

Journal of Computational Biology: <http://mc.manuscriptcentral.com/liebert/jcb>

## A stationary wavelet entropy based clustering approach accurately predicts gene expression

Journal:	<i>Journal of Computational Biology</i>
Manuscript ID:	Draft
Manuscript Type:	Original Paper
Keyword:	algorithms, genetic analysis, genome analysis, GENE EXPRESSION, next generation sequencing
Abstract:	<p>Studying epigenetic landscapes is important to understand the condition for gene regulation. Clustering is a useful approach to study epigenetic landscapes by grouping genes based on their epigenetic conditions. However, classical clustering approaches which often use a representative value of the signals in a fixed-sized window do not fully use the information written in the epigenetic landscapes. Clustering approaches to maximize the information of the epigenetic signals are necessary for better understanding of gene regulatory environments.</p> <p>For effectively clustering of multi-dimensional epigenetic signals, we developed a method called Dewer, which uses the entropy of stationary wavelet of epigenetic signals inside enriched regions for gene clustering. Interestingly, the gene expression levels were highly correlated with the entropy levels of epigenetic signals. Dewer separates genes better than a window-based approach in the assessment using gene expression and achieved a correlation coefficient (CC) above 0.9 without using any training procedure. Our results show that the changes of the epigenetic signals are useful to study gene regulation.</p>

SCHOLARONE™  
Manuscripts

# A stationary wavelet entropy based clustering approach accurately predicts gene expression

Nha Nguyen<sup>1,2</sup>, An Vo<sup>3</sup>, and Kyoung-Jae Won<sup>1,2</sup>

<sup>1</sup>Department of Genetics, <sup>2</sup>Institute for Diabetes, Obesity and Metabolism, School of Medicine, University of Pennsylvania, 3400 Civic Center Blvd, Philadelphia, PA, 19104. <sup>3</sup>The Feinstein Institute for Medical Research, 350 Community Drive, Manhasset, NY 11030, USA  
wonk@mail.med.upenn.edu

**Abstract.** Studying epigenetic landscapes is important to understand the condition for gene regulation. Clustering is a useful approach to study epigenetic landscapes by grouping genes based on their epigenetic conditions. However, classical clustering approaches which often use a representative value of the signals in a fixed-sized window do not fully use the information written in the epigenetic landscapes. Clustering approaches to maximize the information of the epigenetic signals are necessary for better understanding of gene regulatory environments.

For effective clustering of multi-dimensional epigenetic signals, we developed a method called Dewer, which uses the entropy of stationary wavelet of epigenetic signals inside enriched regions for gene clustering. Interestingly, the gene expression levels were highly correlated with the entropy levels of epigenetic signals. Dewer separates genes better than a window-based approach in the assessment using gene expression and achieved a correlation coefficient (CC) above 0.9 without using any train-

1  
2  
3 ing procedure. Our results show that the changes of the epigenetic signals  
4  
5 are useful to study gene regulation.  
6  
7  
8  
9

## 10 11 **1 Introduction**

12  
13  
14 The epigenetic landscape, represented by DNA methylation, modifications to  
15  
16 histones, and other proteins that package the genome, regulates the function of cells  
17  
18 by regulating gene activity (Bernstein et al., 2007; Kouzarides, 2007). To understand  
19  
20 the complex languages of these epigenetic landscapes, gene clustering has been ap-  
21  
22 plied to epigenomic data to study a wide range of biological questions including de-  
23  
24 velopment (Lister et al., 2009; Meissner, 2010; Mikkelsen et al., 2007; Mikkelsen et  
25  
26 al., 2010; Xie et al., 2013; Yu et al., 2013), cancer (Baylin and Jones, 2011; Jones and  
27  
28 Martienssen, 2005; Laird, 2003), and aging (Baylin and Jones, 2011; Li et al., 2011;  
29  
30 Liu et al., 2011). Clustering approaches provided insights into the dynamics of epige-  
31  
32 netic gene changes. Additionally, clustering has been developed to identify co-  
33  
34 occurring histone modification marks or ‘histone codes’ (Barski et al., 2007;  
35  
36 Heintzman et al., 2009; Heintzman et al., 2007; Won et al., 2008; Xie et al., 2013), as  
37  
38 well as to study cell type- or species-specific gene regulation (Ernst et al., 2011; Won  
39  
40 et al., 2013; Yu et al., 2013). Often, clustering approaches have not fully utilized the  
41  
42 characteristics of epigenetic signals because of the difficulties of understanding com-  
43  
44 plex signals that occur in different combinations in epigenetic data. Methodological  
45  
46 development for clustering epigenomic landscapes is required for comprehensive un-  
47  
48 derstanding of gene regulation.  
49  
50  
51  
52

53  
54 To cluster epigenomic data more effectively using their statistical characteris-  
55  
56 tics, we developed a new approach called Dewer. Dewer uses the measurement of  
57  
58  
59  
60

1  
2  
3 entropy of epigenetic signals for clustering. Entropy is the expected value of the in-  
4 formation contained in a random variable (Shannon, 1948). Entropy has been widely  
5 used in many areas including information theory, signal processing and biology  
6 (Daily et al., 2010; Li et al., 2004; Menayo et al., 2014; Shin et al., 2007; Swartz et  
7 al., 1999; Yogesan et al., 1996). Previously, entropy has been applied to genome-wide  
8 datasets to check tissue specificity and co-localization of signals (Sun et al., 2011;  
9 Won et al., 2013). A quantitative method named QDMR also used entropy to identify  
10 differentially methylated regions (DMRs) (Zhang et al., 2011). However, these stud-  
11 ies were limited to studying variation across samples or tissues in a defined genomic  
12 position. Compared with the previous approaches, Dewer employed sophisticated  
13 ways of calculating the entropies in both original and wavelet domain for better dis-  
14 criminative power in clustering genes. Using entropy as a metric, Dewer fully uses the  
15 information contained in the distribution of the multi-dimensional epigenetic data,  
16 rather than condensing the data to mere statistics as in traditional window-based ap-  
17 proaches.

18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36 To apply entropy effectively, Dewer detects areas enriched for multiple histone  
37 marks. The stacked histone marks in a region forms a 2-dimensional (2D) data. Ex-  
38 ploiting 2D spaces is important to fully understand the spatial combinatorial histone  
39 modification patterns for gene regulation. Most of border detectors have been devel-  
40 oped to use only a single mark (Heinz, et al., 2010; Rashid, et al., 2011; Zang, et al.,  
41 2009; Micsinai, et al., 2012; Song and Smith, 2011). To identify enriched regions  
42 from multiple histone modification marks, Dewer upgraded SeqW (Nguyen et al.,  
43 2014b), a method developed by our group for better border detection. Entropy is then  
44 applied to the identified enriched regions. In this paper, we show that both entropy  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 and border detection contribute to the Dewer's discriminative power for gene cluster-  
4  
5 ing.  
6

7 We found that the clusters identified by Dewer have a superior discriminative  
8  
9 power over a window-based approach when we compared the gene expression levels  
10  
11 among the clusters. Interestingly, gene expressions levels were highly correlated with  
12  
13 the entropy levels, suggestive of the importance of the shape of epigenetic signals in  
14  
15 studying gene regulation.  
16  
17

## 18 19 20 **2 Methods**

21  
22 To calculate the entropy from multiple histone modification signals, Dewer uses 2-  
23  
24 directional information retrieved from epigenomic data such as entropy and stationary  
25  
26 wavelet (SW) entropy in the detected areas enriched inside histone modification sig-  
27  
28 nals.  
29  
30

### 31 32 **2.1 Data preprocessing**

33  
34 For clustering analysis, we used activating marks at around transcription start sites  
35  
36 (TSSs): H3K4me1, H3K4me2, H3K4me3, H3K9ac, H3K18ac, H3K27ac, and  
37  
38 H4K91ac from human fetal lung fibroblasts (IMR90) (Bernstein et al., 2010). Tags  
39  
40 were normalized using reads per kilo base per million (RPKM) for 10-bp bins.  
41  
42

43 We also used histone marks enriched in active gene body: H3K36me3 in murine adi-  
44  
45 pocytes (3T3L1) (Mikkelsen et al., 2010) and H3K36me3, H3K79me1 and  
46  
47 H3K79me2 in IMR90 (Bernstein et al., 2010).  
48  
49  
50

### 51 52 **2.2 Identifying segments**

53  
54 We further upgraded SeqW (Nguyen et al., 2014b), our previous border detector, to  
55  
56 detect borders of multiple histone modification signals. To reduce computational  
57  
58  
59  
60

cost, Haar wavelet moving (HWM) (Nguyen et al., 2014a) borrowed from SeqW is still used in Dewer to detect domains, which represent larger regions with epigenetic signals. Inside each domain, instead of detecting border directly in SeqW, Dewer estimates border indirect way by obtaining many segments by assessing signal enrichment against input. Neighboring segments are combined to form enriched regions. A domain usually has a number of enriched regions (Figure 2). Histone modification signals can be written as

$$g(t, m) = M(m)f(t), \quad (1)$$

where  $f(t)$  is a function to normalize tag counts at a genomic position  $t$  of a histone modification mark  $m$ . Compared with SeqW, we used  $M(m)$  which shows relative enrichments across  $m$  histone modification marks. With this new representation in Eq. (1), Dewer can detect borders from multiple marks.

