted binomial mixture model the normalized ChIP signal  $e_i$  is calculated by  
192 dividing the read counts from ChIP-seq by those from the control scaled by the normR  
193 enrichment factor  $f = \frac{\theta_E}{1-\theta_E} \times \frac{1-\theta_B}{\theta_B}$ . To account for noise in low power regions, we  
194 regularize  $e_i$  by adding pseudocounts to the number of ChIP-seq and Input-seq reads resulting  
195 in

$$e_i = \frac{\ln\left(\frac{s_i + \alpha_s}{r_i + \alpha_r} \times \frac{\alpha_r}{\alpha_s}\right)}{\ln(f)}$$

196 where  $\alpha_x = \frac{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i} \times x}{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i}}$  represents the average read count for  $x$  given the normR-fitted  
197 background model.

198 We have implemented normR in C++ and R [34]. normR is available on Bioconductor at  
199 <http://bioconductor.org/packages/normr>.

## 200 **Enrichment Calling in Low and High Signal-To-Noise Ratio Settings with** 201 **enrichR in Primary Human Hepatocytes**

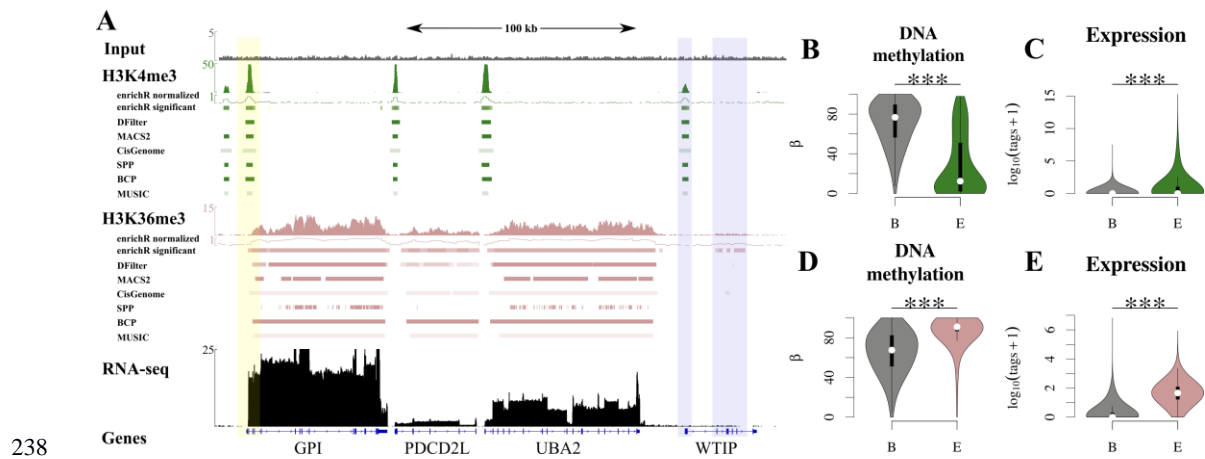
202 To illustrate the enrichment calling based on a robust background estimation, we applied  
203 enrichR to two ChIP-seq experiments against H3K4me3 and H3K36me3 in primary human  
204 hepatocytes. H3K4me3 correlates with promoter activity and DNA-hypomethylation [35-37]  
205 and exhibits a high signal-to-noise ratio (Supplemental Fig. 1). H3K36me3 represents a lower

206 signal-to-noise characteristics (Supplemental Fig. 1) and is associated to transcriptional  
207 elongation in the body of transcribed genes [38] as well as DNA-hypermethylation [39]. We  
208 performed enrichR analyses on the ChIP-seq data against Input-seq (Methods). The  
209 enrichment calls by enrichR were compared to peaks called by six popular peak calling tools  
210 ChIP-seq data: MACS2 [40], DFilter [21], CisGenome [22], SPP [23], BCP [41] and MUSIC  
211 [24].

212 As a first assessment, we inspected the coverage and enrichment/peak calls for H3K4me3 and  
213 H3K36me3 ChIP-seq in the vicinity of the Glucose-6-Phosphate Isomerase gene (GPI, Figure  
214 2A) — a housekeeping gene that is highly expressed in all cell types [42]. GPI was also  
215 expressed in primary human hepatocytes as measured by RNA-seq and showed a  
216 characteristic chromatin signature of transcription, *i.e.* H3K4me3 and H3K36me3 in the  
217 promoter and the gene body, respectively. All tested methods identified these characteristic  
218 enrichments at the GPI locus. Moreover, the promoter of the WTIP gene was detected as  
219 H3K4me3-enriched by all methods. Together with the measured shallow coverage of RNA-  
220 seq reads along its gene body this indicated that WTIP is expressed suggesting a genuine  
221 H3K36me3 enrichment in its gene body. Interestingly, this minute H3K36me3 enrichment  
222 was exclusively recovered by enrichR.

223 Genome-wide enrichR called H3K4me3-enrichment in 142,451 500 base-pair (bp) regions in  
224 primary human hepatocytes, corresponding to 45,522 consecutive regions representing ~3%  
225 of the mappable genome (71.2Mb). The identified regions were characterized by low levels of  
226 DNA methylation (Figure 2B), in line with the idea that H3K4me3 represses DNA  
227 methylation [35-37]. Furthermore, H3K4me3-enriched regions recovered by enrichR showed  
228 a higher density of CAGE-tags than the background (Figure 2C) indicating that they serve as  
229 active transcriptional start sites (TSSs) in this cell type. In fact, enrichR H3K4me3-enriched  
230 regions showed a statistically significant overlap with annotated TSSs (odds-ratio = 25.04,  
231 Fisher's signed exact test,  $P \leq 0.001$ , Supplemental Table 1). Together these observations  
232 support that enrichR identifies *bona fide* H3K4me3-enriched regions.

233 The comparison of enrichR enriched regions to MACS2, DFilter, CisGenome, SPP, BCP and  
234 MUSIC peaks revealed a substantial overlap at  $FDR = 0.1$  indicating that for H3K4me3 in  
235 this dataset all six methods work well, although in terms of covered bp DFilter (39.8Mb) and  
236 CisGenome (38.7Mb) called almost two-fold fewer regions than the other tools  
237 (mean=65.3Mb; Supplemental Note, Supplemental Fig. 2A, Supplemental Table 2).



239 **Figure 2. Enrichment Calling in Low and High Signal-To-Noise Ratio Settings with**  
240 **enrichR in Primary Human Hepatocytes.** (A) Input (grey), H3K4me3 (green, high signal-  
241 to-noise ratio), H3K36me3 (rose, low signal-to-noise ratio) and RNA-seq (black) barplots  
242 indicate coverage proximal to the human Glucose-6-Phosphate Isomerase (GPI, yellow  
243 overlay) locus on chromosome 19 in Primary Human Hepatocytes (PHH). Enrichment calls  
244 are indicated as colored boxes below respective tracks for enrichR, DFilter, MACS2,  
245 CisGenome's SeqPeak and SPP. The WTIP gene (blue overlay) had detectable H3K4me3  
246 enrichment at its promoter and minute H3K36me3 is recovered solely by enrichR. (B-C)  
247 enrichR H3K4me3-enriched regions were DNA-hypomethylated (B) and expressed as  
248 measured by CAGE (C). (D-E) enrichR H3K36me3-enriched regions were DNA-  
249 hypermethylated (D) and expressed as measured by RNA-seq (E).

250 For H3K36me3 enrichR identified 559,560 1 kilo base-pair (kb) windows as enriched,  
251 corresponding to 85,293 consecutive regions representing ~20% of the mappable genome  
252 (599.6Mb). H3K36me3-enriched regions recovered by enrichR showed high levels of DNA  
253 methylation (Figure 2D), in line with the observation that H3K36me3 recruits DNMT3B  
254 leading to *de novo* DNA methylation [39]. Furthermore, these regions showed significantly  
255 higher RNA-seq read coverage than background regions (Wilcoxon-signed-rank test  $P \leq$

256 0.001, Figure 2E), in line with the idea that H3K36me3 covers the gene body of transcribed  
257 genes [38]. Furthermore, enrichR H3K36me3-enriched regions showed a statistically  
258 significant overlap with annotated transcripts (odds-ratio = 17.06, Fisher's signed exact test,  
259  $P \leq 0.001$ , Supplemental Table 1). These results support that enrichR also identifies *bona*  
260 *fide* H3K36me3-enriched regions.

261 When compared to enrichR results, far less H3K36me3-enriched regions were reported by  
262 MACS2 (407.7Mb), BCP (396.5Mb), MUSIC (402.3Mb) and by especially DFilter (87.8Mb),  
263 SPP (25.1Mb) and CisGenome (36.4Mb), even when configured for detection in low signal-  
264 to-noise ratio settings (Methods). Almost all of these regions (MACS2: 399.1Mb; 97.9%,  
265 DFilter: 87.8Mb; 100%; CisGenome: 36.4Mb; 100%; SPP:24.2Mb; 96.7%; BCP:386.8Mb;  
266 97.6%; MUSIC:382.6Mb; 95.1%) were recovered by enrichR which leads to very few  
267 exclusive regions for the benchmark methods (Supplemental Fig. 2B). Regions called  
268 exclusively by enrichR (93.6Mb; 16.7%) were characterized by a median distance of >2kb to  
269 peaks recovered by other methods (Supplemental Fig. 2C). Furthermore, these regions  
270 showed significantly higher DNA-methylation levels and transcriptional activity than  
271 background regions suggesting once more a genuine H3K36me3 enrichment (Wilcoxon-  
272 signed-rank test  $P \leq 0.001$ , Supplemental Fig. 2D-E).

