

## RESPONSE LETTER

### -- Ref1.1.1 – Presentation of in vivo validations --

Reviewer Comment	I understand that the authors tested 102 predicted mouse enhancers (plus 31 human orthologs) in transgenic mice, and had another 151 regions from an independent unpublished effort (Moore, in review) available for comparison. This is an unprecedented effort to assess enhancer predictions in vivo, making a systematic and rigorous comparison between the predictions and the experimental outcomes of the in vivo assays highly interesting. However, I find the presentation in the main text and figures not satisfying and partly confusing. For example, what does "61% predicted active rate versus 70% observed active rate" (page 10) mean? I interpret this statement as 61% of the tested regions were predicted to be positive and 70% of the tested regions were found to be positive - there is no indication if the predicted and observed positives actually agree.
Author Response	We thank the referee for pointing this out. We agree that this sentence is a bit confusing and not quite accurate, and we'll rewrite it. Here we are describing the experimental test result of 62 elements chosen from top, middle and bottom rank of forebrain H3K27ac signal (e.g. how many of them are active in each tier). We made a rough estimation of whether these elements would be active by their overlap with the DHS peaks, but since this estimation is not very relevant, we can remove them to avoid confusion. A rigorous assessment of our model prediction using these experimental data is presented later in the table and ROC/PR curve of Figure 4. Here we are showing that indeed the highest ranking tier has the highest validation rate, and we provide the detail validation result of each element in the supplementary table.
Excerpt From Revised Manuscript	

### -- Ref1.1.2a – Presentation of in vivo validations –

**### Not ready**

Reviewer Comment	It is my understanding that the authors have predictions for different mouse tissues and - for each tested candidate - have a readout of activity across the entire embryo, i.e. all tissues. This should allow the rigorous assessment of the prediction accuracy per tissue in comparison to an appropriate random model that accounts for the overall number of active regions per tissue (I
------------------	---

	assume Fig. 4B and C come close to this, but the corresponding text is confusing - I don't understand what Fig. 4A corresponds to).
Author Response	<p>Indeed each candidate has a readout for all tissues in the embryos. Figure 4B and 4C use this experiment read out to evaluate the prediction. The ROC curve evaluates the false positive rate and true positive rate of our prediction in each tissue at different threshold, and the PR curve evaluates the precision and recall of the prediction. At random, the ROC curve is a diagonal line with AUROC of 0.5, and PR curve is a horizontal line with AUPR equal to the fraction of positives. In 4A the table contains the average ROC/PR of the evaluation results from six different tissues. We are reporting averaged numbers. We have modified the text to make this clearer.</p> <p>We've rewritten the text to make clearer *** Rewrite the text to describe the experiments</p>
Excerpt From Revised Manuscript	<p>The expt have this readout The figure show YYY</p> <p>We evaluated the predictability of our matched filter model for each individual histone marks and DHS, as well as the integrated SVM model (Figure 4). Consistent with previous findings from STARR-seq data, H3K27ac signal is the single best performed histone marks for predicting enhancers, while DHS signal performs well as an independent source. The integrated model, as expected, achieves higher predictability than individual histone marks. We then did similar evaluation using the regulatory elements identified by the transduction-based FIREWACH assay in mouse embryonic stem cells (mESC) [36]. With the same metaprofiles, the predictions are based on epigenetic signals of mESC available from ENCODE website. Again, we observe similar results for individual histone marks and combined SVM model (Figure S16). As the <i>in vivo</i> and FIREWACH assays utilized a single core promoter to validate regulatory regions, the performance of the different models in Figures 4 and S16 are probably underestimated.</p>

	<p><b>Figure 4: Conservation of epigenetic features. The performance of the <i>Drosophila</i> STARR-seq based matched filters and the integrated model for predicting active enhancers identified by transgenic mouse enhancer assays at 6 different tissues in E11.5 mice. A) Average AUROC and AUPR for predicting enhancers by different features and by the integrated model. The weights of the different features in the integrated model is the same as the weights shown in Figure 3 for enhancers. B) The individual ROC curves of each feature and the integrated model for each tissue are shown. C) The individual PR curves of each feature and the integrated model for each tissue are shown.</b></p>
--	--