Assuming a mixture of Gaussians (MoG) for a nucleosome marked by histone modification signals as in (Nguyen et al., 2014a), (Nguyen et al., 2014b) and (Nguyen et al., 2010), we obtain

$$f(t) = \sum_i f_i(t) = \sum_i A_i e^{-(t-\mu_i)^2/(2\sigma_i^2)}, \quad (2)$$

where  $\mu_i$ , and  $\sigma_i$  are the center and the standard deviation of a peak of a nucleosome, respectively. To estimate the borders of MoG signals, Dewer uses the zero-crossing lines across the multi-wavelet scale. A zero-crossing is a point where the sign of a function changes. A zero-crossing line is obtained by connecting the zero-crossing points obtained over the wavelet scale  $s$  (Nguyen et al., 2014a).

Wavelet transform converts signals to the WD using a convolution operator.

$$\begin{aligned} Wg(u, m, s) &= M(m)Wf(u, s) \\ &= M(m) \int_{-\infty}^{\infty} f_i(t) \frac{1}{\sqrt{s}} \psi^* \left( \frac{t-u}{s} \right) dt = M(m)(f_i * \widetilde{\psi}_s)(u), \end{aligned} \quad (3)$$

where  $s$  denotes the scale,  $u$  is the genome position in WD and  $\tilde{\Psi}(t) = \frac{1}{\sqrt{s}} \Psi^* \left( -\frac{t}{s} \right)$ .

As described more in detail in (Nguyen et al., 2014a), we obtain

$$Wg(w, s) = \beta (jw)^n \frac{1}{\sqrt{2\alpha}} e^{-\frac{w^2}{4\alpha}} e^{-i\mu_i w}, \quad (4)$$

The inverse fast Fourier transform (iFFT) of Eq.(4) becomes

$$Wg(u, m, s, n) = \beta \frac{d^n}{du^n} e^{-\alpha(u-\mu_i)^2}, \quad (5)$$

where  $\beta = \frac{A_i \sigma_i s^n \sqrt{2\alpha}}{\sqrt{\Gamma(n+\frac{1}{2})}}$  and  $\alpha = \frac{2}{s^2 + \sigma_i^2}$ .

To detect the borders of a peak, we use the second derivative wavelet ( $n = 2$ ) (DOG2) instead of DOG3 in SeqW, where  $n = 1, 2, 3$  correspond to each order of derivative.

$$Wf(u, s, 2) = -2\alpha\beta \left[ 1 - 2\alpha(u - \mu_i)^2 \right] e^{-\alpha(u-\mu_i)^2}. \quad (6)$$

The zero-crossing points are the parameters of the Gaussians when

$Wf(u_0, s, 2) = 0$ . Then we have

$$u_0 = \mu_i \pm \sqrt{\sigma_i^2 + s^2}. \quad (7)$$

If  $s$  is small compared with  $\sigma_i$ , Eq. (7) becomes

$$u_0 \approx \mu_i \pm \sigma_i. \quad (8)$$

In summary,  $u_0(s)$  in Eq. (7) draws a line over the scale  $s$ . Eq. (8) indicates that the zero-crossing line approximates to  $\sigma$  away from the center of a peak. Eq. (7) and (8) were used to detect the border of the peaks. We regarded a segment as the region defined by the borders of a MoGs. In the regions where nucleosomes are packed, histone modification signals can be represented with a number of segments. To remove falsely predicted segments, each segment is assessed using false discovery rates (FDRs) after considering enrichment against background (input).

### 2.3 False discovery rate (FDR) control

We used the *mattes* and the *mafdr* functions in Matlab R2012a to calculate the statistical significance of the signals in a segment against background (input) after dividing a segment into three sub-segments. P-values were calculated by applying t-test (Huber et al., 2002) on these sub-segments. FDRs were obtained from these p-values (Storey, 2002). Segments with an FDR greater than 0.01 were removed. After this step, adjacent segments were combined to form an enriched region (Figure 2b).

#### 2.4 Gene clustering using entropy and SW entropy

Dewar used Shannon's entropy (Shannon, 1948) for gene clustering. Shannon's entropy has been widely used in signal processing to measure the expected value of the information contents contained in a signal. Entropy  $H(x)$  for a given signal  $x$  can be estimated by

$$H(x) = -\sum_i (P(x_i) \log_2 P(x_i)), \quad (9)$$

where  $P$  represents the histogram of the normalized intensity of histone modification signals inside the detected enriched regions around TSS (1-5Kbps, 5Kbps is default value).

Stationary wavelet transform (SWT)(Mallat, 2009) is a time-frequency analysis that has been widely used in image processing applications such as de-noising (Wang et al., 2003), enhancement (Demirel and Anbarjafari, 2011), segmentation (Deng et al., 2014) and image retrieval (Agarwal et al., 2013). Especially, we used the extracted information from decomposed scales of SWT. SWT (Figure 1) with two levels has eight scales: LL1, LH1, HL1, HH1, LL2, LH2, HL2, and HH2. The stationary wavelet (SW) entropy is obtained from these eight scales. The SW entropy levels are used as the input of clustering. Decomposition part of SWT can be formatted as follows

$$LL_{l+1} = LL_l * L_l(\text{Rows}) * L_l(\text{Columns}),$$

$$\begin{aligned}
LH_{l+1} &= LL_l * L_l(\text{Rows}) * H_l(\text{Columns}), \\
HL_{l+1} &= LL_l * H_l(\text{Rows}) * L_l(\text{Columns}), \\
HH_{l+1} &= LL_l * H_l(\text{Rows}) * H_l(\text{Columns}), \\
L_{l+1} &= L_l(\uparrow 2), \\
H_{l+1} &= H_l(\uparrow 2),
\end{aligned} \tag{10}$$

where  $L_l$  is low pass filter at level  $l$ ,  $H_l$  is high pass filter at level  $l$ ,  $LL_0$  is original image and  $\uparrow 2$  is the up-sampling operation by 2.

Dewer used agglomerative cluster algorithm (Szekely and Rizzo, 2005), a hierarchical clustering approach where each object starts from its own cluster, and pairs of clusters are merged as one move up the hierarchy

### 3 Result

#### 3.1 Dewer overview:

Dewer clusters genes based on the entropy features inside enriched regions. Figure 2a shows the five domains identified using 7 types of histone modification marks (H3K4me1, H3K4me2, H3K4me3, H3K9ac, H3K18ac, H3K27ac, and H4K91ac) in IMR90. A domain contains a number of enriched regions. The derivative wavelet method detects the border of the segments after modeling a segment. P-values and FDRs were calculated to measure the enrichment of epigenetic signals against control input. We applied an FDR cut-off of 0.01 (Figure 2b). The neighboring segments were combined to form a final enriched region. Figure 2c shows two enriched regions around the promoter regions of *Ccn12* and *Mrlp20*. Entropies were calculated in these regions around the TSS. Also, SW entropies were estimated inside eight scales. A hierarchical clustering method is applied to the obtained entropies as

1  
2  
3 the input features. The supplemental document and software are available at  
4 <http://www.naaan.org/nhanguyen/> at Software tab.  
5  
6  
7

### 8 **3.2 Assessing the performance of enriched region detection**

9  
10  
11 Detecting enriched region is an important component as estimating entropy in Dewer.  
12  
13 In this section, we evaluate the performance of enriched region detection.  
14

15  
16 To evaluate the performance of enriched region detection, we compared Dew-  
17  
18 er with previous broad region detectors including SeqW (Nguyen et al., 2014b), Sicer  
19  
20 (Zang et al., 2009), RSEG (Song and Smith, 2011), and QESEQ (Micsinai et al.,  
21  
22 2012). We also included ChromHMM (Ernst et al., 2010) because it annotates the  
23  
24 genome using multiple histone modification marks, even though it was not originally  
25  
26 designed for enriched region detection. For this assessment, we used H3K36me3 in  
27  
28 murine adipocytes (3T3L1) and evaluated using the annotated genes in GENCODE  
29  
30 (Harrow et al., 2012). Only the bodies of the active genes whose gene expressions are  
31  
32 higher than averaged value were used.  
33  
34  
35

36  
37 Figure 3 shows the true positive rates (TPRs) against false positive rates  
38  
39 (FPRs). For this test, only a single H3K36me3 (day -2) was used for QESEQ, Rseq,  
40  
41 Sicer and SeqW and four H3K36me3 were used for Dewer and ChromHMM (day -  
42  
43 2,0,3,7). In our test, QESEQ performed better than Rseq and Sicer, in agreement with  
44  
45 the previous observations in (Micsinai et al., 2012). QESEQ performed better than  
46  
47 ChromHMM using four H3K36me3 marks. It is not surprising because ChromHMM  
48  
49 is originally not designed to predict enriched regions.  
50