273 Next, we studied accuracy of H3K36me3-enrichment peak calls. Because there is no  
274 genome-wide ChIP-seq benchmark set on-hand, we defined a gold standard for each method  
275 based on a consensus vote among the six remaining tools [43] (Supplemental Note): At FDR  
276 0.1 DFilter and CisGenome achieved both highest precision (1.00), while enrichR had the  
277 highest recall (0.997) and BCP had the highest  $F_2$ -score (0.631; Supplemental Table 2).  
278 enrichR which called almost all regions of the five tools combined had a recall-weighted  $F_2$ -  
279 score of 0.533 compensating its menial precision (0.186) at q-value  $\leq 0.1$  with a superior  
280 recall. In fact, enrichR has the highest precision at recall  $\leq 0.9$  indicating that the consensus  
281 vote defined gold standard does not contain many enrichR-exclusive regions at q-value  $\leq$   
282 0.1 (Supplemental Fig. 3). In a second assessment, we studied the validity of tool-specific

283 regions, *i.e.* the peak calls not represented in the gold standard. To this end we defined a  
284 unified gold standard of H3K36me3-enrichment, *i.e.* the union of seven tool-specific gold  
285 standards, and seven sets of tool-specific regions (Supplemental Note). For all methods, the  
286 unified gold standard exhibited a significantly higher enrichment (fold change over Input)  
287 than tool-specific regions for enrichR, MACS2, SPP, BCP and MUSIC (Wilcoxon-signed-  
288 rank test;  $P \leq 0.01$ ; Supplemental Fig. 2F). Among these, enrichR had the most tool-specific  
289 regions (205,064; 36.6%) and showed significantly higher enrichment as well as read  
290 coverage than background regions (Supplemental Fig. 2G). Furthermore, enrichR-specific  
291 regions were remote from unified gold standard regions (median=14Mb; Supplemental Figure  
292 2H) and, yet, still overrepresented in annotated gene bodies (odds-ratio = 13; Supplemental  
293 Table 1).

294 Some ChIP-seq peak callers perform worse when the sequencing depth in the ChIP library is  
295 reduced [44]. To show the robustness of enrichR, we used the unified gold standard to  
296 benchmark all assessed tools on an *in silico* down sampled sequencing library (Supplemental  
297 Note). enrichR and MACS2 called >90% of the gold standard at 50% (30%) of the original  
298 H3K4me3 (H3K36me3) sequencing depth (Supplemental Fig. 4) suggesting that both  
299 methods are specific in even shallow sequenced ChIP libraries.

300 ChIP-seq coverage normalization based on *bona-fide* background regions is also done by  
301 NCIS [17]. For H3K36me3 NCIS estimated a normalization factor that was ~1.5-fold smaller  
302 than  $\theta^*$  and enrichR's  $\theta_B$  was ~2-fold smaller than  $\theta^*$  (Supplemental Fig. 5, Supplemental  
303 Table 3). Thus, enrichR achieved a normalization almost equivalent to NCIS, despite using a  
304 different model.

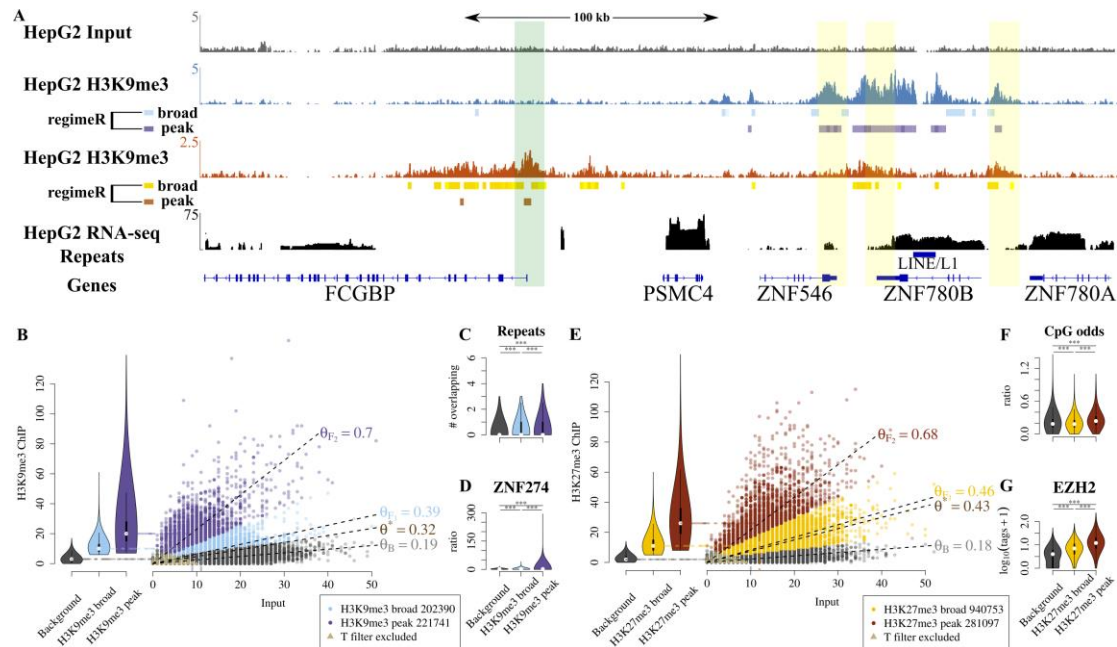
### 305 **Enrichment Regime Identification in H3K27me3 and H3K9me3 in HepG2** 306 **cells with regimeR.**

307 Hitherto discussed was the applicability of normR to a well-studied problem: the  
308 discrimination of enrichment against background. Here, we turn to a problem for which we

309 had found to best of our knowledge no precedent in the literature: the discrimination of  
310 moderate enrichment from high enrichment. We can easily address this problem by increasing  
311 the number of foreground components in normR from one single component to multiple  
312 components (Methods). We refer to this approach as regimeR: In the case of two foreground  
313 components, regimeR discriminates a *peak regime* (high enrichment) and a *broad regime*  
314 (moderate enrichment) over the background. We applied regimeR to H3K9me3 and  
315 H3K27me3 ChIP-seq data from the hepatocarcinoma cell line HepG2 over the control.

316 Figure 3A depicts a representative region on Human chromosome 19 harbouring active and  
317 repressed genes. regimeR segmented the ChIP-seq enrichment into broad and peak regions.  
318 For example, three H3K9me3 peaks flanked by moderate enrichment were detected by  
319 regimeR at the 3'-ends of ZNF546 and ZNF780A/B. Similarly, a H3K27me3-peak within a  
320 H3K27me3-broad domain was identified by regimeR at the “Fc Fragment Of IgG Binding  
321 Protein” gene promoter.

322 For H3K9me3, 14.7% of the HepG2 epigenome got classified into 202,390 broad (47.8%;  
323  $\mu_{\text{ChIP counts}} = 11.27$ ;  $\theta_{F_1} = 0.39$ ) and 221,741 peak regions (52.2%;  $\mu_{\text{ChIP counts}} = 23.75$ ;  
324  $\theta_{F_2} = 0.70$ ; Figure 3B). Both H3K9me3-broad and -peak regions showed a statistically  
325 significant overlap with repetitive DNA elements (Wilcoxon-signed-rank test;  $P \leq 0.001$ ;  
326 Figure 3C, Supplemental Fig. 6A), which is a reported feature of H3K9me3 marked  
327 constitutive heterochromatin [45]. Moreover, H3K9me3-peak regions showed significantly  
328 higher levels of ZNF274 than background and H3K9me3-broad regions (Wilcoxon-signed-  
329 rank test;  $P \leq 0.001$ , Figure 3D), in line with the idea that ZNF274 recruits the H3K9  
330 methyltransferase SETDB1 [46]. Thus H3K9me3-peak regions may coincide with nucleation  
331 sites for heterochromatin assembly at genomic repeat elements.



332

333 **Figure 3. H3K27me3 and H3K9me3 Enrichment Regime Identification in HepG2 cells**  
 334 **with regimeR.** (A) Input (grey), H3K9me3 (blue), H3K27me3 (orange) and RNA-seq  
 335 (black) coverage around a ZNF cluster on chromosome 19 in HepG2 cells. Individual  
 336 regimeR-computed regimes are displayed as boxes below respective tracks. The 5'-ends of  
 337 ZNF genes are marked with high H3K9me3 enrichment (yellow overlay) and the promoter of  
 338 FCGBP is marked by a H3K27me3 peak within a broad H3K27me3 domain (green overlay).  
 339 (B) regimeR identifies broad and peak H3K9me3 enrichment. (C-D) H3K9me3 peaks are  
 340 significantly enriched for repeats (C) and ZNF274 ChIP-seq reads (D) as compared to both  
 341 background and broad regions. (E) regimeR identifies broad and peak H3K27me3  
 342 enrichment. (F-G) H3K27me3 peaks have significantly greater CpG odds (F) and EZH2  
 343 binding (G) as compared to background and broad regions.

344 For H3K27me3, regimeR called 42.4% of the HepG2 epigenome H3K27me3-enriched  
 345 (1,221,850 1kb regions) and subdivided this into 940,753 broad (77%,  $\mu_{\text{ChIP counts}} = 12.03$ ;  
 346  $\theta_{F_1} = 0.46$ ) and 281,097 peak regions (23%,  $\mu_{\text{ChIP counts}} = 29.62$ ;  $\theta_{F_1} = 0.68$  Figure 3E).  
 347 H3K27me3 covered three times more of the genome than H3K9me3, yet, with a lower  
 348 fraction of peak regions than in H3K9me3. Moreover, the vast majority H3K9me3 and  
 349 H3K27me3 regimes were mutually exclusive in HepG2 cells (Supplemental Fig. 6B).  
 350 H3K27me3-peak regions were characterized by a higher CpG odds ratio (CpG-content  
 351 corrected for GC content) than both broad or background regions (Figure 3F, Supplemental  
 352 Fig. 6C). In conjunction with an elevated conservation (Supplemental Fig. 6D) and a

353 statistically significant overlap with annotated TSSs (Fisher's signed exact test;  $P \leq 0.001$ ;  
354 odds ratio = 1.98; Supplemental Table 4) this reaffirms that the TSSs targeted for peak  
355 H3K27me3 levels are high CpG promoters [47]. Similar to H3K9me3-peak regions,  
356 H3K27me3-peak regions were significantly enriched for the enzyme that catalyzes the  
357 modification, *i.e.* EZH2 [48-51] (Wilcoxon signed-rank test;  $P \leq 0.001$ , Figure 3G).  
358 Together these observations suggest that H3K27me3-broad and -peak regions show distinct  
359 characteristics with respect to CpG content, localization and EZH2 levels.