**-- Ref1.1.2b – Presentation of in vivo validations --**

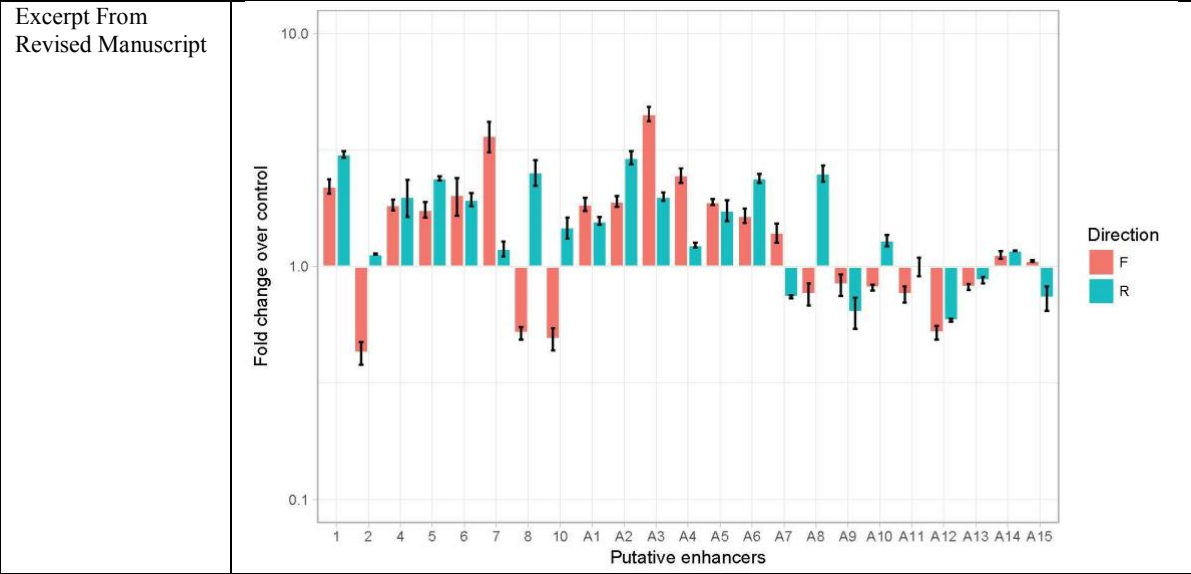
Reviewer Comment	Also, the raw images should be made available either as supplementary information or via a suitable website (e.g. the VISTA database).
Author Response	We've made the raw images of these experimental results available through the VISTA enhancer browser.
Excerpt From Revised Manuscript	

**-- Ref1.2.1 – Validation in human cell lines --**

Reviewer Comment	<p>I find the presentation of the validation in human cell lines confusing and not sufficiently well controlled. Most importantly, the tests for the individual enhancers don't seem to be replicated, such that one cannot draw any statistically sound conclusion about the activity of each putative enhancer. Reported are only two numbers (corresponding to the fold change of gene expression of each enhancer in the forward and reverse orientation) in 4 different cell lines (table S7). These numbers often don't agree well and in some cases, the nature of these numbers is unclear. For example, what does "0. 1.06" or "0, 1.73" (note the "." vs. ",") mean - did the forward experiment fail or was the outcome exactly 0? These validations need</p>
------------------	--

to be performed in triplicates per cell line and construct such that each region's activity can be rigorously assessed, allowing the subsequent assessment of the predictions for each cell line. Alternatively, the cell lines for which replicate experiments cannot be performed should be removed to maintain a minimal quality standard for such validation experiments.

