

MUSIC: Identification of Enriched Regions in ChIP-Seq Experiments using Mapability Correction and Multiscale Signal Processing Framework

[CITATIONS ARE NOT ADDED, YET]

ABSTRACT:

We present MUSIC, a signal processing method for identification of enriched regions in the genome-wide read depth (RD) signal profiles from ChIP-seq experiments. The basic motivation behind MUSIC is twofolds: First, systematic noise introduced by non-uniform read mapability makes it challenging to process the enrichments. Second, many ChIP-Seq assays have a broad spectrum of enrichments, e.g. H3k36me3 marks the active gene bodies whose lengths range from 100 bps to megabases, that makes it necessary to analyze the signal at multiple scales that can be tuned by the user. Motivated by these, MUSIC first utilizes a smoothing and enhancement procedure to correct for non-uniform read mapability. MUSIC then performs median filtering based smoothing to the corrected signal using multiple windows of increasing lengths, i.e. multiple scales. Essentially, at each smoothing scale, the details of the RD signal is smoothed and the candidate enriched regions are revealed as "blobs". The candidate enriched regions are filtered with respect to significance compared to a normalized control signal. The significantly enriched regions from all the scales are merged to generate the final set of enriched regions. We compare the accuracy to other ER identification methods, and show that MUSIC performs favorably in terms of accuracy and reproducibility of ERs compared to other methods. Next, we analyze Polymerase II binding ChIP-Seq data with MUSIC. We show that there is a clear distinction between the expression levels of genes with punctate bound (stalled) and broadly bound (elongating) polymerase at the promoters. This observation indicates that the multiscale analysis can be utilized for evaluating the length scale of enrichment for genomic signals and can reveal important biological information. MUSIC is available for download from Gerstein Lab Github page at: <https://github.com/gersteinlab/MUSIC>

INTRODUCTION:

With the recent advancements in sequencing technologies, chromatin immunoprecipitation based enrichment of the DNA sequences followed by sequencing (ChIP-seq) has become the mainstream experimental method for genome-wide measurement of DNA binding proteins (e.g. transcription factors) and posttranslational modifications of histone proteins, or histone modifications (HMs). Following the sequencing, it is necessary to computationally process the read depth signal profile to analyze the enrichments. Depending on the target of the ChIP-seq assay, the length scale of enrichments can be very different for different experiments. For example, for transcription factor binding, the enrichments are concentrated to a very punctate region of hundreds of nucleotides (cite{Rozowsky,Kharchenko, Fejes,Jothi, MACS reference}). For most HMs, broad enrichments are observed. HMs like H3k9me3 show enrichments that can extend upto megabases. Another interesting example is RNA polymerase II, which binds to the promoters and gene bodies for the purpose of mRNA transcription, whose enrichments can extend over the whole gene bodies or can be punctate and concentrated close to gene promoters. Identification and characterization of the broad enrichments is

Deleted: (Alternative: MUSEQ: Enrichment Analysis from ChIP-Seq Experiments using Mapability Correction and Multiscale Signal Processing Framework?)¶

Formatted: Font: Bold, Italic

Deleted: Chromatin immunoprecipitation followed by sequencing (ChIP-Seq) has become the mainstream experimental methodology for probing genome-wide enrichment of DNA binding proteins and post translational modification of histones, or histone modifications (HMs). Following a ChIP-Seq experiment, the generated reads are mapped and it is necessary to computationally identify the enrichments in the read depth signal. Unlike the transcription factors that show punctate signal enrichments, the enrichments for many ChIP-Seq assays manifest at much larger length scales (e.g., histone modifications and Polymerase II binding) and the spectrum of length scales can vary significantly between different types of ChIP-Seq experiments. This, combined with the systematic noise added to the signal by the non-uniform mapability makes it challenging to process the enrichments in the HM signal profiles.¶ In this paper, we

Deleted: novel

Deleted: identifying and processing enrichments

Deleted: uses a mapability correction procedure coupled with multiscale signal processing bas (... [1])

Deleted: we identify the enriched regions (... [2])

Deleted: we

Deleted: of ERs

Deleted: the

Deleted: to a ground truth

Deleted: ¶ (... [3])

Deleted: for Polymerase II binding. The bro (... [4])

Moved down [1]: state, i.e., stalled or (... [5])

Deleted: We process the SEF pileup signal (... [6])

Deleted: for the protein coding

Deleted: .

Deleted: points to the fact

Deleted: SEF pileup signal can be utilized for

Deleted: of broadness of

Deleted: that

Deleted: ¶

Deleted: There has been a lot of work

Deleted: identification

Deleted: sites, or peaks, from ChIP-Seq (... [7])

Deleted: Identification of the broad enrich (... [8])

Deleted: for most of the HMs signal profiles

the basic first step for understanding the regulatory effects of the HMs and diffuse DNA binding proteins on gene expression as more evidence is brought to light that these epigenetic factors are major driving factors for disease manifestation like cancer.

Comment [O1]: Citations on epigenetics and disease

Several popular methods for identification of broad enrichments include change point identification within the formality of Bayesian inference (BCP), local island identification and clustering (SICER), local thresholding and merging (MACS), and using local Poisson statistics to identify broad enrichments (SPP), Wavelet based smoothing and identification of enriched regions (WaveSeq, Kharpikov et al). In the ENCODE project, the main concentration has been on building integrative unsupervised segmentation models $\{ \text{SegWay} \}$, $\{ \text{ChromHMM} \}$ for using these modifications to annotate and characterize the cis-regulatory elements in the genome and studying the relation between the HM levels at these elements and the gene expression. Although the segmentation methods proved very useful for identifying novel regulatory elements like enhancers, the identification of the enriched regions per HM signal profile and the technical issues like mapability have not been adequately addressed. Finally, some of the HMs are almost never processed in the segmentation methods process simply because they are not well characterized in terms of their function and in terms of how they interact with other HMs.

Formatted: Font color: Background 1

There are two main challenges for identification of broad enrichments in the ChIP-Seq read depth signal profiles. First is that unlike transcription factor binding, HMs enrichments are observed at much larger length scales and the spectrum of enrichment lengths are different for different types of ChIP-Seq experiments. This makes it necessary to identify the enrichments at different scales. A widely used method for identifying the HM signal profiles is smoothing the signal profile with a kernel of constant size and shape and using a null model (e.g., Poisson or negative binomial based) to identify the significantly enriched regions. It is, however, not clear how the kernel size and shape should be selected. The multiscale approaches proposed by the wavelet based methods address this aspect but in those approaches, the selection of the predefined wavelet functions are not justified for their choice. Second, the signal profiles contain systematic noise introduced to the read depth signal by the repeat regions with low mapability. This noise causes discontinuities in the identified enrichments. This becomes an important factor especially in the intergenic regions where a large region of enrichment, which may be a single element like a long repressed region, would get broken into many smaller regions.