360 The observation that both H3K9me3- and H3K27me3-peak regions were associated with  
361 significantly higher levels of their catalyst than broad- and background regions indicates that  
362 they correspond to nucleation sites for heterochromatin assembly. In line with this  
363 observation we found that most H3K9me3-peak regions are either embedded in an H3K9me3  
364 broad domain (43.4%) or at the border of a broad domain (35.1%). The vast majority of  
365 H3K27me3-peak regions were embedded in an H3K27me3 broad domain (82.8%) where both  
366 regimes showed elevated conservation (Supplemental Fig. 6D). On the contrary, H3K9me3-  
367 peaks were less conserved than broad regions further supporting aforementioned idea that  
368 repetitive elements recruit the H3K9me3 methyltransferase.

## 369 **Difference Calling in Primary Human Hepatocytes and HepG2 cells with** 370 **diffR.**

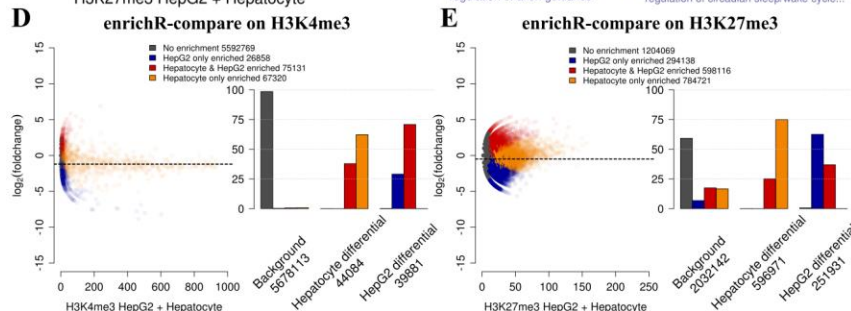
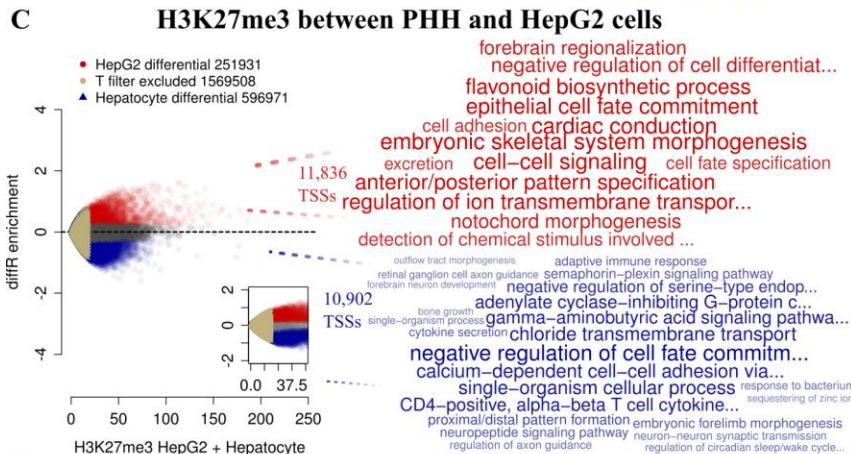
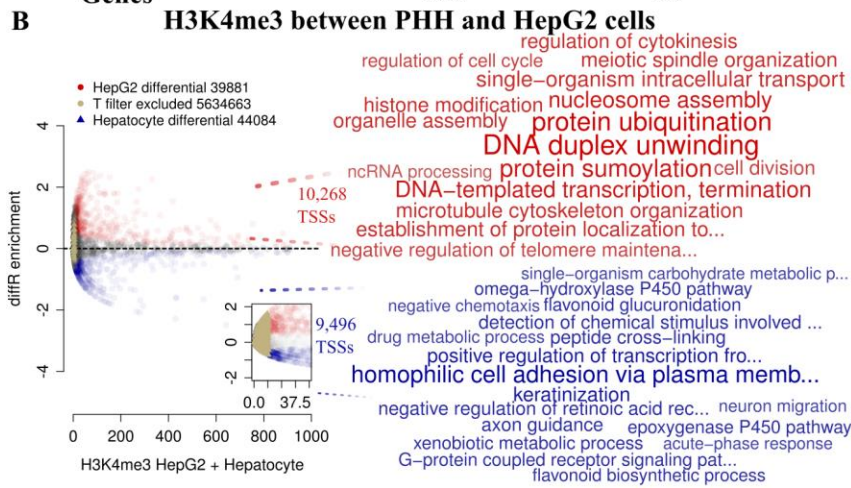
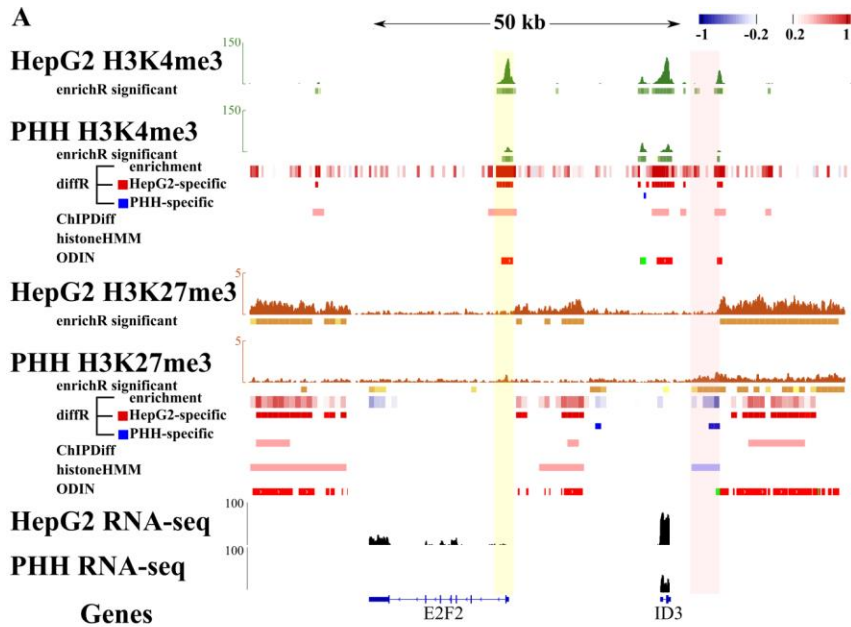
371 In addition to discriminating enrichment from background, another important task consists of  
372 identifying epigenetic alterations between conditions, e.g. healthy versus diseased or between  
373 cell-types. normR can address this problem by calling differential enrichment between ChIP-  
374 seq experiments from two conditions, referred to as "diffR". We applied diffR to H3K4me3  
375 and H3K27me3 ChIP-seq data from primary human hepatocytes (PHH) and the  
376 hepatocarcinoma cell line HepG2 (Methods). We compared the diffR results to those obtained  
377 by calling mutually exclusive enrichment with enrichR on the two conditions separately,  
378 referred to as "enrichR-compare". Additionally, we compared diffR results to three existing



379 tools, namely ChIPDiff [25], histoneHMM [27] and ODIN [26].

380 Visual inspection of a 50kb region on chromosome 19 confirmed that most  
381 H3K4me3/K27me3-enriched regions were common between HepG2 and PHH (Figure 4A).  
382 However, some enrichment was cell-type specific and was called by all methods, e.g. HepG2-  
383 specific H3K27me3-enrichment upstream of E2F2. However, differences in the histone  
384 modification level within mutually enriched regions were apparent, e.g. the increase in  
385 H3K4me3-enrichment at the E2F2 promoter in HepG2 could be identified by diffR, ChIPDiff  
386 and ODIN. E2F transcription factors are important regulators of the cell cycle [52-54]. E2F2  
387 is expressed in HepG2 but not in PHH suggesting that the induction of E2F2 might be linked  
388 to the much higher proliferative potential in HepG2 cells than in PHH. Further downstream of  
389 E2F2, enrichR identified a H3K27me3-differentially enriched domain accompanied by an  
390 emerging H3K4me3 peak in HepG2 cells. Thus, the induction of E2F2 in HepG2 may be  
391 explained by the opening of an enhancer at this region supported by reported binding of RNA  
392 polymerase 2 and CTCF in HepG2 cells [2].

393 For H3K4me3, diffR recovered 59,288 500bp regions (14Mb) as being differentially enriched  
394 between HepG2 and PHH (Figure 4B). Of these, 27,913 regions had a higher enrichment in  
395 HepG2 which overlapped 10,268 TSSs driving genes mainly related to the DNA replication  
396 and cell division. 31,375 PHH-specific H3K4me3 regions upregulated 9,496 TSSs of genes  
397 associated with liver function (P450 pathway) and tissue characteristics (keratinization, cell  
398 adhesion) absent in the HepG2 cell line. For H3K27me3, diffR reported 800,073 1kb regions  
399 (800Mb) as differentially H3K27me3-enriched (Figure 4C). Out of these 215,466 revealed  
400 HepG2-specific repression at 11,836 TSSs of genes regulating morphogenesis and cell-cell  
401 signaling. On the other hand, the 584,607 PHH-specific regions repressed 10,902 TSSs of  
402 genes functioning in cell fate commitment and immune response. Taken together, diffR  
403 uncovered functional differences related to immortalization of liver cells solely based on two  
404 ChIP-seq experiments.