51  
52 Dewer performed even better than SeqW, suggesting that Dewer uses multi-  
53  
54 dimensional epigenetic signals effectively. ZINBA (Rashid et al., 2011) was not in-  
55  
56 cluded in this comparison because of its' exhaustive running time.  
57  
58  
59  
60

1  
2  
3 Additionally, we tested the performance using active marks for gene body  
4 (H3K36me3, H3K79me1 and H3K79me2) in IMR90 (Figure 4). For this test, we  
5 compared Sicer, QESEQ and SeqW as they performed better than other predictors in  
6 our previous test. A single mark (H3K36me3) was used for SeqW, Sicer and QESEQ.  
7 Dewar used H3K36me3, H3K79me1 and H3K79me2. The test also confirmed that  
8 Dewar has a solid performance in detecting the regions enriched for epigenetic signals.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18

### 19 **3.3 Gene clustering using Dewar**

#### 20 **3.3.1 The clusters identified by Dewar well discriminate gene expression levels.**

21 We assessed the performance of gene clustering by comparing Dewar with a  
22 conventional window-based approach. Window sizes have been selected from one to  
23 five Kbps. With Dewar, window sizes only help to narrow selected enriched regions.  
24 Using seven histone modification marks in IMR90, we performed clustering and gener-  
25 ated six to ten clusters. We used gene expression as a surrogate to check if the iden-  
26 tified clusters separated gene groups effectively.  
27  
28

29 After performing clustering, we evaluated the expression levels of the genes for each  
30 cluster and measured if they are different from each other using the student's t-test.  
31 For this test, we also investigated the effectiveness of the two algorithms that Dewar  
32 employed (detection of enriched regions and calculating entropy). Specifically, to as-  
33 sess the advantages of using enriched regions, we compared the performance of using  
34 entropy with/without enriched region detection. We also tested if entropy improves  
35 discriminative power in a given window. Figure 5 compares the averaged p-values of  
36 all clustering pairs in the tested clustering algorithms. For rigorous assessment, we  
37 investigated the averaged p-values while we increased the size of a window around  
38 TSSs and calculated the mean value and the entropy of the signals. Table 1 compares  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the performance of various wavelet filters including Daubechies, Coiflets, Symlets,  
4 Discrete Meyer, Biorthogonal and Reverse Biorthogonal (Mallat, 2009). Eight scales  
5 (LL1, HL1, LH1, HH1, LL2, HL2, LH2 and HH2) of SWT for each filter were tested.  
6  
7 We found that HH2 of Sym10 performed best. Sym10 is a symmetric filter, suggest-  
8 ing that symmetry and high frequency of Sym10 may be suitable for capturing the  
9 characteristics of histone modification signal.  
10  
11

12  
13  
14  
15  
16 The performances were improved in general as we increased the size of a win-  
17 dow around the TSSs except for the case when we used the mean value of the en-  
18 riched regions (Figure 5). This may be caused by the high intensity of the signals in  
19 some local enriched regions. After 5Kbps, all performances went down (not shown  
20 here). The comparison of using entropy and the mean value in a given window clearly  
21 demonstrated that entropy is a good measure for clustering against the window-based  
22 method. Calculating the mean value performed worst in our test. Dewer, which used  
23 both entropy and the enriched region, performed best in this assessment. This test em-  
24 phasizes that both enriched region and entropy provide discriminative power for gene  
25 clustering.  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40

### 41 **3.3.2 Entropy levels correlate with gene expression levels**

42 Next, we studied the association of gene expression with entropy in sixteen clusters.  
43 We also tested the association of gene expression of the clusters obtained using the  
44 mean value in a window. The averaged gene expression levels for all the clusters were  
45 identified for each configuration: the average value of a window (Figure 6), the entro-  
46 py of a window (Figure 7), and Dewer (the entropy in the enriched regions) (Figure 8).  
47  
48 All clusters were sorted based on their averaged feature values.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 We observed gradual increase of gene expression levels in association with the entro-  
4 py levels (Figure 7 and 8). In Figure 6, the expression level of G14 was higher than  
5 that of G15 and G16 in the window-based approach despite its lower mean value level.  
6  
7  
8  
9  
10 When we used entropy in a window (Figure 7), we observed a better correlation than  
11 using the averaged values. Dewer showed the best correlation in our test (Figure 8),  
12 and suggested that both entropy and enriched region detection are useful in clustering.  
13  
14  
15

16 We observed similar phenomena when using 32 clusters. We investigated the  
17 correlation coefficient of the clusters when we used the epigenetic data 1~5Kbps  
18 around TSSs (Figure 9, 10). The CC between expression levels and the entropy levels  
19 reached 0.9 when using Dewer (Figure 10). We further analyzed the contributions of  
20 each histone modification mark for gene expression (Figure 9). When mean values  
21 were used, the contributions of H3K4me1 and H4K91ac for gene expression predic-  
22 tion were very weak (Figure 11, 12). This is mainly because of the low intensity of the  
23 signals. When all marks were combined, the window-based method resulted in a cor-  
24 relation approximately 0.6. The same comparison using the entropy levels resulted in  
25 much stronger correlation in general even when using H3K4me1 only. Even though  
26 entropy of H3K4me2 and H3K4me3 were the highest, their associations with gene  
27 expression were relatively low, especially for H3K4me3. This may be because  
28 H3K4me3 is also enriched for transcriptionally paused genes (Guenther et al., 2007).  
29  
30 We observed that the correlation was the highest when using H3K27ac, consistent  
31 with previous studies (Karlic et al., 2010). H3K9ac and H3K18ac also had high corre-  
32 lation though the signal of H3K18ac was very weak (Figure 11, 12). Using all marks  
33 giving the highest CC may suggest that Dewer makes use of multiple marks effective-  
34 ly by employing entropy.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Figure 13 shows an example of the epigenetic landscapes around the TSS of *Cd151*  
4 and *Orc1* in IMR90. Though *Orc1* had a higher mean value around the TSS, its ex-  
5 pression was much lower compared with *Cd151*. The entropy of *Cd151* was much  
6 higher, which correlated well with their gene expression levels. These examples as  
7 well as our results suggest that entropy is a good marker for evaluating gene expres-  
8 sion  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

#### 20 **4 Discussion**

21  
22 Understanding of the complex nature of epigenomic signals is still a challenging  
23 problem. Gene clustering has been successfully applied to study epigenetic regulation.  
24 In this paper, we introduce a method to improve the clustering performance by using  
25 the signal processing approaches. Dewar effectively uses the nature of epigenetic sig-  
26 nals for clustering by employing the entropy from the SW. We found that the genes  
27 clustered by Dewar better dissected the gene groups in terms of gene expression. In-  
28 terestingly, gene expression levels well correlated with the entropy levels of epigenet-  
29 ic data. Our results showed that signal processing approaches effectively use the char-  
30 acteristics of epigenetic signals for gene clustering.  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

43 Our tests showed a number of advantages of using entropies. Especially, entro-  
44 py can be applied even when the signal intensity is low. This suggests that using en-  
45 tropy or SW entropy is a more robust way to study gene regulation than using mean-  
46 based approaches. Also, entropy can be obtained from multiple histone modification  
47 marks. We found that using the entropy performed best when using all marks together.  
48 This is because the relationships among histone modification marks are well stored in  
49 the entropy of the 2D epigenetic images.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5 Dewer was not designed to predict gene expression using histone modification  
6  
7 (Dong et al., 2012; Karlic et al., 2010; Kumar et al., 2013). We used gene expression  
8  
9 as a surrogate to show the discriminative power of Dewer in clustering. Our results, at  
10  
11 least, suggest that statistical features such as entropy can be a good measure to predict  
12  
13 gene expression.  
14

15  
16 To apply entropy effectively, we restricted entropy to the enriched regions for  
17  
18 histone modification signals. SW entropy in HH2 scale produced the best perfor-  
19  
20 mance, suggesting that the information in high frequency scale is important for gene  
21  
22 clustering. This further suggests that the shape of epigenetic signal is important than  
23  
24 the averaged signals (low frequency) for gene clustering.  
25  
26

27  
28 Clustering analyses have identified various histone codes including bivalent  
29  
30 promoters (Bernstein et al., 2006), poised enhancers (Creyghton et al., 2010; Rada-  
31  
32 Iglesias and Wysocka, 2011), and alternative splicing (Luco and Misteli, 2011; Luco  
33  
34 et al., 2010). Computational approaches have been applied to identify co-enriched  
35  
36 histone modifications (Nguyen et al., 2014a; Rajagopal et al., 2013; Santoni, 2013;  
37  
38 Ucar et al., 2011). While previous approaches identified combinatorial patterns, it is  
39  
40 hard to interpret the clusters obtained using entropy using the same manner. Instead,  
41  
42 we use the information contents residing in epigenetic signals for clustering. The re-  
43  
44 sults provide us with a good measurement to detect gene expression. Though we re-  
45  
46 stricted our study to promoter regions to investigate the relationships with genes, our  
47  
48 approach can easily be applied to distal regulatory regions.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## 5 ACKNOWLEDGEMENT