Deleted: HM

Deleted: broad

Deleted: naïve and

Deleted: obvious

Formatted: Font: Bold, Italic

Comment [O2]: Wavelet based methods are different from ours as they do not utilize the the natural merging of clusters generated by multiscale smoothing approaches

Deleted: there is non-uniform mapability.

Deleted: We utilize a multiscale signal processing approach based on a novel multiscale median filtering decomposition to identify the regions of enrichment in the signal profile generated.

Deleted: an efficient

Deleted: before multiscale decomposition of the signal

Deleted: are based on

Deleted: edge preserving

Deleted: , which has not been applied to processing ChIP-Seq datasets before

Deleted: multiscale

Deleted: procedure

Deleted: features

Deleted: These

Deleted: features

Deleted: .

Deleted: significant enrichment features (SEFs),

Deleted: output of MUSIC algorithm. By pileup of

In this paper, we present MUSIC, a method to identify enriched regions (ERs) in ChIP-Seq experiments. MUSIC uses mapability correction at the nucleotide resolution so as to correct for the spurious loss of signal because of low mapability. Next, MUSIC performs multiscale decomposition of the corrected RD signal. Unlike the wavelet based multiscale approaches that use linear filtering, we take an approach to multiscale decomposition using the non-linear median filtering. Basically, MUSIC exploits the fact that the smoothing with a certain window length removes the small details in the signal (like small peaks and small valleys) and reveals the enriched regions in the signal as “blobs” that can be detected as the regions between consecutive local minima of the smoothed signal.

The enriched regions at each scale are then trimmed and filtered with respect to their significance compared to a normalized control. This procedure yields the scale specific enriched regions (SSERs). Each SSER corresponds to an enriched region at a certain scale. At smaller scales, the SSERs correspond

to the enrichments that are more punctate compared to SSERs at higher scales, which represent the broader enrichments. To identify the ERs, MUSIC merges the SSERs from all the scales.

First, to evaluate the accuracy of ERs, we concentrate on H3k36me3, a well characterized HM that gets enriched on expressed gene bodies, which we use as gold standard. We compare the accuracy of ERs with 5 other methods with respect to accuracy, in terms of consistency with expressed regions, and reproducibility. We show that ERs identified by MUSIC have higher F-measure and higher reproducibility compared to other methods.

Next, we analyze the Polymerase II ChIP-Seq dataset using MUSIC and we demonstrate that the genes with less broad polymerase binding at their promoters have significantly lower expression than the genes that are bound with more broad polymerase at the promoters. We conclude that the binding scale of polymerase at the gene promoters as identified by MUSIC is indicative of its state, i.e., stalled or elongating.

We demonstrate that the genes whose promoters are bound by the small scale polymerase binding, we show that there is significantly lower levels of gene expression compared to .

The broadness of the enrichment of Polymerase ChIP-Seq signal profiles around gene promoters can be used to distinguish the stalled and elongating polymerase. The stalled polymerase shows a punctate binding, whereas the elongating polymerase shows a much more broad binding. In addition, the stalled polymerase does not transcribe the gene, that can be observed as the low or no detectable expression of the downstream gene. For assessing the broadness of binding of polymerase, we computed the scale of enrichment of Polymerase II (using SSERs) around the protein coding gene promoters. We demonstrate that at small scales of polymerase enrichment, the gene expression is significantly smaller compared to the larger scales of enrichment.

RESULTS:

We first present motivation for MUSIC algorithm and lay out the steps of the algorithm. Then we present comparison of MUSIC with other ER identification algorithms. We finally present the joint processing of the polymerase data with gene expression levels.

MUSIC ALGORITHM:

There are two factors that motivate the novel methodology behind MUSIC:

1. Mapability is an important aspect of read mapping and processing. For example, in the repetitive regions the number of uniquely mapable positions decreases significantly. This, depending on the parameters of the mapping algorithm, causes a systematic decrease of signal at repetitive regions and makes it impossible to evaluate whether a decrease in the signal is due to low mapability or a decrease in the modification levels. This becomes problematic especially in the intergenic and intronic regions which contain many repetitive regions. Consequently, the broadly enriched intergenic and intronic regions will be fragmented into many smaller enriched regions. It is worth noting that this problem is less severe for the punctate enrichments like transcription factor binding.

Deleted: SEFs, MUSIC generates a genome-wide signal, which quantifies the broadness of enrichment at each position. This enables us to study the broadness of signal

Deleted: different ChIP-Seq experiment. We show additional utility of SEFs with two applications.

Deleted: SEFs for identification of enriched regions (ERs),

Deleted: ground truth. We build ERs from SEFs and

Deleted: a novel utility of SEF pileup signal for analysis of enrichment of RNA

Deleted: ||

Moved (insertion) [1]

Deleted: .

Formatted: Font color: Background 1

Deleted: close to

Formatted: Font color: Background 1

Deleted: generated

Formatted: Font color: Background 1

Deleted: SEF pileup signal for

Formatted: Font color: Background 1

Formatted: Font color: Background 1

Deleted: ChIP-Seq data from ENCODE project.

Formatted: Font color: Background 1

Formatted: Font color: Background 1

Deleted: joint distribution of broadness

Formatted: Font color: Background 1

Deleted: at promoter, as quantified by SEF pileup signal, versus gene

Formatted: Font color: Background 1

Deleted: shows a bimodal characteristic, where one mode corresponds

Formatted: Font color: Background 1

Deleted: stalled polymerases and other mode corresponds to the elongating polymerases. This showcases a novel application of SEFs and a novel benefit

Formatted: Font color: Background 1

Deleted: MUSIC

Formatted: Font color: Background 1

Deleted: Remaining of the paper is as follows. We first describe the MUSIC algorithm and identification of SEFs. Next we present identification of ERs and compare the accuracy and reproducibility of MUSIC other ER identification methods. Then we fo... [9]

Deleted: a novel application of the enrich... [10]

Deleted: tool

[[INTRODUCE THIS AS IMPULSE NOISE]]

In order to characterize the mapability of different regions, MUSIC generates the genome-wide multi-mapability signal profile. For each position, this profile contains the number of reads (of certain length) that can map from any other position in the genome. In order to gain a perspective on the statistics of multi-mapability signal, we aggregated the signal over different elements. This reveals, as expected, that the protein coding exons and promoter regions show the highest mapability (See Figure S1). The multi-mapability signal is utilized by MUSIC in correction of effects of mapability.

2. The length distribution of ERs for broad enrichments usually have a large variance. This makes it necessary to identify the enrichments for a spectrum of scales. For example, for HMs like H3k36me3, H3k27me3, the ERs can extend from several kilobases to hundreds of kilobases. On the other hand, for HMs like H3k4me3 and H3k27ac, which marks the gene promoters and enhancers, the ERs are around kilobases in length. Another interesting example is the RNA Polymerase II, whose enrichments can extend from less than a kilobase to hundreds of kilobases.