406 **Figure 4. Difference Calling on H3K4me3 and H3K27me3 in Primary Human**  
407 **Hepatocytes (PHH) and HepG2 cells with diffR.** (A) Input (grey), H3K4me3 (green),  
408 H3K27me3 (orange) and RNA-seq (black) coverage around E2F Transcription Factor 2  
409 promoter (E2F2, yellow overlay) locus in Primary Human Hepatocytes (PHH) and HepG2  
410 cells. A region ~40kb upstream of the E2F2 promoter shows significant differential  
411 enrichment for H3K4me3 and H3K27me3 (pink overlay). enrichR-computed enriched regions  
412 displayed as boxes below to respective. Differentially enriched regions are displayed as red  
413 (HepG2 conditional) or blue (PHH conditional) boxes for diffR, ChIPDiff, histoneHMM and  
414 ODIN. (B,C) diffR recovers conditional differences in H3K4me3 (B) and H3K27me3 (C)  
415 enrichments that cover transcriptional start sites (TSSs) driving genes functioning in cell  
416 metabolism and development (wordclouds right panel). (D,E) enrichR-compare identifies  
417 H3K4me3 (D) and H3K27me3 (E) mutually exclusive enrichment between PHH and HepG2  
418 cells, but can not detect differences in histone modification level. (Right panels) diffR regions  
419 fall into enrichR-compare called regions of mutually exclusive enrichment but also resolve  
420 significant differences in ChIP-seq signal not detected by enrichR-compare.

421 Another normR approach can detect conditional differences by calling individual ChIP-seq  
422 enrichment over control for each condition and then identify mutually exclusive enrichment,  
423 referred to as “enrichR-compare”. We used this approach to benchmark results obtained from  
424 diffR. Genome-wide H3K4me3 enrichR-compare analysis revealed that most enriched 500bp  
425 regions were common in HepG2 and PHH (101,989, Figure 4D), while 26,858 were HepG2-  
426 and 67,320 PHH-specific. As expected, the comparison to enrichR-compare also revealed that  
427 by a majority diffR difference calls were either mutual exclusive enrichment or changes in the  
428 level of enrichment (Figure 4D, Supplemental Table 5). For H3K27me3, enrichR-compare  
429 revealed that most H3K27me3-enriched regions were common in HepG2 and PHH (892,254,  
430 Figure 4E), while 294,138 were HepG2- and 784,721 were PHH-specific. Again, diffR was  
431 very specific in capturing both mutual exclusive enrichment and changes in the level of  
432 enrichment (Figure 4E), However, we observed a discrepancy in sensitivity: 58.6% (44%) of  
433 the H3K4me3 (H3K27me3) mutually exclusive regions were not called by diffR leading to  
434 contradictory results (Supplemental Fig. 7A,B; Supplemental Table 5). Interestingly, most of  
435 the discrepancies were attributed to a more strict P-value filter to eliminate low power (i.e.  
436 low count) regions in the two-sided binomial test in diffR (Methods). By applying the diffR

437 P-value filter to enrichR-compare, results became substantially more concordant , e.g. 2.99%  
438 (319) false negatives for H3K4me3 in HepG2 cells (Supplemental Fig. 7 C,D, Supplemental  
439 Table 5).

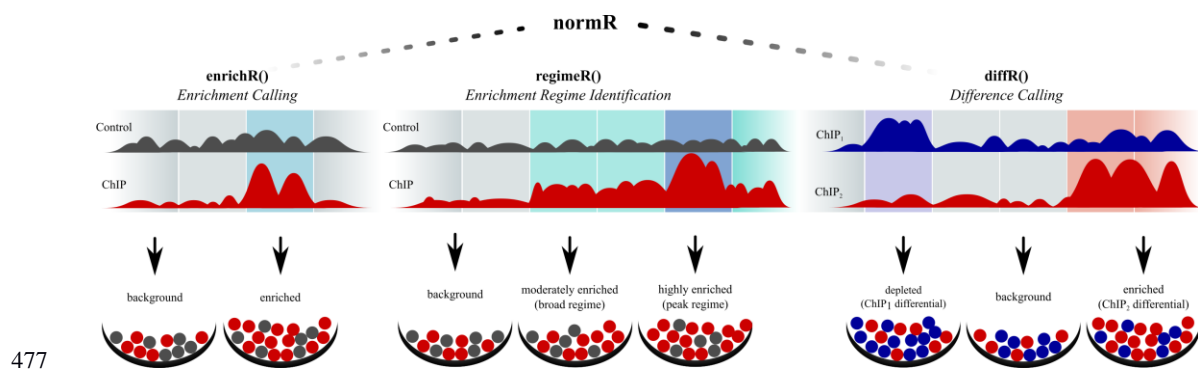
440 In addition, some discrepancies between diffR and enrichR-compare may be attributed to  
441 Copy Number Variations (CNVs) in HepG2 cells which are prevalent in immortalized cell  
442 types [55,56]. To alleviate this problem we ran diffR on HepG2 and PHH Input tracks with 20  
443 and 50kb windows (Supplemental Fig. 8). Assuming that there are no CNVs in the PHH data,  
444 diffR recovered 91% of 6,487 windows (odds-ratio=112.7) which overlap 80 annotated  
445 amplifications in HepG2 [2] (13% of genome; median(*length*)=163kb). Nevertheless, diffR  
446 failed to detect 88% of 249 windows (odds-ratio=40.8) that overlap 170 annotated very short  
447 heterozygous and homozygous deletions (6% of genome; median(*length*)=9kb). Despite this,  
448 the discrepancies between enrichR-compare and diffR were partially removed when filtering  
449 results for diffR called CNVs (Supplemental Fig. 7E,F, Supplemental Table 5) to a similar  
450 extend than filtering for experimentally validated CNVs (Supplemental Fig.7G,H,  
451 Supplemental Table 5).

452 Next, we compared genome-wide diffR results to those obtained from ChIPDiff,  
453 histoneHMM and ODIN. To this end we once more defined a gold standard based among a  
454 consensus vote among the tools (Supplemental Note): ChIPDiff was most precise ( $\mu_{\text{Precision}} =$   
455 0.70) and diffR had the highest recall ( $\mu_{\text{Recall}} = 0.80$ ) together with the best F1-scores  
456 ( $\mu_{\text{F1-score}} = 0.50$ ; Supplemental Table 6). A unified gold standard of all tool-specific gold  
457 standards revealed that most tool-specific regions were called by diffR (28.9Mb) and ODIN  
458 (25.4Mb) for H3K4me3 and by ODIN (701.7Mb) and histoneHMM (689.1Mb) for  
459 H3K27me3 (Supplemental Table 7). Turning to absolute fold changes, the unified gold  
460 standard showed highest levels together with diffR, ChIPDiff and histoneHMM  
461 (Supplemental Fig. 9A,B). In terms of read coverage, diffR- and ODIN-specific regions had  
462 highest counts (Supplemental Fig. 9C,D). In conclusion, diffR identified conditional  
463 differences for, both, H3K4me3 and H3K27me3 which were supported by a good classifier

464 performance, a high absolute fold change as well as an inference-adequate read coverage  
465 eliminating low power regions.

## 466 Discussion

467 In summary, we present an extendable methodology called “normR” that enables the  
468 extensive analysis of ChIP-seq data in epigenetic studies (Fig. 5). By modeling foreground  
469 and background jointly, normalization and enrichment calling are performed simultaneously.  
470 The implicit modeling of the effect of enrichment on the overall read statistics increases the  
471 sensitivity in detecting shallow differences in ChIP enrichment even in low signal-to-noise  
472 ratio data. Furthermore, we demonstrated the suitability of the normR approach for the  
473 identification of distinct epigenetic enrichment regimes in hepatocarcinoma cells and the  
474 quantification of conditional epigenetic differences between hepatocarcinoma cells and their  
475 tissue-of-origin. We envision how normR enrichment calling augments today’s epigenetic  
476 analyses ranging from clustering [30] to visualization [31].



478 **Figure 5. The normR Approach: A Robust and Broadly Applicable Methodology for**  
479 **Normalization and Difference Calling in ChIP-seq Data.** The translational normR

480 methodology allows for the calling of ChIP enrichment over a user-specified control, the  
481 identification of distinct ChIP enrichment regimes and the quantification of differences in  
482 ChIP signal level between two conditions.

483 Firstly we used normR to call enrichment in high (H3K4me3) and low (H3K36me3) signal-  
484 to-noise ratio ChIP-seq data, referred to as “enrichR”. Auxiliary information such as DNA  
485 methylation and expression supported the enrichR-based classification. Given the difficulty

486 inherent in the ill-defined problem represented by ChIP-seq analysis we introduce a novel  
487 binary classifier statistic that defines a gold standard based on a consensus vote among seven  
488 published ChIP-seq peak callers. Our findings indicated that enrichR performs equally well as  
489 previously described approaches in ChIP-seq tracks with high signal-to-noise-ratio such as  
490 H3K4me3. Furthermore, enrichR outperformed existing tools in the detection of low levels of  
491 genuine enrichment in low signal-to-noise ratio data such as H3K36me3. We attribute the  
492 superior performance in the latter scenario to our sensitive normalization technique which  
493 accounts not only for varying sequencing depth but specifically addresses the effect of ChIP  
494 enrichment on the overall read statistics. The sensitive enrichR approach is an asset in future  
495 studies on epigenetic signatures and segmentations.

496 Secondly normR was used to facilitate the discrimination of peak- and broad-regions against  
497 background in a single analysis, referred to as “regimeR”. The analysis of H3K9me3 and  
498 H3K27me3 in HepG2 cells revealed that there exist distinct characteristics of peak- and broad  
499 regions in these heterochromatic marks. Specifically, H3K9me3 peaks were enriched for  
500 ZNF274 at repetitive elements. High enrichment of H3K9me3 at these sites can be explained  
501 by the recruitment of the H3K9 methyltransferase SETDB1 by ZNF274 [46]. H3K27me3  
502 peaks were found within broad H3K27me3 domains at conserved CpG-dense regions bound  
503 by EZH2, supporting the idea of CpG-enriched polycomb recruitment sites [57]. Taken  
504 together, our regimeR-based study suggests that H3K9me3 and H3K27me3 peaks correspond  
505 to nucleation sites for heterochromatin assembly. In the future, regimeR will prove useful in  
506 studies of heterogeneity in cellular epigenetic markings to identify regions of promiscuous  
507 protein binding.