**Author Response**  
 We acknowledge the referee's comment. In the revised manuscript, we describe the details of the human cell line validation experiments to make it more clear. The original experiment tested each enhancer in all four cell lines in replicates for both forward and reverse orientation. The read out of each experiment was normalized to the control. The numbers in the table represent fold change over controls, where 0 occurs when the number of positive cells is less than that of control according to FlowJo gating. Based on the referee's suggestion, we performed another set of triplicate experiments on these randomly selected putative enhancers in H1-hESC (23 out of 25 elements successfully went through PCR and transduction. We remove the two elements for which the experiments cannot be performed based on the comment). The triplicate experiment read out is consistent with our previous report. We show the result of each replicate in the supplementary table XX and a supplementary figure is provided to visualize the data. As the figure shows, the validation experiments are highly reproducible, with the correlation between each pair of replicate being 0.9 and above.



**-- Ref1.2.2 – Validation in human cell lines --**

Reviewer	The same applies to the two statements in the main text
----------	---

Comment	(page 11): "a few elements showed significantly higher levels of gene expression in one of the orientations" and "even though some of the elements were preferentially active in one of the cell lines". Both statements are not sufficiently supported by data: neither has a systematic comparison been done, nor are the data on which these statements are based replicated. These experiments need to be performed according to minimal quality standards or the statements need to be removed.
Author Response	<p>Here we are describing part of the experiment result rather than making strong statement about the directionality of general enhancer activity. As shown in the figure above, we find that some elements (eg, 7, 8 and A8) have significant different fold change (compared to control) for different directions, and the results are based on three replicates. However, as we are not trying to make strong statement about the directionality of enhancers, we agree to remove this description and present the raw data to the readers.</p> <p>As we clarified under section 1.2.1, the experiments are done in replicates and are normalized under the control.</p>
Excerpt From Revised Manuscript	We will remove these two sentences in case of any confusion.

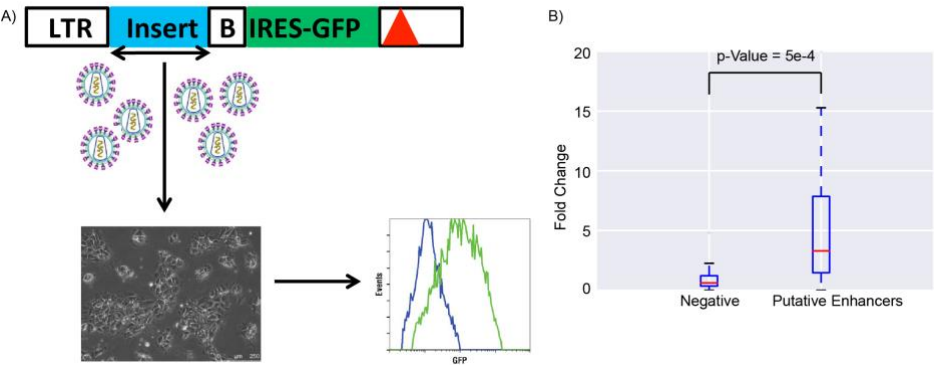
### -- Ref1.2.3 – Validation in human cell lines --

Reviewer Comment	The presentation is also confusing: for example, figure 5 and the main text state that the Oct4 promoter is used, but also that a "housekeeping promoter is used" (page 11). Figure 5 shows an IRES-GFP construct, which is typically used in combination with a selection marker, yet no such marker is shown and the methods don't indicate selection (which would distort enhancer activity measurements). The authors should also comment on the LTRs' promoter function and if this could influence their results.
Author Response	We have made changes to the description in the manuscript so it is clearer. A minimal basal Oct4 promoter was used in the SIN HIV vector since a primary focus of the work was DNA elements active in hESC. IRES-eGFP was used downstream of the DNA elements to allow flow cytometric analysis of positive cells after cell transduction. The presence of a selectable marker gene would have needlessly increased the size of the vector, which would be problematic for some of the longer elements. IRES was used so that there would be eGFP translation/readout even if transcription began within the element itself, several kb upstream of eGFP start codon. To address concerns regarding the HIV LTR, figure 5 now shows SIN HIV vector structure after genomic