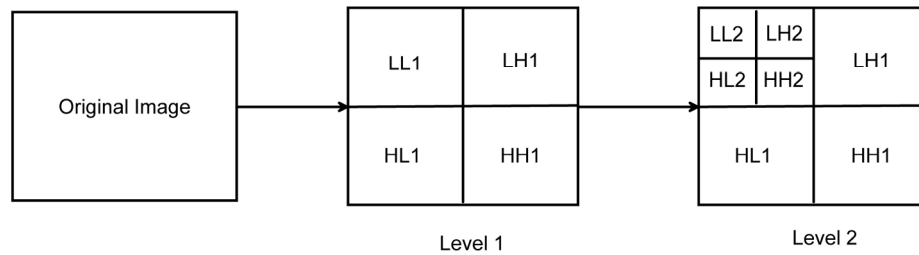
This work is supported by R21-DK098769 and P30-DK19525 from the National Institutes of Diabetes, and Digestive and Kidney Diseases and the Diabetes Research Center at the University of Pennsylvania

## 6 References

- Agarwal, S, Verma, AK, Singh, P. 2013. Content Based Image Retrieval using Discrete Wavelet Transform and Edge Histogram Descriptor. In *Information Systems and Computer Networks (ISCON), 2013 International Conference on*. 19-23.
- Barski, A, Cuddapah, S, Cui, K, et al. 2007. High-resolution profiling of histone methylations in the human genome. *Cell* 129,823-837.
- Baylin, SB, Jones, PA. 2011. A decade of exploring the cancer epigenome - biological and translational implications. *Nature reviews*. *Cancer* 11,726-734.
- Bernstein, BE, Meissner, A, Lander, ES. 2007. The mammalian epigenome. *Cell* 128,669-681.
- Bernstein, BE, Mikkelsen, TS, Xie, X, et al. 2006. A bivalent chromatin structure marks key developmental genes in embryonic stem cells. *Cell* 125,315-326.
- Bernstein, BE, Stamatoyannopoulos, JA, Costello, JF, et al. 2010. The NIH Roadmap Epigenomics Mapping Consortium. *Nat Biotechnol* 28,1045-1048.
- Creyghton, MP, Cheng, AW, Welstead, GG, et al. 2010. Histone H3K27ac separates active from poised enhancers and predicts developmental state. *Proc Natl Acad Sci U S A* 107,21931-21936.
- Daily, K, Rigor, P, Christley, S, et al. 2010. Data structures and compression algorithms for high-throughput sequencing technologies. *Bmc Bioinformatics* 11,514.
- Demirel, H, Anbarjafari, G. 2011. IMAGE Resolution Enhancement by Using Discrete and Stationary Wavelet Decomposition. *Ieee Transactions on Image Processing* 20,1458-1460.
- Deng, J, Ban, YF, Liu, JS, et al. 2014. Hierarchical Segmentation of Multitemporal RADARSAT-2 SAR Data Using Stationary Wavelet Transform and Algebraic Multigrid Method. *Ieee Transactions on Geoscience and Remote Sensing* 52,4353-4363.
- Dong, X, Greven, MC, Kundaje, A, et al. 2012. Modeling gene expression using chromatin features in various cellular contexts. *Genome biology* 13,R53.
- Ernst, J, Kheradpour, P, Mikkelsen, TS, et al. 2011. Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* 473,43-49.
- Ernst, J, Plasterer, HL, Simon, I, et al. 2010. Integrating multiple evidence sources to predict transcription factor binding in the human genome. *Genome Res* 20,526-536.
- Guenther, MG, Levine, SS, Boyer, LA, et al. 2007. A chromatin landmark and transcription initiation at most promoters in human cells. *Cell* 130,77-88.
- Harrow, J, Frankish, A, Gonzalez, JM, et al. 2012. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome research* 22,1760-1774.
- Heintzman, ND, Hon, GC, Hawkins, RD, et al. 2009. Histone modifications at human enhancers reflect global cell-type-specific gene expression. *Nature*.
- Heintzman, ND, Stuart, RK, Hon, G, et al. 2007. Distinct and predictive chromatin signatures of transcriptional promoters and enhancers in the human genome. *Nat Genet* 39,311-318.
- Huber, W, von Heydebreck, A, Sultmann, H, et al. 2002. Variance stabilization applied to microarray data calibration and to the quantification of differential expression. *Bioinformatics* 18 Suppl 1,S96-104.
- Jones, PA, Martienssen, R. 2005. A blueprint for a Human Epigenome Project: the AACR Human Epigenome Workshop. *Cancer research* 65,11241-11246.
- Karlic, R, Chung, HR, Lasserre, J, et al. 2010. Histone modification levels are predictive for gene expression. *Proc Natl Acad Sci U S A* 107,2926-2931.
- Kouzarides, T. 2007. Chromatin modifications and their function. *Cell* 128,693-705.
- Kumar, V, Muratani, M, Rayan, NA, et al. 2013. Uniform, optimal signal processing of mapped deep-sequencing data. *Nature biotechnology* 31,615-622.

- 1  
2  
3 Laird, PW. 2003. The power and the promise of DNA methylation markers. *Nature reviews. Cancer* 3,253-266.
- 4 Li, H, Zhang, K, Jiang, T. 2004. Minimum entropy clustering and applications to gene expression analysis.  
5 Proceedings / IEEE Computational Systems Bioinformatics Conference, CSB. IEEE Computational  
6 Systems Bioinformatics Conference,142-151.
- 7 Li, Y, Daniel, M, Tollefsbol, TO. 2011. Epigenetic regulation of caloric restriction in aging. *BMC medicine* 9,98.
- 8 Lister, R, Pelizzola, M, Dowen, RH, et al. 2009. Human DNA methylomes at base resolution show widespread  
9 epigenomic differences. *Nature* 462,315-322.
- 10 Liu, L, van Groen, T, Kadish, I, et al. 2011. Insufficient DNA methylation affects healthy aging and promotes age-  
11 related health problems. *Clinical epigenetics* 2,349-360.
- 12 Luco, RF, Misteli, T. 2011. More than a splicing code: integrating the role of RNA, chromatin and non-coding  
13 RNA in alternative splicing regulation. *Current opinion in genetics & development* 21,366-372.
- 14 Luco, RF, Pan, Q, Tominaga, K, et al. 2010. Regulation of alternative splicing by histone modifications. *Science*  
15 327,996-1000.
- 16 Mallat, SG. 2009. *A wavelet tour of signal processing : the sparse way*. Elsevier/Academic Press, Amsterdam ;  
17 Boston.
- 18 Meissner, A. 2010. Epigenetic modifications in pluripotent and differentiated cells. *Nature biotechnology* 28,1079-  
19 1088.
- 20 Menayo, R, Encarnacion, A, Gea, GM, et al. 2014. Sample entropy-based analysis of differential and traditional  
21 training effects on dynamic balance in healthy people. *Journal of motor behavior* 46,73-82.
- 22 Micsinai, M, Parisi, F, Strino, F, et al. 2012. Picking ChIP-seq peak detectors for analyzing chromatin  
23 modification experiments. *Nucleic Acids Res.*
- 24 Mikkelsen, TS, Ku, M, Jaffe, DB, et al. 2007. Genome-wide maps of chromatin state in pluripotent and lineage-  
25 committed cells. *Nature* 448,553-560.
- 26 Mikkelsen, TS, Xu, Z, Zhang, X, et al. 2010. Comparative epigenomic analysis of murine and human  
27 adipogenesis. *Cell* 143,156-169.
- 28 Nguyen, N, Huang, H, Oraintara, S, et al. 2010. Mass spectrometry data processing using zero-crossing lines in  
29 multi-scale of Gaussian derivative wavelet. *Bioinformatics* 26,i659-665.
- 30 Nguyen, N, Vo, A, Won, KJ. 2014a. A wavelet-based method to exploit epigenomic language in the regulatory  
31 region. *Bioinformatics*.
- 32 Nguyen, N, Vo, A, Won, KJ. 2014b. A wavelet approach to detect enriched regions and explore epigenomic  
33 landscapes. *Journal of Computational Biology Accepted*.
- 34 Rada-Iglesias, A, Wysocka, J. 2011. Epigenomics of human embryonic stem cells and induced pluripotent stem  
35 cells: insights into pluripotency and implications for disease. *Genome Med* 3,36.
- 36 Rajagopal, N, Xie, W, Li, Y, et al. 2013. RFECS: a random-forest based algorithm for enhancer identification from  
37 chromatin state. *PLoS computational biology* 9,e1002968.
- 38 Rashid, NU, Giresi, PG, Ibrahim, JG, et al. 2011. ZINBA integrates local covariates with DNA-seq data to identify  
39 broad and narrow regions of enrichment, even within amplified genomic regions. *Genome Biol* 12,R67.
- 40 Santoni, FA. 2013. EMdeCODE: a novel algorithm capable of reading words of epigenetic code to predict  
41 enhancers and retroviral integration sites and to identify H3R2me1 as a distinctive mark of coding versus  
42 non-coding genes. *Nucleic acids research* 41,e48.
- 43 Shannon, CE. 1948. A Mathematical Theory of Communication. *Bell System Technical Journal* 27,379-423.
- 44 Shin, JH, Park, CH, Yang, YJ, et al. 2007. Entropy-based analysis of the non-linear relationship between gene  
45 expression profiles of amplified and non-amplified RNA. *International journal of molecular medicine*  
46 20,905-912.
- 47 Song, Q, Smith, AD. 2011. Identifying dispersed epigenomic domains from ChIP-Seq data. *Bioinformatics*  
48 27,870-871.
- 49 Storey, JD. 2002. A direct approach to false discovery rates. *Journal of the Royal Statistical Society Series B-  
50 Statistical Methodology* 64,479-498.
- 51 Sun, H, Wu, J, Wickramasinghe, P, et al. 2011. Genome-wide mapping of RNA Pol-II promoter usage in mouse  
52 tissues by ChIP-seq. *Nucleic acids research* 39,190-201.
- 53 Swartz, JB, Rothenberg, SJ, Teklehaimanot, S, et al. 1999. Comparison of the entropy technique with two other  
54 techniques for detecting disease clustering using data from children with high blood lead levels.  
55 *American journal of epidemiology* 149,750-760.
- 56 Szekely, GJ, Rizzo, ML. 2005. Hierarchical clustering via joint between-within distances: Extending Ward's  
57 minimum variance method. *Journal of Classification* 22,151-183.
- 58 Ucar, D, Hu, Q, Tan, K. 2011. Combinatorial chromatin modification patterns in the human genome revealed by  
59 subspace clustering. *Nucleic acids research* 39,4063-4075.
- 60