508 Finally we presented normR for the direct comparison of two ChIP-seq experiments, referred  
509 to as “diffR”. Our diffR-based comparison of H3K4me3 and heterochromatic H3K27me3  
510 between HepG2 cells and PHH revealed conditional differences associated to cell function  
511 and immortalization, e.g. a potential E2F2 enhancer region made accessible in HepG2 cells.  
512 Interestingly, H3K27me3 covered a smaller fraction of the HepG2 genome as compared to

513 PHH. Using a statistic of mutually exclusive enrichment by enrichR-compare and consensus  
514 votes among previously developed difference callers, we showed that diffR performs  
515 outstandingly in the detection of conditional differences in ChIP-seq data. Furthermore, we  
516 could show that diffR's accuracy can be increased by incorporation of CNV information, as  
517 measured experimentally or by using diffR on two Input experiments. In the future, a more  
518 principled approach of the joint modelling of conditional ChIP-seq tracks together with their  
519 control is desirable.

520 Taken together normR proved as a versatile and sensitive toolbox for the discrimination of  
521 enrichment against background ("enrichR"), the unprecedented detection of enrichment  
522 regimes such as peaks and broad enrichment ("regimeR") and the direct quantification of  
523 differences between two conditions ("diffR"). We anticipate that normR will be applied to all  
524 enrichment based sequencing technologies like MeDIP-seq and HiC. In fact, a derivate of  
525 normR has recently been used to identify co-localizing histone modifications in a novel  
526 reChIP-seq data set [43] where the background estimation is complicated by the presence of  
527 enrichment in the control experiment. In the future, an automated determination of the  
528 number of enrichment components in the normR model will be adjuvant in studying  
529 epigenomic heterogeneity in conjunction with recently reported single cell ChIP-seq data [58].

## 530 **Methods**

### 531 **The normR Methods**

532 Given two vectors of integers  $r$  (control) and  $s$  (treatment) of identical length  $n$ , we model the  
533 read counts from the ChIP and control by a binomial  $m$ -mixture model:

$$k_i \sim \text{Categorical}(\pi)$$

$$N_i = s_i + r_i \mid k_i = j \sim \text{Binomial}(N_j, \theta_j)$$

534 with  $i = 1, \dots, n$  and  $\sum \pi_j = 1$ ;  $\pi_j \in [0,1]$ ;  $j = 1, \dots, m$ . Given this model, normR follows a  
535 two step procedure: (i) The mixture model is fit by expectation maximization (EM; [28])

536 using the likelihood function,

$$\mathcal{L} = P(s_i, r_i | \pi, \theta, N_i) = \prod_{i=1}^n \binom{N_i}{s_i} \sum_{j=1}^m \pi_j \times \theta_j^{s_i} \times (1 - \theta_j)^{r_i};$$

537 and (ii) each entry  $(r_i, s_i)$  is tested for significance against a fitted background to component  
538 to label enriched regions.

539 In a preprocessing stage, the vectors  $r$  and  $s$  are filtered for entries where  $r = s = 0$  because  
540 no assertion about their enrichment state can be made. Secondly, a map of unique  $(r, s)$  tuples  
541 is created to reduce the number of computations needed which improves runtime substantially.

542 In the first mode fitting step, the EM is initialized with  $\pi$  sampled from  $U(0,1)$  and  $\theta$  sampled  
543 from  $U(0.001, \theta^*)$ . Upon convergence with  $\varepsilon \leq 0.001$ , an enrichment factor (average fold  
544 enrichment)  $f_j = \frac{\theta_j}{1-\theta_j} * \frac{1-\theta_B}{\theta_B}$  is computed for each mixture component  $j \neq B$ , where  $\theta_B$  the

545 smallest of  $\{\theta_1, \dots, \theta_m\}$  (the closest to  $\theta^*$ ) in the case of enrichment (difference) calling. The  
546 EM is run 10 times per default to find the fit with greatest  $\mathcal{L}$ . In the second step, every  $(r_i, s_i)$   
547 is tested for significance against the background component. Resulting P-values are filtered  
548 using the T method [32] (P-value threshold 0.0001 per default) to take into account the  
549 discreteness of P-values for a correct estimation of the proportion of true null hypotheses. T-  
550 filtered P-values are transformed to q-values for FDR correction [33]. Additionally, a

551 normalized enrichment  $e_i$  is calculated for every entry  $(r, s)_i$  with  $e_i = \frac{\ln\left(\frac{s_i + \alpha_s}{r_i + \alpha_r} \times \frac{\alpha_r}{\alpha_s}\right)}{\ln(f_j)}$  where

552  $\alpha_r = \frac{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i} \times r_i}{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i}}$  and  $\alpha_s = \frac{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i} \times s_i}{\sum_i \theta_B^{s_i} \times (1-\theta_B)^{r_i}}$  represent a model specific pseudo count for

553 control and treatment, respectively. The normalized enrichment can be written to bigWig or  
554 bedGraph format for convenient display in a genome browser of choice, e.g. UCSC genome  
555 browser [59] or Integrative Genomics Viewer [60].

556 In the case of enrichment calling two components (background, enrichment) are fit with the  
557 enrichR subroutine of the normR package. Herein, the background model  $\theta_B$  is set to the  
558 mixture component with smallest  $\theta$ . For difference calling, three components (background,



559 control enriched, treatment enriched) are fit with the diffR subroutine for  $r$  (condition 1) and  $s$   
560 (condition 2) counts. The background model is set to  $\theta_B$  closest to  $\theta^*$ . The diffR T method  
561 uses the maximal threshold estimated from P-values for  $\theta_B$  fit for either  $(r, s)$  or the label-  
562 switched  $(s, r)$ . For regime calling, the regimeR subroutine fits an arbitrary number of  
563 components representing background plus a fixed number of enrichment regimes. Identically  
564 to enrichment calling, the background model is set to the mixture component with smallest  $\theta$ .  
565 In a second step, every significantly enriched bin passing the P-value filter (see above) is  
566 assigned to an enrichment regime by *Maximum A Posteriori*.

567 Note that by nature the binomial mixture model assumes the independence between regions  
568 which is valid for a sufficiently large bin size (*i.e.* fragment length). Consequently, the usage  
569 of a binomial mixture model improves computational runtime. The normR algorithm is  
570 implemented in C++ and R. A ready-to-use R-package can be obtained from  
571 <http://bioconductor.org/packages/normr> where also a tutorial on use cases can be found.

## 572 **ChIP-seq in primary human hepatocytes and HepG2 cells**

573 HepG2 cells and human hepatocytes, obtained from donors after written consent by tissue  
574 resection and perfusion [61], have been fixed in for 5 minutes in 1% formaldehyde.  
575 Formaldehyde has been quenched using 125 mM glycine and cells have been washed in PBS,  
576 pelleted and snap-frozen in liquid nitrogen. Five (human hepatocytes) to ten (HepG2) million  
577 cells have been processed for chromatin preparation, using the NEXSON protocol, as  
578 previously described [62]. After chromatin sonication, samples have been quality controlled  
579 to check chromatin recovery and fragment size distribution as previously described.

580 Prior ChIP, chromatin has been diluted 1:2 in the ChIP buffer H from the Diagenode Auto  
581 histone ChIP-seq kit (C01010022), supplemented with protease inhibitor cocktail. Chromatin  
582 from 100,000 to 500,000 cells has been incubated with one microgram of the following  
583 antibodies: H3K4me3 (C15410003), H3K36me3 (C15410192), H3K9me3 (C15410193),  
584 H3K27me3 (C15410195), all from Diagenode. ChIP has been performed using the automated

585 platform SX-8G IP-Star (Diagenode), with the following parameters: “indirect ChIP”, 200 µl  
586 ChIP volume, 14 hours of antibody incubation, 4 hours of beads incubation, and 5 minutes  
587 beads washes. After the DNA elution from the beads, samples were collected, RNaseA-  
588 treated, de-proteinized and decrosslinked overnight at 65 °C. Input samples have been  
589 prepared by taking 1% of the starting chromatin before ChIP and by decrosslinking it together  
590 with the ChIP samples. DNA has been manually purified using the Qiagen minElute columns.  
591 Libraries from 2 to 10 ng of purified DNA have been prepared using the NEBNext Ultra  
592 DNA library preparation kit (NEB, E7370S) following manufacturer’s instruction and  
593 skipping the size selection. Libraries have been sequenced paired-end, with a read length of  
594 50 bp, on an Illumina HiSeq 2500 (version 3 chemistry).

### 595 **RNA-seq in primary human hepatocytes and HepG2 cells**

596 Trizol extration was used for preparation of Total RNA according to the manufacturer’s  
597 guidelines and as described in [63]. An Agilent Bioanalyzer (Agilent, Santa Clara, USA) was  
598 used to check RNA integrity following the manufacturer’s guidelines.

599 Strand-specific sequencing libraries for mRNA and total-RNA were constructed for the  
600 HepG2 cells and human hepatocytes using the TruSeq stranded Total RNA kit (Illumina Inc,  
601 San Diego, USA) starting from 500 ng of the total RNA of the samples. Illumina HiSeq2000  
602 was used to perform the sequencing (101-nucleotide paired-end reads for each library)  
603 resulting in the creation of about 100 million reads per library.

604 The reads were aligned to the NCBI 37.1 version of human genome using TopHat v2.0.11  
605 [64] in the settings “--library-type fr-firststrand” and “--b2-very-sensitive”.

606 Reads mapping to genes were counted using htseq-count from HTSeq-0.6.1p1 [65] in '-f bam  
607 -s reverse -m union -a 20' setting. Annotation file for running htseq-count was downloaded  
608 from GENCODE release 19 (GRCh37.p13).