Motivated by these facts, we designed MUSIC to account for the effects of mapability and to be scale sensitive. In essence, MUSIC first corrects the RD signal from ChIP experiment for the mapability. MUSIC then computes the multiscale decomposition of the signal by smoothing the signal with multiple smoothing window lengths. In the process of smoothing, fine details in the signal are removed and the broad enrichments are revealed as “blobs” in the smoothed signal, which are detected as the regions between consecutive local minima of the smoothed signal.

[[ASPECT OF MERGING OF BLOBS: Technical citations]]

These regions are then filtered with respect to significance computed in comparison with the control signal to generate the scale specific enriched regions, SSERs. The SSERs at small scales represent the small enrichments in the signal and the vice versa for SSERs at large scales. With multiple scales, MUSIC can detect SSERs within a spectrum of lengths that can be tuned by adjusting the starting and ending scale levels to be processed by MUSIC.

Figure XX shows the flowchart of MUSIC (See Methods for more details.) Here we summarize each step briefly. The input to MUSIC are the sets of reads from the ChIP and control samples (Steps 1 and 2), and the set of smoothing window lengths to be used in multiscale analysis. MUSIC first preprocesses the reads and filters the PCR duplicates for both samples. Then MUSIC computes a scaling factor using linear regression between the ChIP and control signal profiles. The slope of the regression is used as a normalization factor for control.

Then, in Step 3, the ChIP and normalized control signal profiles are generated, and the ChIP profile is smoothed and corrected with respect to mapability using the multi-mapability profile. The correction can be formulated as following:

Deleted: This signal profile is similar to other mapability maps computed previously.

Deleted: that

Deleted: smallest

Deleted:) with (Figure SXX), with average

Deleted: 1.2 reads mapping at these positions on average. We computed multi-

Deleted: profiles for several read lengths and they are available for download with MUSIC (See Methods).

Comment [O3]: cite

Deleted: It is worth noting that although there are computational methods that aim at assigning the reads uniquely to repeat regions by resolving the multi-mapping reads, the underlying algorithms are too compute intensive for the purpose of enrichment identification, where the computation power should be allocated for identification of the enriched features.¶

Deleted: couple of

Deleted: couple

Deleted: profile

Deleted: at

Deleted: scale levels, then uses the decomposition to compute the significant enrichment features.

Deleted: can be

Deleted: enriched features

Deleted: [[Each scale level is represented with smoothing window length where higher scales are associated with longer smoothing windows and thus are associated with broader features.]]¶

Deleted: . The

Deleted: of each step can be found in Methods.

Deleted: $\hat{x}_i =$
Compare the signal value at i with
the median signal at highly mapable positions
 $\max\{x_i, \text{median}(\{x_a\}_{a \in [i-l/2, i+l/2]} \mid m_a < \bar{m}_{\text{exonic}})\}$ ¶
Median of the signal values at highly mapable
positions around i
Where

Compare the signal value at i with
the median signal at highly mapable positions

$$\tilde{x}_i = \max[x_i, \text{median}(\{x_a\}_{a \in [i-l_c/2, i+l_c/2] \mid m_a < \overline{m}_{\text{exonic}}})]$$

Median of the signal values at highly mapable
positions around i

where x_i and \tilde{x}_i are the uncorrected and corrected signal values, respectively, at position i , m_a is the value of multi-mapability profile at position a , l_c is the length of median filter utilized in correction which is by default set to 2000 base pairs, and $\overline{m}_{\text{exonic}}$ is the average multi-mapability signal value over the exonic regions, which we identified as the most mapable regions in the genome (See Fig S1). In summary, MUSIC first generates the mapability corrected signal profile from the ChIP-seq signal (See Methods), where for each position i , MUSIC computes the median of the signal values at highly mapable positions (multi-mapability signal smaller than 1.2) within l_c vicinity of i . Then MUSIC compares this value with the signal value at i and assigns the maximum of them to the corrected value. The basic idea behind this correction is that since we know that mapability causes loss of signal, if the signal value at i is higher than its vicinity, then it is highly likely that the mapability did not affect the signal value at i . It should be noted that maximum filtering, also known as dilation in image processing, is used for feature enhancement in images.

MUSIC then performs median filtering to the mapability corrected ChIP profile to compute multiscale decomposition of ChIP signal at multiple scales (Step 4.) For this, MUSIC uses window lengths beginning with l_{start} and ending at l_{end} . The window length is increased multiplicatively, thus, the window lengths that are used in multiscale filtering follow a geometric series:

$$\{l_{\text{start}}, l_{\text{start}} \times \sigma, l_{\text{start}} \times \sigma^2, \dots, l_{\text{end}}\}$$

where σ is the multiplicative factor between consecutive window lengths. This parameter is set to 1.5 by default.

Compared to the kernel based linear filters (which are also used in the wavelet based multiscale decompositions), median filtering has two advantages. First, at low noise levels, median smoothing preserves the edges, i.e. sharpness of increase and decrease of the RD signal at the ends of enriched regions, in the signal better than the linear filters. Secondly, the median smoothing is more tolerant to the burst or impulse noise compared to the linear filters. This is important for the enriched region identification since the systematic noise added by multi-mapability can be viewed as an impulse noise.

For each smoothed signal at each scale, MUSIC identifies all the local extrema, i.e., local minima and local maxima (Step 4 in Fig. 1). The regions between the consecutive local minima are marked as the candidate enriched regions. For each enriched region, MUSIC computes the fraction of the maximum of smoothed RD signal (at the corresponding scale) to the maximum of the unsmoothed ChIP signal within the boundaries of the enriched region. If this fraction is smaller than the smoothed versus unsmoothed signal ratio threshold (denoted by γ), MUSIC discards this enrichment feature (Refer to Methods.) This way, MUSIC avoids overmerging of the enriched regions.

Deleted: l

Deleted: identified

Deleted: .

Deleted: first

Deleted: 2000 bp

Deleted: For a given scale s with smoothing window length l_s , the median smoothed signal is formulated as

Deleted: x_i^s

Deleted: filtered value at position i . The

Deleted: is utilized extensively in signal processing as an edge preserving

Deleted: filter. MUSIC utilizes

Deleted: filtering in a novel application in multiscale decomposition.

Deleted: decomposition

Deleted: enrichment features. It is worth noting that these features have exactly one local maxima in them.

Deleted: enrichment feature

Deleted: ChIP

Deleted: enrichment feature.