- 1  
2  
3 Wang, XH, Istepanian, RSH, Song, YH. 2003. Microarray image enhancement by denoising using stationary  
4 wavelet transform. *Ieee Transactions on Nanobioscience* 2,184-189.
- 5 Won, KJ, Chepelev, I, Ren, B, et al. 2008. Prediction of regulatory elements in mammalian genomes using  
6 chromatin signatures. *Bmc Bioinformatics* 9,547.
- 7 Won, KJ, Zhang, X, Wang, T, et al. 2013. Comparative annotation of functional regions in the human genome  
8 using epigenomic data. *Nucleic acids research* 41,4423-4432.
- 9 Xie, R, Everett, LJ, Lim, HW, et al. 2013. Dynamic chromatin remodeling mediated by polycomb proteins  
10 orchestrates pancreatic differentiation of human embryonic stem cells. *Cell Stem Cell* 12,224-237.
- 11 Yogesan, K, Jorgensen, T, Albrechtsen, F, et al. 1996. Entropy-based texture analysis of chromatin structure in  
12 advanced prostate cancer. *Cytometry* 24,268-276.
- 13 Yu, P, Xiao, S, Xin, X, et al. 2013. Spatiotemporal clustering of the epigenome reveals rules of dynamic gene  
14 regulation. *Genome research* 23,352-364.
- 15 Zang, C, Schones, DE, Zeng, C, et al. 2009. A clustering approach for identification of enriched domains from  
16 histone modification ChIP-Seq data. *Bioinformatics* 25,1952-1958.
- 17 Zhang, Y, Liu, H, Lv, J, et al. 2011. QDMR: a quantitative method for identification of differentially methylated  
18 regions by entropy. *Nucleic acids research* 39,e58.
- 19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



18 Figure 1. Stationary wavelet transform with two levels. LL1, LH1, HL1, HH1, LL2, LH2, HL2, and HH2 are the  
19 eight scales. L means low frequency. H means high frequency. LH1 mean low frequency in horizontal and  
20 high frequency in vertical at level 1.  
21 174x50mm (300 x 300 DPI)

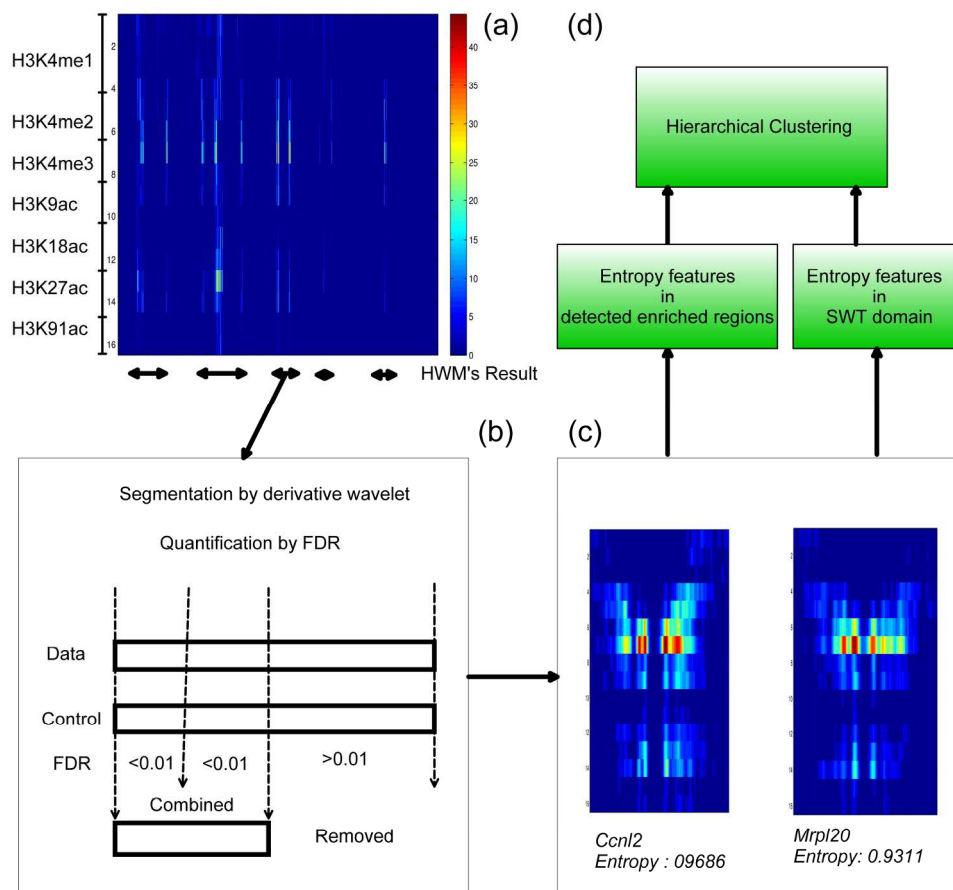


Figure 2. The Dewey procedure. (a) Five domains were identified by HWM. A domain is composed of a set of segments enriched for histone modification signals. (b) Each segment is evaluated using FDR. Adjacent segments are combined to form an enriched region (c) Entropy features are estimated in the obtained enriched regions. Stationary wavelet entropy features are also calculated in SWT domain. (d) Both the entropy and stationary wavelet entropy features are used for clustering.  
191x184mm (300 x 300 DPI)

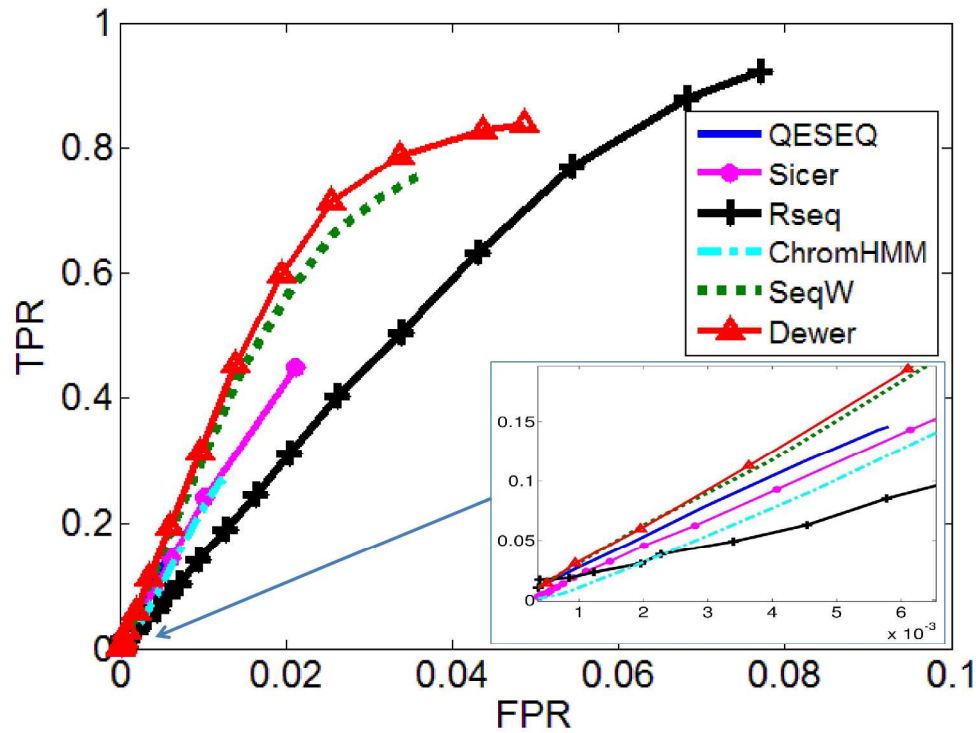


Figure 3. Performance comparison for enriched region detection using 3T3L1 datasets. H3K36me3 data from 3T3L1 cells were used for each detector. For ChromHMM and Dewer, H3K36me3 in four time points during adipogenesis were used. Dewer outperformed other predictors in this test.  
220x167mm (300 x 300 DPI)