## 609 Quantification of reads

610 Paired-end reads from Input, H3K4me3, H3K27me3, H3K36me3 and H3K9me3 ChIP-seq for  
611 primary human hepatocytes and HepG2 cells were mapped with bwa (version 0.6.2) against  
612 hg19. Fragment coverage tracks for browser display were generated with deepTools [66] in  
613 25 bp windows (-bs 25) considering only first reads in a properly mapped pair (--samFlag 66)  
614 with a mapping quality of at least 20 (--MinMappingQuality 20) and normalized to the  
615 effective genome size (--normalizeTo1x 2451960000):

```
616 bamCoverage -bam in.bam -o out.bw -of bigwig -bs 25 \  
617 --samFlag 66 --minMappingQuality 20 --normalizeTo1x 2451960000
```

618 For enrichment and peak calling, only regions on regular autosomes (chr1-chr22; 2.9Gb) were  
619 used:

```
620 require(GenomeInfoDb)  
621 genome <- fetchExtendedChromInfoFromUCSC("hg19")  
622 genome <- genome[which(!genome$circular &  
623 genome$SequenceRole=="assembled-molecule"), 1:2]  
624 genome <- genome[grep("X|Y|M", genome[, 1], invert=T), ]  
625  
626 require(GenomicRanges)  
627 genome.gr <- GRanges(  
628 seqnames = genome[, 1],  
629 ranges = IRanges(start = 1, end = genome[, 2]),  
630 seqinfo = Seqinfo(  
631 seqnames = genome[,1],  
632 seqlengths = genome[,2],  
633 genome = "hg19"))  
634 )  
635 }
```

636 For paired end data, we considered only reads with a mapping quality of at least 20  
637 (mapqual=20). We regarded midpoints of properly mapped fragments (midpoint = TRUE)  
638 that were non-duplicated (filteredFlag=1024) and within 100 to 220 bp in length  
639 (tlenFilter=c(100,220)) in 500 (1,000) bp windows for H3K4me3 (H3K27me3/K36me3/  
640 K9me3) with normR's countConfigPairedEnd function:

```
641 require(normr)
```

```
642 countConfig <- countConfigPairedEnd(  
643   binsize = 500, #1000  
644   mapqual = 20,  
645   midpoint = TRUE,  
646   filteredFlag = 1024,  
647   tlenFilter = c(100,220)  
648   shift = 0  
649 )
```

650 HepG2 CAGE data was downloaded from GSM849335 [67]. Primary human hepatocyte  
651 CAGE data was downloaded from CAGE  
652 [http://fantom.gsc.riken.jp/5/datafiles/latest/basic/human.primary\\_cell.hCAGE/Hepatocyte%2](http://fantom.gsc.riken.jp/5/datafiles/latest/basic/human.primary_cell.hCAGE/Hepatocyte%252c%2520donor2.CNhs12349.11603-120I1.hg19.nobarcodes.bam)  
653 [52c%2520donor2.CNhs12349.11603-120I1.hg19.nobarcodes.bam](http://fantom.gsc.riken.jp/5/datafiles/latest/basic/human.primary_cell.hCAGE/Hepatocyte%252c%2520donor2.CNhs12349.11603-120I1.hg19.nobarcodes.bam) (Fantom5 [68]) Reads with  
654 mapping quality of at least 20 were counted with bamsignals  
655 (<http://bioconductor.org/packages/bamsignals>):

```
656 require(bamsignals)  
657 cage <- bamProfile(  
658   bampath = "Cage.bam",  
659   gr = genome.gr,  
660   binsize = 500, #1000  
661   mapqual = 20  
662 )
```

663 EZH2 ChIP-seq alignments (GSM1003576) and the respective control alignment  
664 (GSM733780) were downloaded from the UCSC encode repository ([2]  
665 [hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeBroadHistone/](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeBroadHistone/)). For these  
666 single end data, we shifted reads by 100 bp in 3' direction (shift=100) and counted in 500  
667 (1,000) bp bins:

```
668 countConfig <- countConfigSingleEnd(  
669   binsize = 500, #1000  
670   mapqual = 20,  
671   filteredFlag = 1024,  
672   shift = 100  
673 )
```

## 674 **Enrichment calling with enrichR**

675 Read counts in H3K4me3 and H3K36me3 were modeled with 2 components in enrichR and

676 the fitted background components were used for significance tests. Bins with q-value  $\leq 0.05$   
677 (H3K4me3) and q-value  $\leq 0.1$  (H3K27me3/K36me3/K9me3) were called enriched and  
678 exported to bed tracks for display:

```
679 enrichment <- enrichR(  
680   treatment = "ChIP.bam",  
681   control = "Input.bam",  
682   genome = genome,  
683   countConfig = countConfig,  
684   procs = 24  
685 )  
686 exportR(  
687   x = enrichment,  
688   filename = "enriched.bed",  
689   type = "bed",  
690   fdr = 0.05 #0.1  
691 )
```

## 692 **DNA-methylation in primary human hepatocytes and HepG2 cells**

693 For whole-genome bisulfite sequencing we produced two types of NGS libraries to achieve  
694 even read coverage. Firstly, we used 100ng of DNA with the TruSeq DNA methylation kit  
695 (Illumina, San Diego, USA) according to the manufacturer's protocol. The second type was  
696 done as previously described [43]. Briefly, 2  $\mu$ g of DNA were sheared using a Bioruptor NGS  
697 device (Diagenode, Liege, Belgium) and cleaned-up using Ampure beads XP (Beckman  
698 Coulter, Brea, USA). Then samples were subjected to end-repair, A-tailing and adaptor  
699 ligation steps using components of the TruSeq DNA PCR-Free Library Preparation Kit  
700 (Illumina). After bisulfite conversion involving the Zymo Gold kit (Zymo, Irvine, USA) the  
701 libraries were PCR amplified for 10-12 cycles. The amplified libraries were purified using  
702 Ampure beads XP and sequenced on three lanes of V3 paired-end flow cells (2x 100bp).  
703 Reads were mapped using BWA [69] and methylation levels were called with Bis-SNP37 [70].

704 Beta values were calculated for each bin and weighted by coverage and number of CpGs  $M$  in

705 that region:  $\beta = \frac{\sum_{j=0}^M \text{ReadCount}_j * \text{FractionMethylated}_j}{\sum_{j=0}^M \text{ReadCount}_j}$ . Only regions with at least 2 CpGs

706 covered by reads were reported.

## 707 **Transcription Start Site Definition**

708 54,763 promoters (extend 750bp down- and upstream of TSS) of 54,849 GENCODE genes

709 [71] obtained by using GenomicFeatures R package [72]:

```
710 require(GenomicFeatures)
711 gencode <- loadDb("data/gencode.v19.annotation.transcriptDb.sqlite")
712 genes <- genes(gencode)
713 proms <- unique(promoters(genes, upstream=750, downstream=750))
```

## 714 **MACS, DFilter, CisGenome, SPP, BCP and MUSIC Peak Calling**

715 Peaks were called with MACS2 [40] (v2.1.0.20150731), DFilter [21] (v1.6), CisGenome [22],  
716 SPP [23], BCP [41] (v1.1) and MUSIC [24]. A FDR threshold of 0.1 was used. To compare  
717 called peaks by above methods to enrichR called regions, overlap of peaks with 500 bp (1,000  
718 bp) windows was calculated for H3K4me3 (H3K36me3). See Supplemental Note for details.

## 719 **Normalization Factor Comparison with NCIS**

720 NCIS [17] was run in R to calculate the normalization factor for comparison to enrichR's  
721 normalization factor:

```
722 require(NCIS)
723 ncis <- NCIS(
724   chip.data = "ChIP.bed",
725   input.data = "Control.bed",
726   data.type = "BED",
727   chr.vec = genome[,1],
728   chr.len.vec = genome[,2]
729 )
730 ncis.norm <- ncis$est
```

## 731 **Regime calling with regimeR**

732 Read counts in H3K27me3 and H3K9me3 in HepG2 cells were modeled in regimeR with 3  
733 components (background, moderate enrichment, high enrichment). Bins with FDR q-value  $\leq$   
734 0.1 were called enriched and assigned to an enrichment component by *Maximum A Posteriori*

735 and exported to bed using normR's exportR function:

```
736 regimes <- regimeR(  
737   treatment = "ChIP.bam",  
738   control = "Input.bam",  
739   genome = genome,  
740   models = 3,  
741   countConfig = countConfig,  
742   procs = 24  
743 )  
744 exportR(  
745   x = regimes,  
746   filename = "regimes.bed",  
747   type = "bed",  
748   fdr = 0.05 #0.1  
749 )
```

## 750 **Difference calling with diffR**

751 Read counts in H3K4me3 and H3K27me3 in primary human hepatocytes (control) and  
752 HepG2 cells (treatment) were modeled in diffR with 3 components (background/no difference,  
753 treatment-enriched, control-enriched) and the mixture component with  $\theta_j$  closest to  $\theta^*$  was  
754 used as background for a two-sided significance test. Bins with q-value  $\leq 0.05$  (0.1) for  
755 H3K4me3 (H3K27me3) were called differentially enriched and assigned to treatment or  
756 control by *Maximum A Posteriori*.

```
757 diffs <- diffR(  
758   treatment = "ChIP1.bam",  
759   control = "ChIP2.bam",  
760   genome = genome,  
761   countConfig = countConfig,  
762   procs = 24  
763 )  
764 exportR(  
765   x = diffs,  
766   filename = "differences.bed",  
767   type = "bed",  
768   fdr = 0.05 #0.1  
769 )
```

770 To analyze differentially enriched regions for precision and recall, *mutually exclusive*

771 *enrichment* in control (treatment) was obtained by considering `enrichR()` calls present only in  
772 control (treatment) with respect to treatment (control). For a fair comparison, only significant  
773 regions with a posterior of  $\geq 0.50$  were considered.