Deleted: maximum

Deleted: smoothing

Deleted: removes, at large scales,

Deleted: features with local enrichment

The regions identified from the consecutive minima are rough and it is necessary to identify the location of densest signal enrichment within each region. To achieve this, MUSIC performs a Poisson background based thresholding and p-value minimization to trim the ends and identifies the densest, or most compact, regions of signal enrichment in the enrichment feature. Step 5 in Fig. 1 illustrates the trimmed ends of the candidate enriched regions. Finally, MUSIC computes the binomial p-value for each trimmed enrichment feature and filters out those whose p-values are larger than 0.05. We refer to the remaining regions as the scale specific enriched regions (SSERs).

[[SSERs at different scales can overlap, however the SSERs are not necessarily transferred all the way to the highest scale of processing.]]

The basic assumption is that SSERs contain all the information about the enrichments in the signal over a spectrum of smoothing scales. MUSIC utilizes the SSERs for processing the enrichments in the signal.

Identification of ERs

MUSIC utilizes SSERs to identify enriched regions in the genome. The ERs, unlike SSERs, is a set of non-overlapping regions that are enriched. For this, the candidate ERs are computed by merging the union of the SEFs identified from all the decompositions (Step 6 in Fig. 1). MUSIC then evaluates the quality of the ERs with respect to concordance of the signal levels on positive and negative strands. MUSIC computes the amount of signal mapping to positive and negative strand in each ER and filters out the ERs for which the counts of reads that map to positive and negative strand within a factor of 2 of each other (See Methods.)

For each of the remaining ERs, MUSIC computes a binomial p-value using the number of reads in the ChIP and control samples. Since the features have different lengths, the all the counts are normalized to l_{pval} window length (See Methods):

$$p(i, j) = B\left(\frac{n_{chip}}{(j - i + 1)} \times l_{pval}, \frac{n_{control}}{(j - i + 1)} \times l_{pval}\right)$$

Where n_{chip} and $n_{control}$ are the read counts in ChIP and control samples within the ER starting at nucleotide position i and ending at j . The multiple hypothesis correction is performed by the Benjamini-Hochberg procedure [cite{XXX}]. The q-values computed from the correction are thresholded with respect to 0.05 for identification of the significant ERs.

Evaluation of the Enrichment Broadness

The scale dependence of SSERs is a useful property for evaluating the “broadness” of enrichments. Each SSER represents a locally enriched region at a certain scale. Therefore, a position that is covered by many SSERs has a broader enrichment than a position that is covered by small number of SSERs. Following this basic observation, MUSIC generates the SSERs pileup signal profile, a genome-wide profile that is generated by counting the number of SSERs covering each position, which quantifies the broadness of enrichment at each position in the genome.

Deleted: features

Deleted: feature

Deleted: (

Deleted: XX). For a feature starting at position i and ending at position j the trimming operation identifies new start and end coordinates, i' and j' , can be formulated as following: $\mathbb{1}(i', j') = \underset{i < a < b < j}{\operatorname{argmin}}_{a, b} (p(a, b) \mid \exists a' \in (a, b) \text{ s.t. } x_{a'} > \tau)$ $\mathbb{1}$ where τ is the threshold identified from Poisson model, and $p(a, b)$ represents the binomial p-value of the region within coordinates a and b , and i' and j' represent the start and end of

Deleted: feature,

Moved down [2]: respectively.

Deleted: (Refer to Methods).

Deleted: enrichment features as significant enrichment features (SEFs)

Deleted: SEFs

Deleted: local

Deleted: , therefore

Deleted: SEFs

Deleted: Next, we will present the utility of SEFs with two applications.

Moved (insertion) [3]

Moved (insertion) [4]

Formatted: Font color: Background 1

Moved (insertion) [5]

Formatted: Font color: Background 1

Deleted: with SEF Pileup Signal

Deleted: SEFs

Deleted: SEF

Deleted: smoothing window length.

Deleted: SEFs

Deleted: SEFs

Deleted: SEF

Deleted: SEFs

Deleted: also

Formatted: Font: Not Bold, Not Italic

To illustrate this, we processed multiple ChIP-Seq datasets (CTCF, Polymerase 2, and several HMs) from ENCODE project for K562 cell line from with MUSIC with smoothing window lengths starting from 100 bp to 2.5 megabase with $\sigma = 1.5$ (Total of 25 scales) and generated the SSER pileup signal for chromosome 1. Figure 2 shows the distribution of SEF pileup signal for different datasets. In this plot, we mapped the value of SSER pileup signal to its corresponding smoothing window length which is also shown in the x-axis. As expected, CTCF, a punctate binding transcription factor, shows the least broad enrichments compared to other datasets. H3k4me3 and H3k4me1, active promoter and enhancer HM marks, show broader enrichments than CTCF. H3k36me3 and H3k27me3, which mark active and repressed gene bodies, show broader enrichments and finally H3k9me3, an HM associated with large heterochromatin domains, shows the broadest enrichments. Another interesting observation is that H3k4me3, H3k4me1, and H3k36me3 have maxima at certain scales, which indicates that these HMs get enriched at specific length scales that are observed very frequently. Finally RNA Polymerase II signal profiles show a high frequency of enrichments at small scales that shows more gradual decrease in frequency as the scale increase.

Comparison with Other ER Identification Methods:

In order to evaluate the accuracy of the enriched, we compared the ERs from MUSIC with 5 other algorithms that identify ERs from ChIP-Seq data: BCP, SPP, MACS, SICER, and PeakRanger. We ran all the algorithms using H3k36me3, and H3k27me3 ChIP-Seq datasets for GM12878 and K562 cell lines from ENCODE project. H3k36me3 correlates well with expressed transcript regions and this allows us to build a gold standard set for H3k36me3 as the bodies of expressed transcripts. We downloaded the transcript quantifications (in RPKMs) from ENCODE RNA-seq dashboard and thresholded the expression levels of the transcripts and filtered the transcripts with low expression. The expressed transcripts are then merged to generate the final set of expressed regions. Rather than selecting one expression threshold, we selected thresholds between 0 and 1 units of RPKM increasing with steps of 0.01 so as to evaluate the accuracy of peak calls against multiple gold standard sets identified at different levels of expression.

[[Parameter selection for this comparison]]

Accuracy Measures:

To measure the accuracy of ERs, we used sensitivity (the fraction of the coverage of correctly predicted ERs to the coverage of the gold standard set) and positive predictive value (the fraction of the coverage of correctly predicted ERs to the coverage of identified ERs). In order to combine the sensitivity and PPV into one accuracy measurement, we used F-measure, which is the harmonic mean of sensitivity and positive predictive value (See Methods). Having one measure of accuracy enables us to easily compare the accuracy of methods with changing RPKM thresholds.