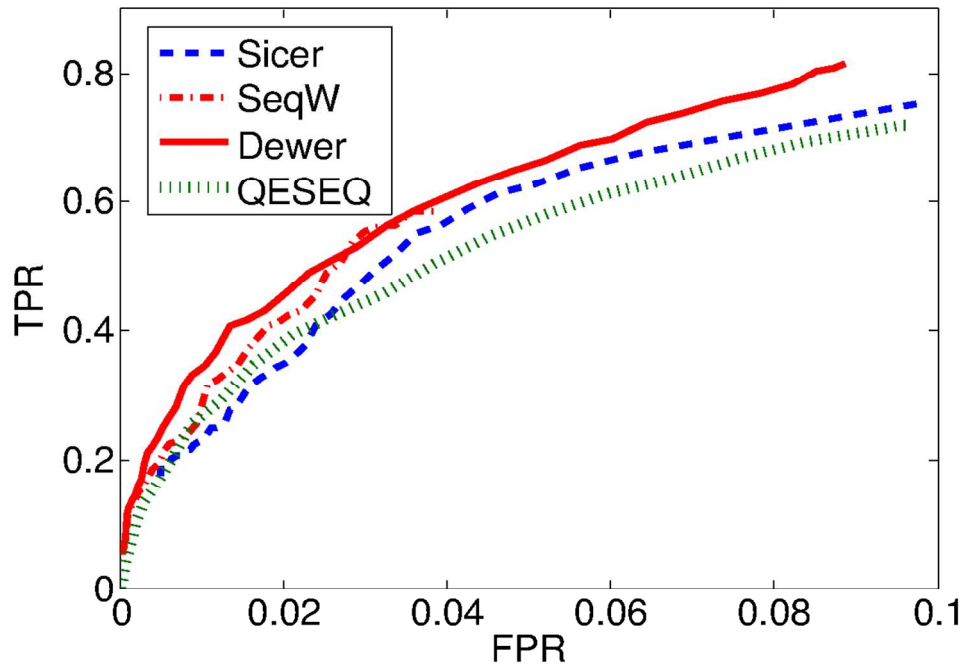


Figure 4. Assessing the performance of detecting enriched region using IMR90 datasets. SeqW performed Sicer and QESEQ when using H3K36me3. Dewer using H3K36me3, H3K79me1 and H3K79me2 performed best.  
167x124mm (204 x 196 DPI)

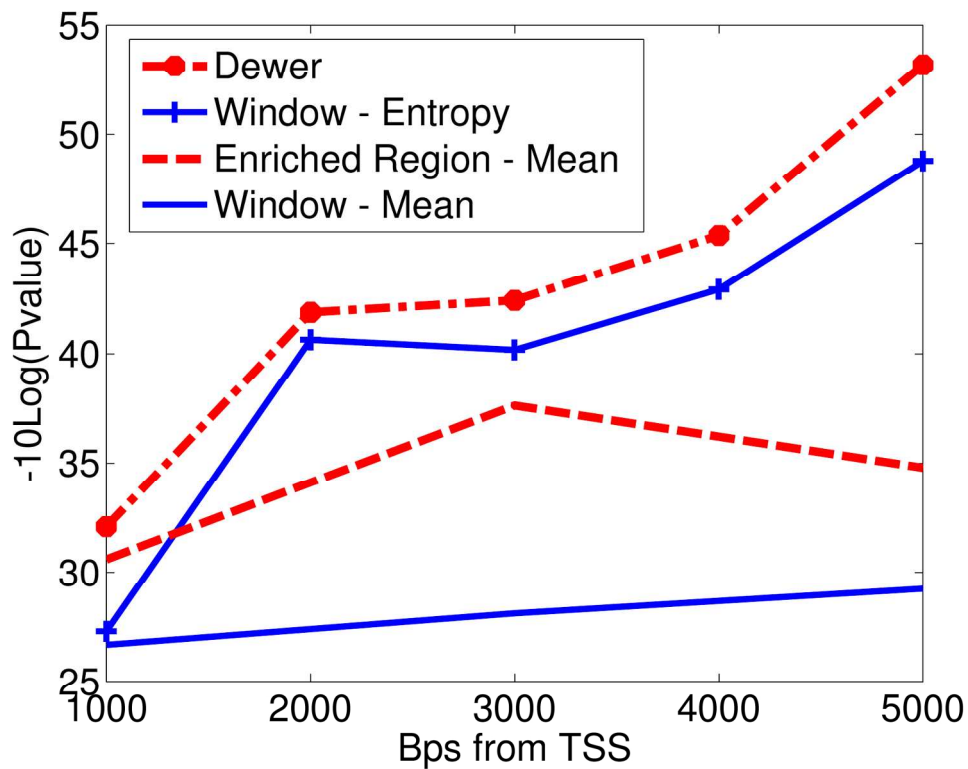


Figure 5. Comparison of the algorithms Dewer used for its clustering. We compared the algorithms of Dewer (entropy and enriched regions). We used entropy in a window and used the mean value in the enriched region. We also compared the performance using the classical mean value in a window for gene clustering. 168x137mm (300 x 300 DPI)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

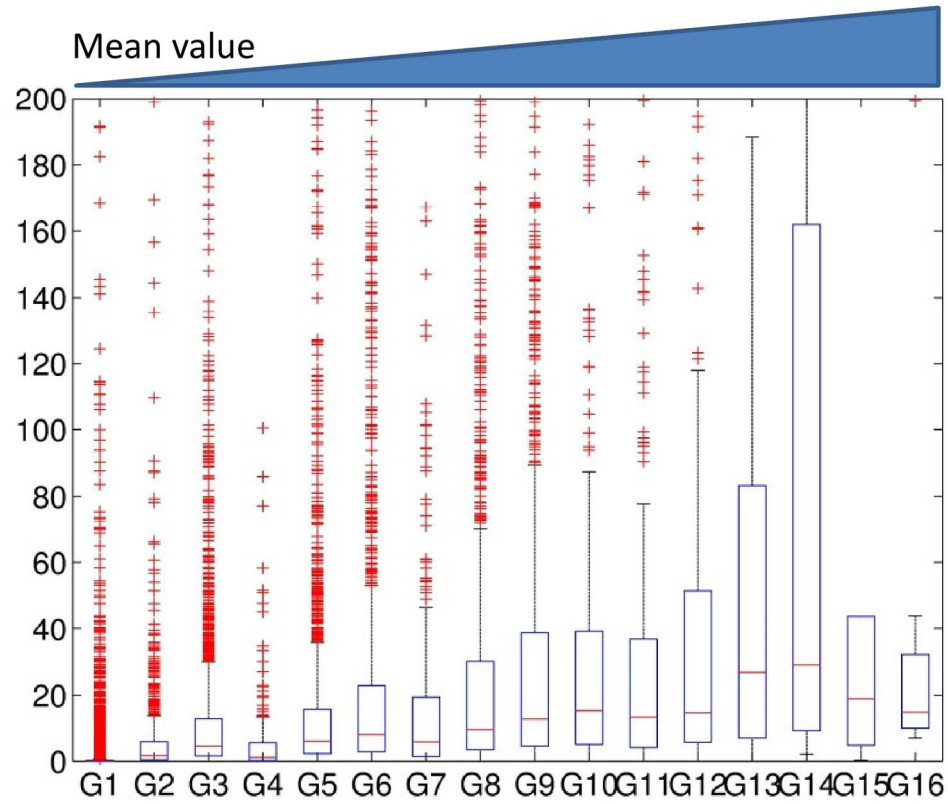
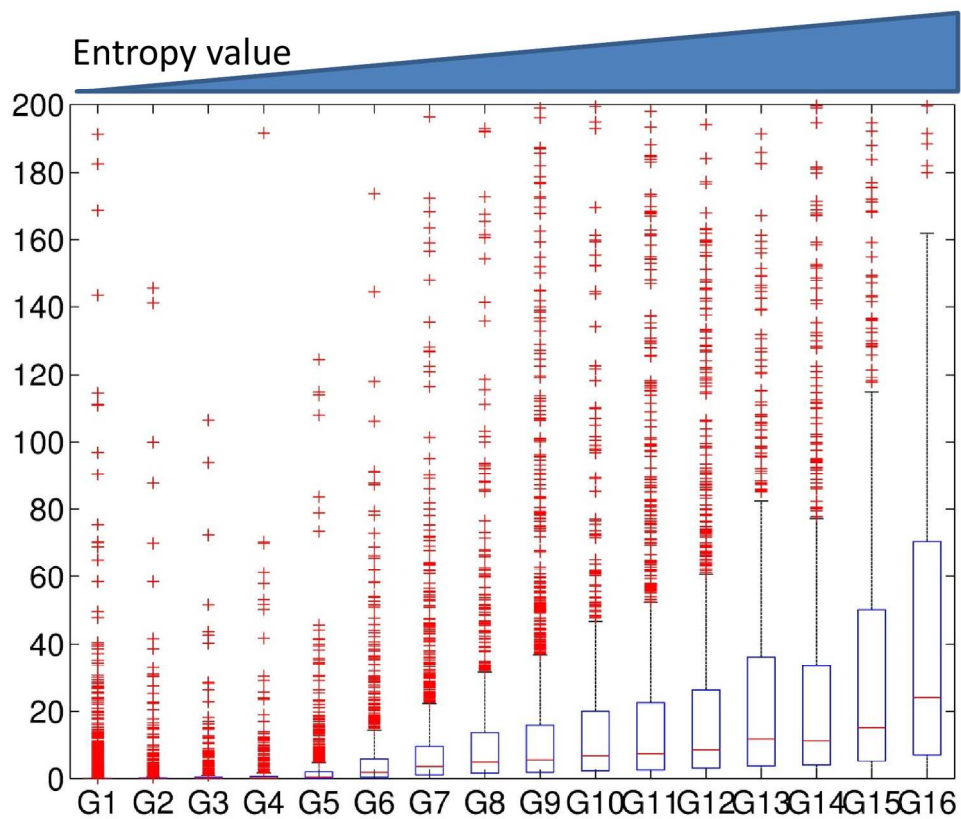