## 774 **Gene Ontology Analysis**

775 We used `topGO` [73] on gene ontology “Biological Process” (BP) with algorithms “classic”  
776 (algorithm=”classic”) and “elim” (algorithm=”elim”) for statistics “fisher” (statistic=”fisher”)  
777 and “ks” (statistic=”ks”) for GENCODE gene IDs mapped to Ensembl gene IDs. The “ks”  
778 statistic allows for supplying a score for each entity. We used the `diffR` calculated q-value as  
779 score. We retained only top 1,000 (n=1000) GO terms ordered by “elim” algorithm and  
780 ranked by “classic” algorithm calculated P-values:

```
781 require(topGO)
782
783 #get GO annotated Ensembl Genes
784 go2ensembl <- annFUN.org(ontology, mapping="org.Hs.eg.db", ID="ensembl")
785
786 #get GENCODE genes and filter these for the ones in gene universe
787 gencode <- loadDb("data/gencode.v19.annotation.transcriptDb.sqlite")
788 gene.universe <- intersect(
789   unique(GenomicFeatures::genes(gencode)$genes),
790   unique(unlist(go2ensembl))
791 )
792
793 #set diffR pvalue as score for differentially modified TSSs
794 idx <- gene.universe %in% diffTSSs
795 allGenes <- 1-as.integer(idx)
796 names(allGenes) <- gene.universe
797 allGenes[idx] <- pvals[diffTSSs %in% gene.universe]
798 goData <- new("topGOdata",
799   description="diffR differential TSS histone marking study (scored)",
800   ontology="BP",
801   allGenes=allGenes, geneSel=function(p) { return(p <= 0.05) },
802   annot=annFUN.GO2genes, GO2genes=go2ensembl, #GO mapping for ensembl IDs
803   nodeSize=10
804 )
805
806 #testing
```



```
807 resultFisher <- runTest(goData, algorithm="classic", statistic="fisher")
808 resultKS <- runTest(goData, algorithm="classic", statistic="ks")
809 resultKS.elim <- runTest(goData, algorithm="elim", statistic="ks")
810
811 #compile results
812 resDf <- GenTable(goData,
813   classicFisher = resultFisher,
814   classicKS = resultKS,
815   elimKS = resultKS.elim,
816   orderBy = "elimKS",
817   ranksOf = "classicFisher",
818   topNodes=1000
819 )
```

## 820 **ChIPDiff, histoneHMM and ODIN Difference**

821 Differences for H3K4me3 (H3K27me3) between Hepatocytes and HepG2 cells were called  
822 with ChIPDiff [25], histoneHMM (v1.6) [27] and ODIN (v0.4) [26]. A FDR threshold of 0.1  
823 was used. To compare called peaks by above methods to diffR called regions, overlap of  
824 peaks with 500 bp (1,000 bp) windows was calculated for H3K4me3 (H3K27me3). See  
825 Supplemental Note for details.

## 826 **HepG2 Genotyping**

827 HepG2 genotype information for hg19 was generated by ENCODE/HudsonAlpha  
828 (GSM999286) and downloaded from UCSC ([http://hgdownload.cse.ucsc.edu/goldenPath](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeHaibGenotype/wgEncodeHaibGenotypeHepg2RegionsRep1.be)  
829 [/hg19/encodeDCC/wgEncodeHaibGenotype/wgEncodeHaibGenotypeHepg2RegionsRep1.be](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeHaibGenotype/wgEncodeHaibGenotypeHepg2RegionsRep1.be)  
830 [dLogR.gz](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeHaibGenotype/wgEncodeHaibGenotypeHepg2RegionsRep1.be)).

## 831 **Data Access**

832 H3K4me3, H3K9me3, H3K27me3, H3K36me3 ChIP-seq and Input data for primary  
833 human hepatocytes have been deposited at the “European Genome-Phenome  
834 Archive” under the accession EGAS00001002080. H3K4me3, H3K9me3, H3K27me3,  
835 H3K36me3 ChIP-seq and Input data for HepG2 have been deposited at the European

836 Nucleotide Archive under the accession PRJEB7356.

## 837 **Acknowledgements**

838 We thank the DEEP consortium for their extended help and support. This work was  
839 supported by the Bundesministerium für Bildung und Forschung 'Deutsches Epigenom  
840 Programm' [01KU1216C] (PR, JGH, JW, TM and HRC).

841 *Author Contributions:* JH and HRC developed the methodology. JH implemented the  
842 algorithm, deployed the package and performed analyses. KG and CC provided the primary  
843 human hepatocyte samples. LA and NL performed ChIP-seq experiments. GG performed  
844 whole genome bisulfite sequencing. AS performed RNA-seq experiments. JGH, PR, JW, TM  
845 and HRC supervised the experiments and analyses and acquired funding. JH and HRC wrote  
846 with the help of TM the manuscript. All authors read and approved the manuscript.

## 847 **Disclosure Declaration**

848 *Ethics approval and consent to participate:* T.B.D.

849 *Consent for publication:* T.B.D.

850 *Competing interests:* We declare no competing interests.

## 851 **References**

- 852 1. Johnson DS, Mortazavi A, Myers RM, Wold B. Genome-wide mapping of in vivo  
853 protein-DNA interactions. *Science*. 2007;316:1497–502.
- 854 2. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the  
855 human genome. *Nature*. 2012;489:57–74.
- 856 3. Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen  
857 A, et al. Integrative analysis of 111 reference human epigenomes. *Nature*.  
858 2015;518:317–30.
- 859 4. Thomas-Chollier M, Darbo E, Herrmann C, Defrance M, Thieffry D, van Helden J. A

- 860 complete workflow for the analysis of full-size ChIP-seq (and similar) data sets using  
861 peak-motifs. *Nat Protoc.* 2012;7:1551–68.
- 862 5. Mammana A, Chung H-R. Chromatin segmentation based on a probabilistic model for  
863 read counts explains a large portion of the epigenome. *Genome Biol.* 2015;16:151.
- 864 6. Perner J, Lasserre J, Kinkley S, Vingron M, Chung H-R. Inference of interactions  
865 between chromatin modifiers and histone modifications: from ChIP-Seq data to  
866 chromatin-signaling. *Nucleic Acids Res.* 2014;42:13689–95.
- 867 7. Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, et al. Distinct and  
868 predictive chromatin signatures of transcriptional promoters and enhancers in the  
869 human genome. *Nat. Genet.* 2007;39:311–8.
- 870 8. Heintzman ND, Hon GC, Hawkins RD, Kheradpour P, Stark A, Harp LF, et al. Histone  
871 modifications at human enhancers reflect global cell-type-specific gene expression.  
872 *Nature.* 2009;459:108–12.
- 873 9. Tsankov AM, Gu H, Akopian V, Ziller MJ, Donaghey J, Amit I, et al. Transcription factor  
874 binding dynamics during human ES cell differentiation. *Nature.* 2015;518:344–9.
- 875 10. Karlič R, Chung H-R, Lasserre J, Vlahoviček K, Vingron M. Histone modification levels  
876 are predictive for gene expression. *Proc. Natl. Acad. Sci. U.S.A.* 2010;107:2926–31.
- 877 11. Bonn S, Zinzen RP, Girardot C, Gustafson EH, Perez-Gonzalez A, Delhomme N, et al.  
878 Tissue-specific analysis of chromatin state identifies temporal signatures of enhancer  
879 activity during embryonic development. *Nat. Genet.* 2012;44:148–56.
- 880 12. Sims D, Sudbery I, Iltott NE, Heger A, Ponting CP. Sequencing depth and coverage: key  
881 considerations in genomic analyses. *Nat Rev Genet.* 2014;15:121–32.
- 882 13. Vega VB, Cheung E, Palanisamy N, Sung W-K. Inherent signals in sequencing-based  
883 Chromatin-ImmunoPrecipitation control libraries. *PLoS ONE.* 2009;4:e5241.
- 884 14. Jain D, Baldi S, Zabel A, Straub T, Becker PB. Active promoters give rise to false  
885 positive “Phantom Peaks” in ChIP-seq experiments. *Nucleic Acids Res.* 2015;43:6959–68.
- 886 15. Meyer CA, Liu XS. Identifying and mitigating bias in next-generation sequencing  
887 methods for chromatin biology. *Nat Rev Genet.* 2014;15:709–21.
- 888 16. Flensburg C, Kinkel SA, Keniry A, Blewitt ME, Oshlack A. A comparison of control  
889 samples for ChIP-seq of histone modifications. *Front Genet.* 2014;5:329.
- 890 17. Liang K, Keles S. Normalization of ChIP-seq data with control. *BMC Bioinformatics.*  
891 2012;13:199.
- 892 18. Xu H, Handoko L, Wei X, Ye C, Sheng J, Wei C-L, et al. A signal-noise model for  
893 significance analysis of ChIP-seq with negative control. *Bioinformatics.* 2010;26:1199–  
894 204.
- 895 19. Diaz A, Park K, Lim DA, Song JS. Normalization, bias correction, and peak calling for

- 896   ChIP-seq. *Stat Appl Genet Mol Biol*. 2012;11:Article9.
- 897   20. Zhang Y, Liu T, Meyer C, Eeckhoutte J, Johnson D, Bernstein B, et al. Model-based  
898   Analysis of ChIP-Seq (MACS). *Genome Biol*. 2008;9:R137.
- 899   21. Kumar V, Muratani M, Rayan NA, Kraus P, Lufkin T, Ng HH, et al. Uniform, optimal  
900   signal processing of mapped deep-sequencing data. *Nat. Biotechnol*. 2013;31:615–22.
- 901   22. Ji H, Jiang H, Ma W, Johnson DS, Myers RM, Wong WH. An integrated software  
902   system for analyzing ChIP-chip and ChIP-seq data. *Nat Biotech*. Nature Publishing Group;  
903   2008;26:1293–300.
- 904   23. Kharchenko PV, Tolstorukov MY, Park PJ. Design and analysis of ChIP-seq  
905   experiments for DNA-binding proteins. *Nat. Biotechnol*. 2008;26:1351–9.
- 906   24. Harmanci A, Rozowsky J, Gerstein M. MUSIC: identification of enriched regions in  
907   ChIP-Seq experiments using a mappability-corrected multiscale signal processing  
908   framework. *Genome Biol*. 2014;15:474.
- 909   25. Xu H, Wei C-L, Lin F, Sung W-K. An HMM approach to genome-wide identification of  
910   differential histone modification sites from ChIP-seq data. *Bioinformatics*.  
911   2008;24:2344–9.
- 912   26. Allhoff M, Seré K, Chauvistré H, Lin Q, Zenke M, Costa IG. Detecting differential peaks  
913   in ChIP-seq signals with ODIN. *Bioinformatics*. 2014;30:3467–75.
- 914   27. Heinig M, Colomé-Tatché M, Taudt A, Rintisch C, Schafer S, Pravenec M, et al.  
915   histoneHMM: Differential analysis of histone modifications with broad genomic  
916   footprints. *BMC Bioinformatics*. 2015;16:60.
- 917   28. Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data via  
918   the EM Algorithm on JSTOR. *Journal of the royal statistical society ....* 1977.
- 919   29. Ernst J, Kellis M. ChromHMM: automating chromatin-state discovery and  
920   characterization. *Nat. Methods*. 2012;9:215–6.
- 921   30. Nair NU, Kumar S, Moret BME, Bucher P. Probabilistic partitioning methods to find  
922   significant patterns in ChIP-Seq data. *Bioinformatics*. 2014;30:2406–13.
- 923   31. Ramírez F, Dünder F, Diehl S, Grüning BA, Manke T. deepTools: a flexible platform for  
924   exploring deep-sequencing data. *Nucleic Acids Res*. 2014;42:W187–91.
- 925   32. Dialsingh I, Austin SR, Altman NS. Estimating the proportion of true null hypotheses  
926   when the statistics are discrete. *Bioinformatics*. 2015;31:2303–9.
- 927   33. Storey JD. A direct approach to false discovery rates. *Journal of the Royal Statistical*  
928   *Society: Series B (Statistical Methodology)*. 2002;64:479–98.
- 929   34. Team RC. R: A Language and Environment for Statistical Computing. Vienna, Austria:  
930   R Foundation for Statistical Computing; 2015.
- 931   35. Ooi SKT, Qiu C, Bernstein E, Li K, Jia D, Yang Z, et al. DNMT3L connects unmethylated