Figure 3a and b shows the F-measure of the H3k36me3 peak calls for different methods with respect to the changing RPKM cutoffs used to identify expressed regions. MUSIC has higher F-measure than all the other methods for GM12878 at all expression cutoffs, followed by BCP. For K562, MUSIC has higher F-measure than all other methods for expression cutoffs smaller than 0.8 then falls slightly below BCP. It

Deleted: for

Deleted: scales

Deleted: scale

Deleted: scale with exponentially increasing smoothing window lengths of 1.5

Deleted: SEF

Deleted: SEF

Deleted: and relatively high

Deleted: of enrichments at large scales

Deleted: ¶

Identification of ERs Using SEFs¶
MUSIC utilizes SEFs to identify enriched regions in the genome. The ERs, unlike SEFs

Moved up [3]: , is a set of non-overlapping regions that are enriched. For this, the candidate ERs are computed by merging the union of the SEFs identified from all the decompositions (Step 6 in Fig. 1). MUSIC then evaluates the quality of the ERs with respect to concordance of the signal levels on positive and negative strands. MUSIC computes the amount of signal mapping to positive and negative strand in each ER and filters out the ERs for which the

Deleted: fraction of total signal on positive strand to that on negative strand (or vice versa) is less than 0.5. We observed that this filter removes many spurious ERs for the HMs with relatively less broad enrichments (See Methods).¶

Moved up [4]: For each of the remaining ERs, MUSIC computes a binomial p-value using the number of reads in the ChIP and control samples. Since the features have different lengths, the all the counts are normalized to l_{pval} window length (See Methods).¶

Formatted: Font color: Background 1

Deleted: $p(i, j) = Bin$

Moved up [5]: $(\frac{n_{chip}}{j-i+1} \times l_{pval} \times \frac{n_{control}}{j-i+1} \times l_{pval})$ ¶
Where n_{chip} and $n_{control}$ are the read counts in ChIP and control samples within the ER starting at nucleotide position i and ending at j .

Formatted: Font color: Background 1

Deleted: The p-values are corrected by Benjamini-Hochberg procedure(cite{XXX}). The final corrected p-values are thresholded with respect to 0.05 for identification of significant ERs. MUSIC can be utilized to determine ERs from precompute(... [11]

Comment [04]: Add citations to this

Deleted: ground truth

Deleted: ground truth

Formatted: Left, Don't keep with next

Deleted: ground truth

Deleted: ¶

should be noted that RPKM cutoff of 0.8 is a very stringent threshold for identifying expressed transcripts.

For assessing the importance of mapability correction, we ran ER identification without mapability correction and computed the F-measure of the ERs. Fig 3c shows the F-measure versus RPKM threshold. Using mapability map significantly increases the accuracy of peak calls and shows the importance of utilizing the mapability correction in ER identification.

[[DISCUSS RESULTS ON REPRODUCIBILITY]]

Analysis of Polymerase II Enrichments and Gene Expression Levels

Next, in order to illustrate a novel utility for the SEFs identified by MUSIC, we concentrated on the Polymerase II binding data from ENCODE project. Polymerase shows distinct patterns of binding such that the depending on the state of polymerase, i.e., elongating or stalled, the ChIP-Seq enrichment becomes more broad and more punctate for elongating and stalled polymerase, respectively. In addition, the stalled and elongating polymerase can be distinguished by comparing the detected amount of transcription at the polymerase binding.

For evaluating the relation between the expression and the enrichment broadness as measured by SEF pileup signal, we processed and computed the SEF pileup signal (100 bases to 2.5 megabases) using the ChIP-Seq dataset for RNA polymerase II (Pol2b) from ENCODE project. For each protein coding gene, we computed the maximum value of the SEF pileup signal within the promoters. This gives us, at each gene, an estimate of the broadness of polymerase binding at the promoter. Next, we also quantified the gene expression levels in RPKMs. Finally, we plotted the joint distribution of SEF signal and gene expression level for each gene which is plotted in Fig. 5. Visual inspection of this plot reveals two components: The maximum of one component can be located at SEF pileup signal at 9 and log expression (log expression level at 2. This component can be associated with actively transcribed genes. The maximum of other component is located at SEF pileup signal at 9 and log expression level at 0. Although the maximum does not have a distinguishable local maximum, it can be spotted by looking at the distribution from two different orientations, as in Fig. 5a and 5b.

DISCUSSION:

We present a novel method, MUSIC, for the identification of enriched regions in ChIP-Seq experiments. Although MUSIC can be used to identify enrichments in any ChIP-Seq experiment, we concentrated on identifying histone ChIP-Seq experiments in this paper. MUSIC utilizes a multiscale decomposition of the ChIP-seq signal profile in conjunction with a novel mapability correction for mediating the effects of the data. Mapability is an important aspect of peak calling from next generation sequencing data especially for identifying the broad domains of enrichment since the read depth profiles are highly correlated with the mapability map. We showed that MUSIC outperforms other methods in terms of accuracy of H3k36me3 peaks in comparison with the expressed transcripts identified from the expression data from ENCODE project.

Deleted: [[TIME AND MEMORY USAGE COMPARISONS]]¶

Deleted: using SEFs

Deleted: SEFs

Deleted: [[FOLLOWING IS THE AGGREGATION PLOT OF Pol2s2 data: 4 quadrants in the expression/broadness plane]]¶

Deleted: filtering

Deleted: correction

Deleted: MUSIC, to our knowledge, is the first peak caller that takes mapability into account for identifying broad domains of enrichment at nucleotide level.

An important advantage of MUSIC is that the users can specify the scales that they would like to concentrate on, which is done using the begin and end scale parameters for the multiscale filtering. We believe this customizability will prove very useful for processing the datasets generated using ChIP-Seq experiments for which broad binding profiles are observed.

Deleted: peak calling scale

Deleted: process

Deleted: can be easily

Deleted: decomposition, which sets the scale at which the algorithm identifies the enrichment features. To our knowledge, other peak calling methods do not present an intuitive way to set the scale at which the peaks are called

Deleted: As with all algorithms, MUSIC has

[[There is no mode for one sample analysis, which is reasonable bc ...]]

There are several limitations of MUSIC. Currently, MUSIC cannot be directly used on genomes with high chromosomal aberrations, i.e., copy number variations. Although the Poisson background model partly compensates for this by modeling the read distribution over a large window, the current significance estimation by binomial p-value computation does not correct for these effects and can therefore generate spuriously high number of peak calls on regions with high copy numbers. This is a limitation that is vital to perform epigenomics analyses on the genomes with extreme copy number variations like cancer samples.

Deleted: is not addressed by many peak callers and

We believe that MUSIC is an important tool for identification and analysis of enrichments in ChIP-Seq datasets. **METHODS:**

Deleted: analysis for

Deleted: ¶
[[MUSIC allows changing the smoothing scale levels used for analyses in an intuitive manner, which is a novel feature]]¶

We describe signal processing pipeline underlying MUSIC in more detail.