Figure 6. The association of gene expression levels with the mean levels when using window method. We generated 16 clusters were using a window-based approach. Clusters were sorted based on their strength of the features. The expression levels for each cluster were investigated.  
164x138mm (300 x 300 DPI)



36 Figure 7. The association of gene expression levels with the entropy levels when using window method. We  
37 generated 16 clusters were using a window-based approach. Clusters were sorted based on their strength of  
38 the features. The expression levels for each cluster were investigated.  
39 161x139mm (300 x 300 DPI)

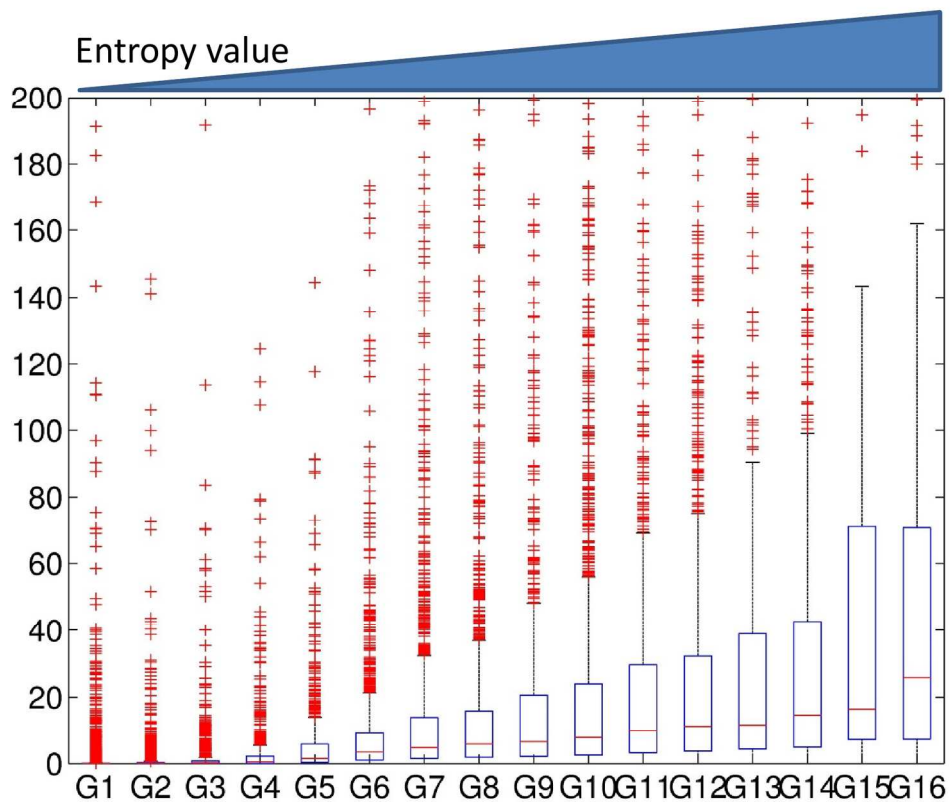


Figure 8. The association of gene expression levels with the entropy levels when using Dewar method. We generated 16 clusters were using Dewar. Clusters were sorted based on their strength of the features. The expression levels for each cluster were investigated.  
165x138mm (300 x 300 DPI)

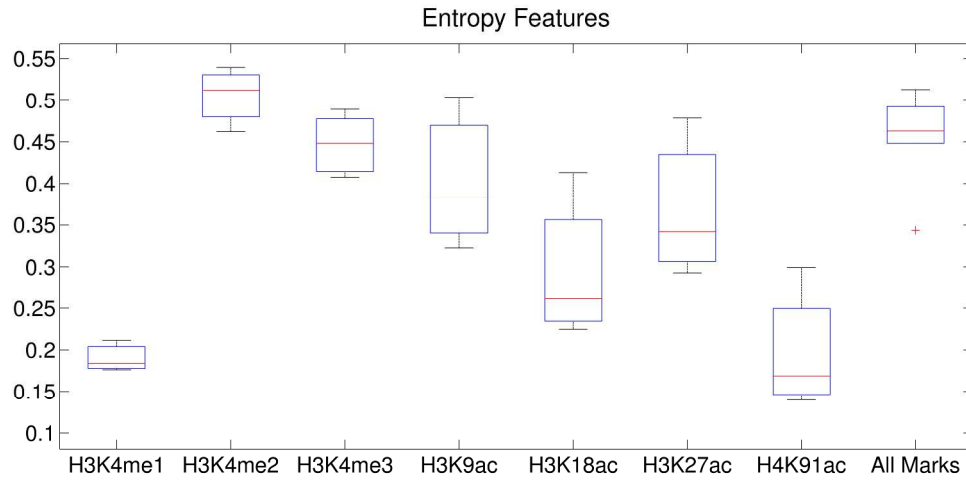


Figure 9. The entropy value of each histone modification mark using Dewey.  
258x128mm (300 x 300 DPI)

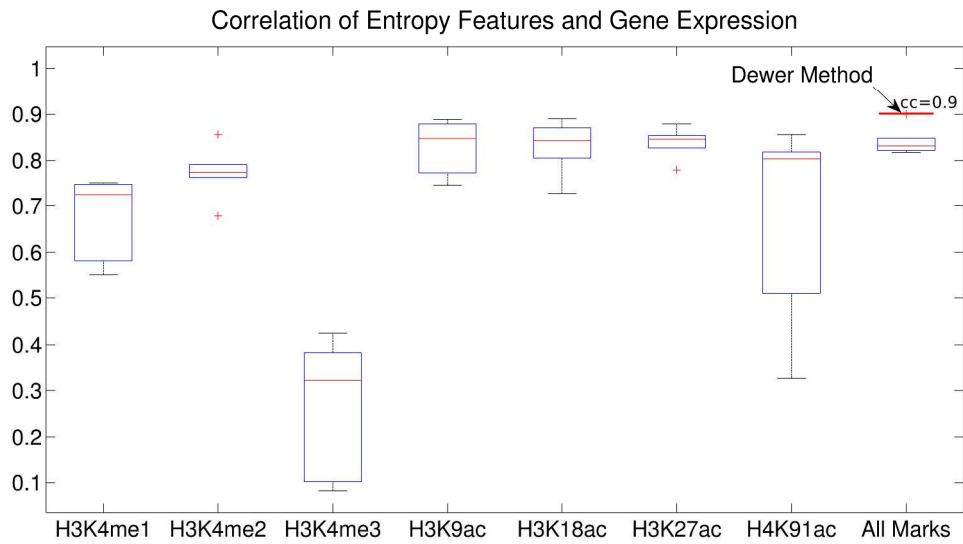


Figure 10. The entropy value and its association for each histone modification mark. CCs between expression levels and the entropy levels for the clusters were calculated when using the epigenetic data 1~5Kbps around the TSSs. The CC of Dewey reached to 0.9.  
258x146mm (300 x 300 DPI)