- 932 lysine 4 of histone H3 to de novo methylation of DNA. *Nature*. 2007;448:714–7.
- 933 36. Long HK, Sims D, Heger A, Blackledge NP, Kutter C, Wright ML, et al. Epigenetic  
934 conservation at gene regulatory elements revealed by non-methylated DNA profiling in  
935 seven vertebrates. *Elife*. 2013;2:e00348.
- 936 37. Hu J-L, Zhou BO, Zhang R-R, Zhang K-L, Zhou J-Q, Xu G-L. The N-terminus of histone  
937 H3 is required for de novo DNA methylation in chromatin. *Proc. Natl. Acad. Sci. U.S.A.*  
938 2009;106:22187–92.
- 939 38. Kim A, Kiefer CM, Dean A. Distinctive signatures of histone methylation in  
940 transcribed coding and noncoding human beta-globin sequences. *Mol. Cell. Biol.*  
941 2007;27:1271–9.
- 942 39. Baubec T, Colombo DF, Wirbelauer C, Schmidt J, Burger L, Krebs AR, et al. Genomic  
943 profiling of DNA methyltransferases reveals a role for DNMT3B in genic methylation.  
944 *Nature*. 2015;520:243–7.
- 945 40. Feng J, Liu T, Qin B, Zhang Y, Liu XS. Identifying ChIP-seq enrichment using MACS. *Nat*  
946 *Protoc*. 2012;7:1728–40.
- 947 41. Xing H, Mo Y, Liao W, Zhang MQ. Genome-wide localization of protein-DNA binding  
948 and histone modification by a Bayesian change-point method with ChIP-seq data. *PLoS*  
949 *Comput. Biol.* 2012;8:e1002613.
- 950 42. Eisenberg E, Levanon EY. Human housekeeping genes, revisited. *Trends Genet.*  
951 2013;29:569–74.
- 952 43. Kinkley S, Helmuth J, Polansky JK, Dunkel I, Gasparoni G, Fröhler S, et al. reChIP-seq  
953 reveals widespread bivalency of H3K4me3 and H3K27me3 in CD4(+) memory T cells. *Nat*  
954 *Commun*. 2016;7:12514.
- 955 44. Teytelman L, Thurtle DM, Rine J, van Oudenaarden A. Highly expressed loci are  
956 vulnerable to misleading ChIP localization of multiple unrelated proteins. *Proc. Natl.*  
957 *Acad. Sci. U.S.A.* 2013;110:18602–7.
- 958 45. Wang J, Jia ST, Jia S. New Insights into the Regulation of Heterochromatin. *Trends*  
959 *Genet*. 2016;32:284–94.
- 960 46. Frietze S, O'Geen H, Blahnik KR, Jin VX, Farnham PJ. ZNF274 recruits the histone  
961 methyltransferase SETDB1 to the 3' ends of ZNF genes. *PLoS ONE*. 2010;5:e15082.
- 962 47. Saxonov S, Berg P, Brutlag DL. A genome-wide analysis of CpG dinucleotides in the  
963 human genome distinguishes two distinct classes of promoters. *Proc. Natl. Acad. Sci.*  
964 *U.S.A.* 2006;103:1412–7.
- 965 48. Müller J, Hart CM, Francis NJ, Vargas ML, Sengupta A, Wild B, et al. Histone  
966 methyltransferase activity of a Drosophila Polycomb group repressor complex. *Cell*.  
967 2002;111:197–208.
- 968 49. Kuzmichev A, Nishioka K, Erdjument-Bromage H, Tempst P, Reinberg D. Histone

- 969 methyltransferase activity associated with a human multiprotein complex containing the  
970 Enhancer of Zeste protein. *Genes Dev.* 2002;16:2893–905.
- 971 50. Cao R, Wang L, Wang H, Xia L, Erdjument-Bromage H, Tempst P, et al. Role of histone  
972 H3 lysine 27 methylation in Polycomb-group silencing. *Science.* 2002;298:1039–43.
- 973 51. Czermin B, Melfi R, McCabe D, Seitz V, Imhof A, Pirrotta V. Drosophila enhancer of  
974 Zeste/ESC complexes have a histone H3 methyltransferase activity that marks  
975 chromosomal Polycomb sites. *Cell.* 2002;111:185–96.
- 976 52. Sardet C, Vidal M, Cobrinik D, Geng Y, Onufryk C, Chen A, et al. E2F-4 and E2F-5, two  
977 members of the E2F family, are expressed in the early phases of the cell cycle. *Proc. Natl.*  
978 *Acad. Sci. U.S.A.* 1995;92:2403–7.
- 979 53. Sylvestre Y, De Guire V, Querido E, Mukhopadhyay UK, Bourdeau V, Major F, et al. An  
980 E2F/miR-20a autoregulatory feedback loop. *J. Biol. Chem.* 2007;282:2135–43.
- 981 54. Ramboer E, De Craene B, De Kock J, Vanhaecke T, Berx G, Rogiers V, et al. Strategies  
982 for immortalization of primary hepatocytes. *J. Hepatol.* 2014;61:925–43.
- 983 55. Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, et al. Origins and  
984 functional impact of copy number variation in the human genome. *Nature.*  
985 2010;464:704–12.
- 986 56. Shirley MD, Baugher JD, Stevens EL, Tang Z, Gerry N, Beiswanger CM, et al.  
987 Chromosomal variation in lymphoblastoid cell lines. *Hum. Mutat.* 2012;33:1075–86.
- 988 57. Tanay A, O'Donnell AH, Damelin M, Bestor TH. Hyperconserved CpG domains  
989 underlie Polycomb-binding sites. *Proc. Natl. Acad. Sci. U.S.A.* 2007;104:5521–6.
- 990 58. Rotem A, Ram O, Shores N, Sperling RA, Goren A, Weitz DA, et al. Single-cell ChIP-  
991 seq reveals cell subpopulations defined by chromatin state. *Nat. Biotechnol.* 2015.
- 992 59. Speir ML, Zweig AS, Rosenbloom KR, Raney BJ, Paten B, Nejad P, et al. The UCSC  
993 Genome Browser database: 2016 update. *Nucleic Acids Res.* 2016;44:D717–25.
- 994 60. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, et al.  
995 Integrative genomics viewer. *Nat. Biotechnol.* 2011;29:24–6.
- 996 61. Godoy P, Hewitt NJ, Albrecht U, Andersen ME, Ansari N, Bhattacharya S, et al. Recent  
997 advances in 2D and 3D in vitro systems using primary hepatocytes, alternative  
998 hepatocyte sources and non-parenchymal liver cells and their use in investigating  
999 mechanisms of hepatotoxicity, cell signaling and ADME. *Arch. Toxicol.* 2013;87:1315–  
1000 530.
- 1001 62. Arrigoni L, Richter AS, Betancourt E, Bruder K, Diehl S, Manke T, et al. Standardizing  
1002 chromatin research: a simple and universal method for ChIP-seq. *Nucleic Acids Res.*  
1003 2016;44:e67.
- 1004 63. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas MA, et al.  
1005 Transcriptome and genome sequencing uncovers functional variation in humans.

- 1006 Nature. 2013;501:506–11.
- 1007 64. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R. TopHat2: accurate alignment of  
1008 transcriptomes in the presence of insertions, deletions and gene fusions. *Genome ...*  
1009 2013.
- 1010 65. Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-  
1011 throughput sequencing data. *Bioinformatics*. 2015;31:166–9.
- 1012 66. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al. deepTools2: a  
1013 next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.*  
1014 2016;44:W160–5.
- 1015 67. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, et al. Landscape of  
1016 transcription in human cells. *Nature*. 2012;489:101–8.
- 1017 68. FANTOM Consortium and the RIKEN PMI and CLST (DGT), Forrest ARR, Kawaji H,  
1018 Rehli M, Baillie JK, de Hoon MJL, et al. A promoter-level mammalian expression atlas.  
1019 *Nature*. 2014;507:462–70.
- 1020 69. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler  
1021 transform. *Bioinformatics*. 2010;26:589–95.
- 1022 70. Liu Y, Siegmund KD, Laird PW, Berman BP. Bis-SNP: combined DNA methylation and  
1023 SNP calling for Bisulfite-seq data. *Genome Biol*. 2011;13:R61–1.
- 1024 71. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al.  
1025 GENCODE: the reference human genome annotation for The ENCODE Project. *Genome*  
1026 *Res*. 2012;22:1760–74.
- 1027 72. Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, et al. Software  
1028 for computing and annotating genomic ranges. *PLoS Comput. Biol*. 2013;9:e1003118.
- 1029 73. Aibar S, Fontanillo C, Droste C, Las Rivas De J. Functional Gene Networks: R/Bioc  
1030 package to generate and analyse gene networks derived from functional enrichment and  
1031 clustering. *Bioinformatics*. 2015;31:1686–8.
- 1032