Control Scaling Value Computation:

Mapability Correction and Enrichment Feature Enhancement: Given the read depth signal at each nucleotide position, MUSIC generates the per nucleotide multi-mapability signal and corrects for the mapability based loss of signal using following filtering:

$$\tilde{x}_i = \max[x_i, \text{median}(\{x_a\}_{a \in [i-l_c/2, i+l_c/2]} \mid m_a < \overline{m}_{\text{exonic}})]$$

Deleted: $\max[x_i, \text{median}(\{x_a\}_{a \in [i-l/2, i+l/2]} \mid m_a < 1.2)]$

Where x_i is the signal value at nucleotide position i , $\text{median}(\{x_i\})$ is the median of the set $\{x_i\}$, m_a is the value of the multi-mapability profile at the position a , and l_c is the window length used in mapability aware filtering. Using this filtering, MUSIC infers the signal values for positions with low mapability using the median of the values at nearby positions with multi-mapability signal lower than 1.2. We selected this value since it is the smallest multi-mapability signal profile value, i.e. most mapable, over exons and promoters as shown in Fig S1. We set the window length l_c to 2000 bps from observations. This window length depends on the distribution of length of the non-mapable region lengths. Different window lengths did not seem to affect the results too much for our tests on human genome.

Deleted: l

Deleted: xxx

Deleted: l

This filtering is inspired from the dilation operation in image processing, which is a morphological filter and has been used, in combination with other filters, for image enhancement. In our experiments, we observed that the operation defined above tends to enhance the significant features and does not change the significance of the background regions.

Multiscale Enrichment Feature Identification: Multiscale signal processing has been used in the context of wavelet transform [1] to process ChIP-Seq data and for peak calling. In this paper, we are using a more general form of multiscale filtering, namely the multiscale decomposition [2]. MUSIC utilizes a median filtering based smoothing for generating a multiscale decomposition. We selected to use median filtering since it has many applications in signal processing for performing signal smoothing with edge preserving. Given a window length, i.e. the scale, median filtering can be formulated as:

$$x_i^s = \text{median} \left(\{ \tilde{x}_a \}_{a \in [i - \frac{l_s}{2}, i + \frac{l_s}{2}]} \right), l_s \in (l_{begin}, \dots, l_{end})$$

Where x_i^s is the i^{th} value of the decomposition at scale level s for which the smoothing window length is l_s , and \tilde{x} is the mapability corrected signal profile. The window length l_s is chosen from a geometric series with the factor σ to make sure that the larger scales do not dominate the generated features.

The multiscale decomposition enables automatic identification of blobs in the signal profiles at different scales with very small computational requirement. MUSIC uses a fast and efficient method to implement the median filtering by storing the histogram of the signal values in the window and processes only the new and obsolete signal values that enter and leave, respectively, the current window to update the histogram when moved to the next window.

It should be noted that any type of smoothing filter can be used, e.g., Gaussian, Triangular, Rectangular, etc., to generate the multiscale decomposition for peak calling.

Identification of Scale Specific Enriched Regions: After the multiscale decomposition, MUSIC identifies all the local minima points in the decomposition. MUSIC utilizes regions between minima points as the regions of enrichment. For this, MUSIC computes the derivative of the signal at each point as the difference between consecutive values:

$$x'_i{}^s = (x_i^s - x_{i-1}^s)$$

where $x'_i{}^s$ is the derivative of the smoothed signal x_i^s . MUSIC assigns the local extrema points at the points where the derivative changes sign:

$$I_{min} = \{i \mid x'_i{}^s < 0 \text{ and } x'_{i-1}{}^s > 0\}$$

$$I_{max} = \{i \mid x'_i{}^s > 0 \text{ and } x'_{i-1}{}^s < 0\}$$

Where I_{min} and I_{max} are the sets of positions of minima and maxima of x_i^s , respectively. The scale specific candidate enriched regions of x_i^s are identified as the regions between the consecutive minima.

Comparison of Smoothed Signal in Candidate Enriched Regions: For the candidate enriched regions in each smoothing scale, MUSIC uses the value of smoothed signal levels and unsmoothed signal levels for assessing the quality of enriched region. A scale specific candidate enriched region is filtered if the ratio

Deleted: , l_{end}

Deleted: , that MUSIC uses has an exponential increase

Formatted: Font color: Background 1

Deleted: MUSIC currently supports three filters: Median, Gaussian, and Mean. To our knowledge MUSIC is the first algorithm to utilize a non-linear median filtering for multiscale feature identification for processing genomic signal profiles.

Formatted: Font color: Background 1

Deleted: **[[ADD PER SCALE SIGNAL LOSS FILTER: Maximum signal smoothing threshold (denoted by γ)]]**

Deleted: **Enrichment Features**

Deleted: We refer to these regions as *enrichment features*.

Moved (insertion) [2]

Deleted: **Feature**

of the maximum of smoothed signal to the maximum of the unsmoothed signal within the candidate region is higher than the smoothing ratio threshold, γ . In other words, MUSIC removes the candidate enriched region $[i, j]$ at scale s , if

$$\frac{\max(\{x_a^s\}_{a \in [i, j]})}{\max(\{x_a\}_{a \in [i, j]})} < \gamma.$$

This test is designed as a simple and efficient check to evaluate whether the signal within the candidate region identified at the scale level s is severely smoothed. This way MUSIC efficiently detects and avoids overmerging of consecutive regions that have high signal enrichment and are close to each other. In addition, MUSIC removes the enriched regions whose signal levels are severely smoothed. By default γ is set to 0.25.

Candidate Enriched Region End Trimming using Poisson Distribution Model: MUSIC trims the ends of the candidate enriched regions using a Poisson null model for the signal distribution. For this, MUSIC divides genome into 1 megabase windows and for each 1 megabase window estimates the mean of all the values. Using this as the mean parameter μ of the Poisson distribution, MUSIC selects a threshold that satisfies 5% false positive rate:

$$\tau = \operatorname{argmin}_t \{F_{X_\mu}(t) > 0.95\}, X_\mu \sim \text{Poisson}(\mu)$$

Where F_{X_μ} represents the cumulative distribution function of X_μ , which is distributed as Poisson with mean μ . For a feature with start and end at positions i and j , respectively, the trimmed end coordinates are given as:

$$i' = \operatorname{argmin}_{a \in [i, j]} (x_a > \tau)$$

$$j' = \operatorname{argmax}_{a \in [i, j]} (x_a > \tau)$$

Where i' and j' are the trimmed start and end coordinates, respectively. The features that do not pass the threshold are removed from the candidate peak list.