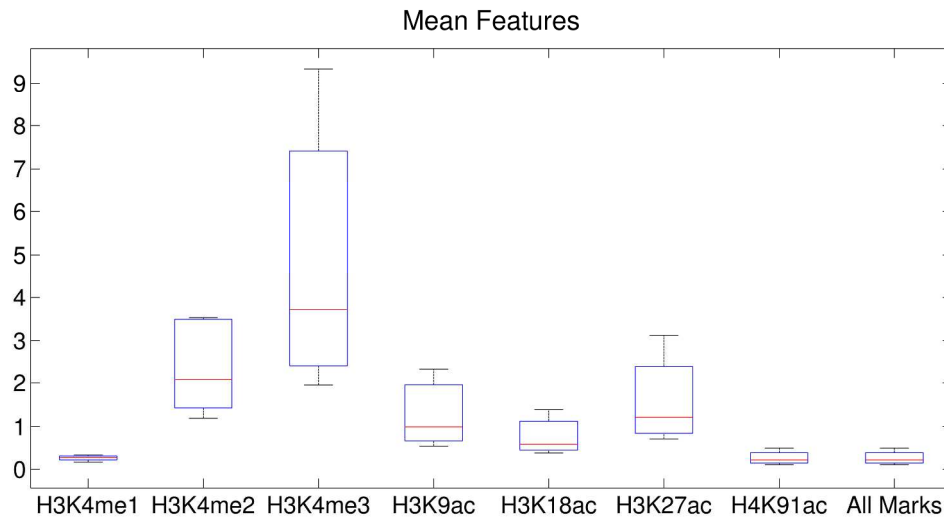


Figure 11. The mean value of each histone modification mark using window method.  
240x130mm (300 x 300 DPI)

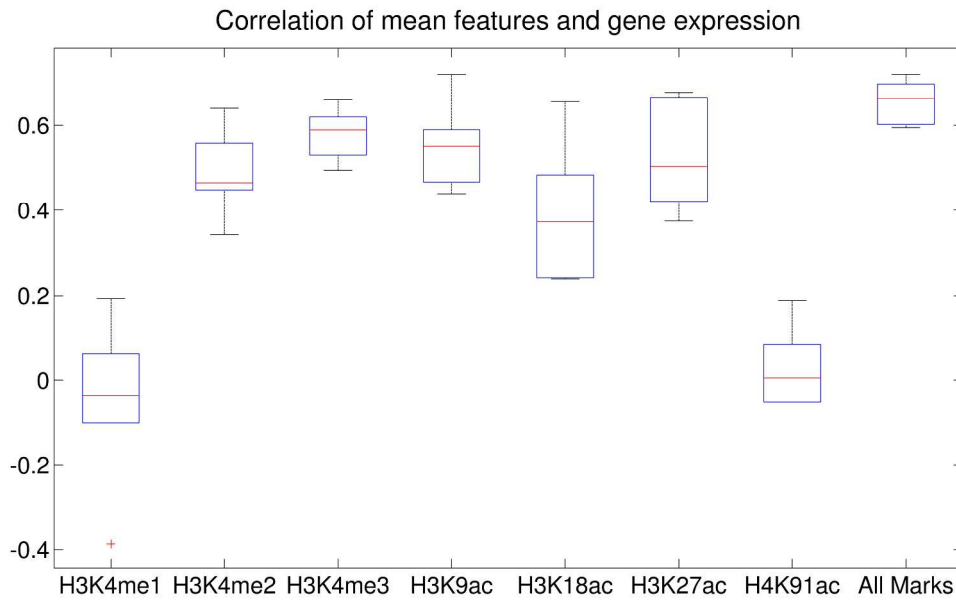


Figure 12. The mean value and its association for each histone modification mark. CCs between expression levels and the mean levels for the clusters were calculated when using the epigenetic data 1~5Kbps around the TSSs.  
241x150mm (300 x 300 DPI)

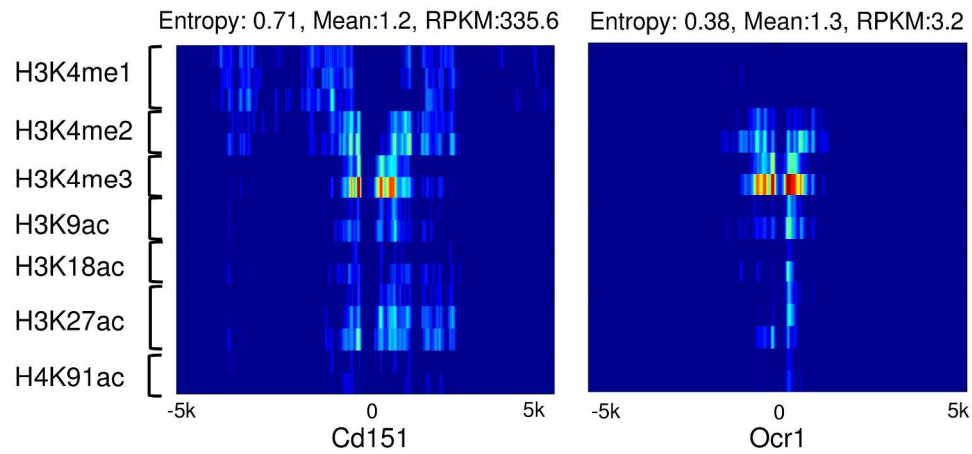


Figure 13. The epigenetic landscapes around the promoter of Cd151 and Ocr1. The entropy, mean value and expression levels (RPKM) are compared in association with epigenetic landscapes.  
238x114mm (300 x 300 DPI)

Table 1: Values of  $-10\log(\text{Pvalue})$  eight SWT scales and many wavelet filters have been calculated. In eight SWT scales, L means low frequency while H means high frequency. 1 and 2 are levels of SWT. Daubechies (db1-db10), Coiflets (coif1- coif5), Symlets (sym2- sym11), Discrete Meyer (dmey), Biorthogonal (bior1.1- bior2.8) and Reverse Biorthogonal (rbio1.1- rbio2.8) are the used wavelet filters.

Scale \ Wavelet	LL1	LH1	HL1	HH1	LL2	LH2	HL2	HH2
db1	36.2	39.6	42.9	45.0	48.1	46.1	37.9	38.0
db2	27.5	34.7	37.8	44.2	33.8	35.9	36.1	41.7
db3	37.4	35.2	40.5	39.1	30.9	33.3	42.6	39.8
db4	38.5	48.7	38.9	29.4	40.7	39.1	32.3	39.0
db5	49.2	47.9	43.4	43.6	36.6	38.1	45.7	40.3
db6	37.6	38.2	35.8	44.5	31.1	26.9	39.6	34.3
db7	38.6	39.3	41.9	40.7	29.1	32.5	35.1	37.1
db8	31.5	36.2	36.2	40.1	28.6	30.6	29.9	36.4
db9	32.5	30.0	42.8	39.5	34.4	37.6	41.3	36.4
db10	34.3	32.7	44.9	43.9	21.7	40.2	45.2	44.6
sym2	27.5	34.7	37.8	44.2	33.8	35.9	36.1	41.7
sym3	37.4	35.2	40.5	39.1	30.9	33.3	42.6	39.8
sym4	45.5	38.9	34.3	35.2	45.8	35.5	43.2	31.6
sym5	32.4	32.3	39.5	43.4	33.1	40.6	42.0	40.1
sym6	35.5	31.8	34.1	36.8	43.4	37.6	43.0	39.4
sym7	47.6	39.3	32.4	41.6	33.2	29.1	34.6	32.1
sym8	32.1	30.0	34.5	41.5	44.5	39.3	48.7	43.9
sym9	45.3	34.0	38.5	34.0	39.6	32.4	44.7	45.5
sym10	30.3	41.2	43.8	40.6	25.7	31.0	36.9	61.5
sym11	30.2	38.9	34.1	42.2	29.5	38.2	42.1	43.8
coif1	35.5	43.0	34.0	33.3	39.1	29.6	35.2	40.2
coif2	32.0	35.0	41.4	41.6	43.2	41.3	43.5	41.4
coif3	35.5	36.6	45.6	35.9	31.3	32.4	43.0	35.3
coif4	43.0	39.1	37.4	51.5	45.4	34.5	36.6	38.8
coif5	38.3	38.6	45.9	34.2	34.6	35.1	39.5	40.6
dmey	45.9	35.7	46.3	33.2	33.5	40.1	36.4	45.8
bior1.1	36.2	39.6	42.9	45.0	48.1	46.1	37.9	38.0
bior1.3	39.6	40.7	43.5	45.0	33.8	35.0	37.5	37.8
bior1.5	43.6	36.9	31.5	45.0	41.9	39.7	50.4	40.0
bior2.2	41.7	41.6	37.5	42.0	40.3	36.3	52.1	38.5
bior2.4	36.1	43.6	28.9	42.0	40.4	33.3	38.0	38.3
bior2.6	34.0	32.7	43.4	42.0	39.4	33.9	51.3	43.5
bior2.8	37.2	35.4	43.5	42.0	34.3	36.0	35.3	43.6
rbio1.1	36.2	39.6	42.9	45.0	48.1	46.1	37.9	38.0

rbio1.3	36.2	38.4	37.9	36.3	48.1	31.3	38.7	44.3
rbio1.5	36.2	35.6	41.9	35.6	48.1	35.4	34.5	50.9
rbio2.2	39.4	45.1	41.9	37.3	28.2	27.2	36.3	36.8
rbio2.4	39.4	36.7	39.3	36.2	28.2	48.4	38.5	42.4
rbio2.6	39.4	39.6	42.0	36.7	28.2	35.3	31.9	32.8
rbio2.8	39.4	36.0	36.8	44.0	28.2	25.9	45.8	32.7