Enriched Region End Trimming via p-value Minimization: Then MUSIC further fine-tunes the ends of the merged features using a novel p-value minimization using the chip and control profiles. For a given region, MUSIC starts thresholding the signal at the ends of the region and identifies the signal height at which the p-value of the region is minimized. This maximizes the compactness of the merged feature regions. The end-refined merged feature regions are the candidate regions of enrichment before p-value computation.

$$i' = \operatorname{argmin}_{a \in [i, j]} (p(a, j \mid l_{pval} = (j - a + 1)))$$

Deleted: features first

Deleted: first

Deleted: $\operatorname{argmin}_a (x_a > \tau), a \in (i, j)$

Deleted: $\operatorname{argmax}_a (x_a > \tau), a \in (i, j)$

Deleted: Feature

Deleted: trimming

Deleted: minimization

Deleted: $\operatorname{argmin}_a (p(a, j \mid l_{pval} = (j - a + 1))), a \in (i, j)$

$$j' = \underset{a}{\operatorname{argmin}}(p(i', a \mid l_{pval} = (a - i' + 1))), a \in [i', j]$$

where $p(a, b \mid l_{pval})$ represents the p-value for the peak starting at a and ending at b with the length of p-value window given by l_{pval} (Refer to p-value computation.)

Enriched Region Merging: After the features are identified, MUSIC ~~takes the union~~ all the features and identifies where the clumps of features are. This is done basically by identifying the positions that are covered by at least 1 feature.

Per Strand Concordance Filter: For each ER, MUSIC computes the total signal on positive and negative strands and filters out the enriched regions for which there is high discordance between the signals:

$$\min\left(\frac{\sum_i x_i^+}{\sum_i x_i^-}, \frac{\sum_i x_i^-}{\sum_i x_i^+}\right) < 0.5$$

where $\sum_i x_i^+$ and $\sum_i x_i^-$ is the total signal on the positive and negative strand within the start and end coordinates of the ER, respectively.

P-value computation and FDR Estimation: We use one-tailed binomial test to compute the p-values for each end-refined merged feature region. We first count the number of reads in the chip sample (n_{chip}) and control sample ($n_{control}$) that overlap with the region, then compute one tailed p-value as:

$$p = \sum_{r=n'_{chip}+1}^{n'_{chip}+n'_{control}} \binom{n'_{chip} + n'_{control}}{r} 0.5^{(n'_{chip}+n'_{control})}$$

Where n'_{chip} and $n'_{control}$ are the normalized read counts for the region:

$$n'_{chip} = \frac{n_{chip}}{l_{chip}} \times l_{pval}$$

$$n'_{control} = \frac{n_{control}}{l_{control}} \times l_{pval}$$

Where l_{pval} is the length of the p-value computation window and p refers to the p-value value for the peak. Larger values of l_{pval} increase the significance of regions (See parameter selection). We correct for the p-values using the Benjamini-Hochberg procedure to generate the corrected p-values, i.e., q-values:

$$q_i = p_i \times \frac{N_{peaks}}{i}$$

where N_{peaks} is the total number of peak regions and i is the rank of the peak in the peak list sorted with respect to increasing p-value. By default, MUSIC uses q-value cutoff of 0.05. The filtered peaks are reported in BED format with their q-values in the score field.

Deleted: Feature

Formatted: Font color: Background 1

Deleted: merges

Formatted: Font color: Background 1

Deleted: peaks

Deleted: Correction

Multi-Mapability Signal Generation: MUSIC can generate per nucleotide multi-mapability signal profiles. For this it is required to have a read mapping program installed on the system. Currently MUSIC uses bowtie2^[XXX], a very popular fast read mapping algorithm, by default. MUSIC first fragments all the chromosomes to the read length of interest, maps all the fragments to the genome using bowtie2 with 2 mismatches and reporting of maximum of 5 multimapping positions per fragment. Then MUSIC uses the mapped reads to build the mapability signal profile. The regions with high signal corresponds to regions with low mapability. Then MUSIC processes the mapability profile to store space since it does not require the whole mapability signal profiles. We generated mapability maps for hg19 genome for read lengths of 36, 50, 76, 100, and 200 bps that are available for download with MUSIC.

Parameter Selection for Benchmarking: There are 3 parameters associated with MUSIC, starting scale window length, ending scale window length, and the p-value computation window length. [In order to select the parameters while comparing with other algorithms,](#)

Deleted: *[[HOW ARE THE THRESHOLDS SELECTED: Add the training F-measure plots.]]*

Deleted:

[In order to select the window lengths for broad scale peak calls \(H3k36me3, H3k27me3\), we used H1HESC human datasets as training dataset \(since we used K562 and GM12878 for benchmarking\) and ran MUSIC with a large range of parameter sets for the three of F-measure versus percentage overlap between H3k36me3 and H3k27me3 peak calls. This is necessary because we observed that the F-measure increases as the window scales are increased for H3k36me3 dataset. We chose the parameter set that has yields highest F-measure while the overlap percentage is below 1 percent. This parameter set turned out to be \$l_{base}=1100\$ bps, \$l_{end}=14000\$ bps, \$l_{pval}=1750\$ bps.](#)

Formatted: Font color: Background 1

Accuracy Measures: For evaluating the accuracy of H3k36me3 peak calls, we computed sensitivity, positive predictive values:

$$Sensitivity = \frac{covg(P \cap G)}{covg(G)}$$

$$PPV = \frac{covg(P \cap G)}{covg(P)}$$

Where $covg(P)$ is the coverage of peaks, $covg(G)$ is the coverage of expressed gene bodies and $covg(P \cap G)$ is the coverage of the overlap between expressed gene bodies and peaks. We combined these two accuracy measures to compute F-measure, computed as:

$$F - measure = \frac{2 \times Sensitivity \times PPV}{(Sensitivity + PPV)}$$

For H3k4me3 peaks, we used all the promoters (TSS of the transcript ± 2500 bps). For these, we use a slightly different approach to compute sensitivity and PPV:

$$Sensitivity = \frac{\#(S \cap P)}{\#(S)}$$

$$PPV = \frac{\#(P \cap S)}{\#(P)}$$

Where $\#(S)$, $\#(P)$, $\#(P \cap S)$ represent number of active promoters, number of peaks, and number of peaks that overlap with active promoters, respectively.

Datasets and Data Processing: We downloaded ENCODE ChIP-Seq from UCSC genome browser. The RNA-seq expression quantifications are downloaded from ENCODE RNA Dashboard. For the transcript quantifications, we used the average RPKM values for the transcripts from two replicates that satisfied the reproducibility criteria that iDR smaller than 0.1.

SUPPLEMENTARY MATERIAL

Mapability is an important factor for processing genome wide signals. This stems from the fact that the signal levels at region with low mapability will show a systematic decrease at the nucleotide resolution. We used the multi-mapability signal profiles generated by MUSIC (See Methods) and aggregated the

signal on different regions (Fig. S1). Promoters and the regions downstream of TSS into the first exon show significantly higher mapability compared to random regions, regions that are upstream into the intergenic side of the genes show significantly lower mapability compared to . In addition, introns show slightly higher mapability compared to random regions and exons show are much more mapable than random regions. Transcription start sites and mid points of exons show almost the same amount of average multi-mapability, 1.2 reads.

Comparison of H3k4me3 ER accuracy with Other Methods:

For H3k4me3, we used the active promoter identification accuracy per top set of peaks of each method for comparison. Although we did not have a negative set for H3k4me3 peaks, unlike H3k36me3, since H3k4me3 is predominantly associated with promoters, we assumed that the top peaks from peak calling will be enriched in active promoters. Starting from the top peaks (sorted with respect to the score reported by each method), we computed the F-measure for promoter identification for each method with changing fraction of coverage of top peaks for the top 30 megabases of the peaks. This way we can evaluate the accuracy of peak calls with changing peak rank. For each peak caller, we sorted the peaks with respect to the reported score. MUSIC tends to perform as one of the best (with MACS) for the accuracy of the top peaks.

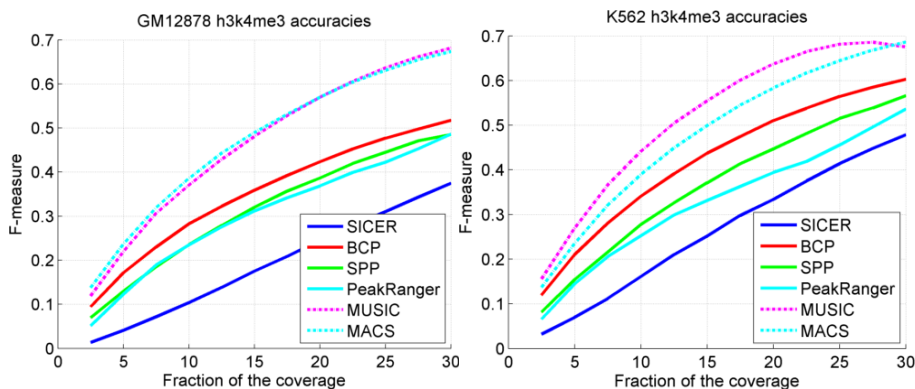
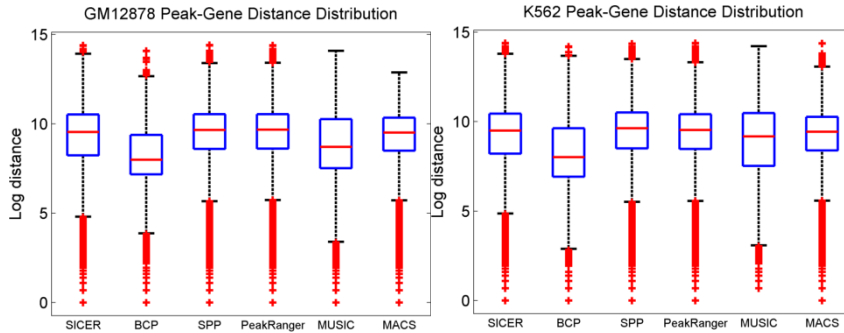


Figure 1: F-measure vs coverage of H3k4me3 peaks in GM12878 (left) and K562 (right)

Next we evaluated the distribution of the distance between the ends of gene bodies and the ends of H3k36me3 peaks to identify whether the peak ends match with the annotated ends of the genes. Figure XXX shows the distribution of smallest peak end to gene end distance for all the peaks for all the methods. The median values are highlighted in the plots to compare the methods with each other.

MUSIC has the second smallest median value following BCP.



We also evaluated the reproducibility of the peaks generated by the peak callers. We used the replicates generated by ENCODE with the same HM datasets to assess reproducibility of peak calling. Figure XXX shows the average of fraction of the overlapping regions to the total coverage of each replicate. MUSIC has higher reproducibility for H3k27me3 and H3k36me3 than all other methods except for K562 H3k36me3 dataset, where BCP has slightly higher reproducibility than MUSIC. For K562, MUSIC has highest reproducibility for H3k27me3. For H3k36me3, BCP has slightly higher reproducibility than MUSIC. Overall, MUSIC has higher or comparable reproducibility with respect to other peak callers.

Formatted: Font: 14 pt

Formatted: Font: 11 pt

Page 1: [1] Deleted Ozgun 2/22/2014 3:08:00 PM

uses a mapability correction procedure coupled with multiscale signal processing based approach to identify the significant enrichment features (SEFs) that represent the significant enrichments at different length scales in the signal. By piling up the SEFs, MUSIC generates a genome-wide signal that can be utilized for quantifying the broadness of enrichment at each location. We show the utility of SEFs and SEF pileup signal within two applications

Page 1: [2] Deleted Ozgun 2/22/2014 3:08:00 PM

we identify the enriched regions (ERs) using the SEFs and

Page 1: [3] Deleted Ozgun 2/22/2014 3:08:00 PM

Second, in order to showcase a novel application of the SEF pileup signal, we concentrate on processing the

Page 1: [4] Deleted Ozgun 2/22/2014 3:08:00 PM

for Polymerase II binding. The broadness of enrichments in signal profiles for polymerase binding can be used as an indicator of the polymerase

Page 1: [5] Moved to page 3 (Move #1) Ozgun 2/22/2014 3:08:00 PM

state, i.e., stalled or elongating.

Page 1: [6] Deleted Ozgun 2/22/2014 3:08:00 PM

We process the SEF pileup signal at protein coding gene promoters and demonstrate the bimodality of the joint distribution of signal broadness at the promoter versus the gene

Page 1: [7] Deleted Ozgun 2/22/2014 3:08:00 PM

sites, or peaks, from ChIP-Seq experiments

Page 1: [8] Deleted Ozgun 2/22/2014 3:08:00 PM

Identification of the broad enrichments in the read depth signal profiles, however, did not receive the same amount of attention. The

Page 3: [9] Deleted Ozgun 2/22/2014 3:08:00 PM

Remaining of the paper is as follows. We first describe the MUSIC algorithm and identification of SEFs. Next we present identification of ERs and compare the accuracy and reproducibility of MUSIC other ER identification methods. Then we focus on joint processing the signal profiles for Polymerase II and gene expression levels. Finally, we present the algorithmic details of MUSIC in Methods.

Page 3: [10] Deleted Ozgun 2/22/2014 3:08:00 PM

a novel application of the enrichment features identified by MUSIC to

The p-values are corrected by Benjamini-Hochberg procedure\cite{XXX}. The final corrected p-values are thresholded with respect to 0.05 for identification of significant ERs. MUSIC can be utilized to determine ERs from precomputed SEFs (“-get_ERs_per_SEFs”), or identify ERs from reads (“-get_peaks_per_reads”).