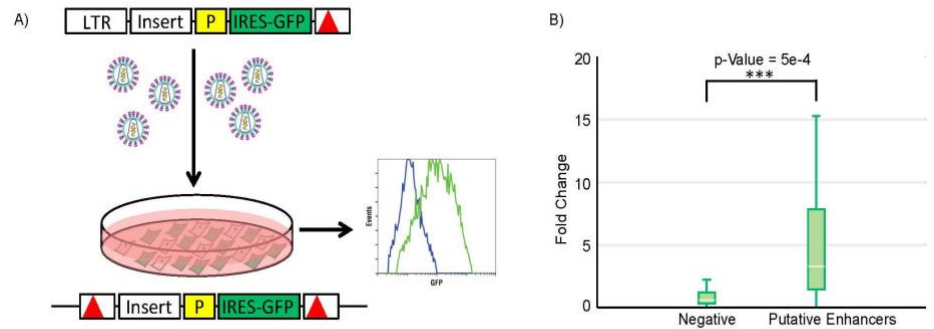
integration, with the duplication of ~400 bp deletion of the U3 portion of the LTR. This essentially renders the LTR inactive. However, to take into account possible residual activity (and any activity of the basal Oct4 promoter), all of the transduction data is normalized to that of EV, tested on the same cells.

Excerpt From Revised Manuscript

Original figure  
Figure 5



Revised  
Figure 5



Will add revised figure caption

**-- Ref1.3.1 – Prediction algorithm --**

Reviewer Comment

The brief description of the metaprofile-based predictions on page 6 suggests optimization steps that are not well explained and could break cross-validation if performed incorrectly. Specifically, the authors state that they “scanned with multiple matched filters with templates that vary in width between the two maxima in the double peak” (page 6). How many such templates are used and how many parameters does this add to the model? In the 10-fold cross-validation, are these templates exclusively derived from the training set or are they created prior to cross-validation (which would break it!)?

Author Response	<p>We see this wasn't clear</p> <p>we've modified ...</p> <p>Thanks for the referee's comment. In the original text:</p> <p>“Due to the aforementioned variability in the double peak pattern, the H3K27ac signal track is scanned with multiple matched filters with templates that vary in width between the two maxima in the double peak and the highest matched filter score with these matched filters is used to rate the regulatory potential of this region(see Methods). The dependent profiles are then used on the same region with the matched filter to score the corresponding genomic tracks.”</p> <p>During the ten fold cross validation with a single histone mark, the profiles are created with 90% of the STARR-seq positives and 10% of the positives are used for testing the accuracy of the model. With the main SVM model within the manuscript, 6 different matched filter profiles are created with 90% of the STARR-seq positives and 10% of the positives are used for testing the accuracy of the SVM model.. We have modified the manuscript [[supplement]] to make clearer</p> <p>We did the corss validation properly</p>
Excerpt From Revised Manuscript	<p>To discuss with ANS</p> <p>During the ten fold cross validation with a single histone mark, the profiles are created with 90% of the STARR-seq positives and 10% of the positives are used for testing the accuracy of the model. With the main SVM model within the manuscript, 6 different matched filter profiles are created with 90% of the STARR-seq positives and 10% of the positives are used for testing the accuracy of the SVM model.</p>

**-- Ref1.3.2 – Prediction algorithm --**

Reviewer Comment	<p>I also note that the result that H3K27ac has the highest predictive value and that DHS is partly redundant to H3K27ac is highly confounded by 1. the choosing of templates based on H3K27ac and subsequent application to the other histone modifications (page 12, top paragrah) and 2. the fact that the metaprofile with the two maxima and the dip in-between (plus its width) already captures the DHS signal, which is complementary.</p>
Author Response	<p>Thanks to the referee for the comment. Indeed we show that H3K27ac has the highest predictive value and that DHS is partly</p>

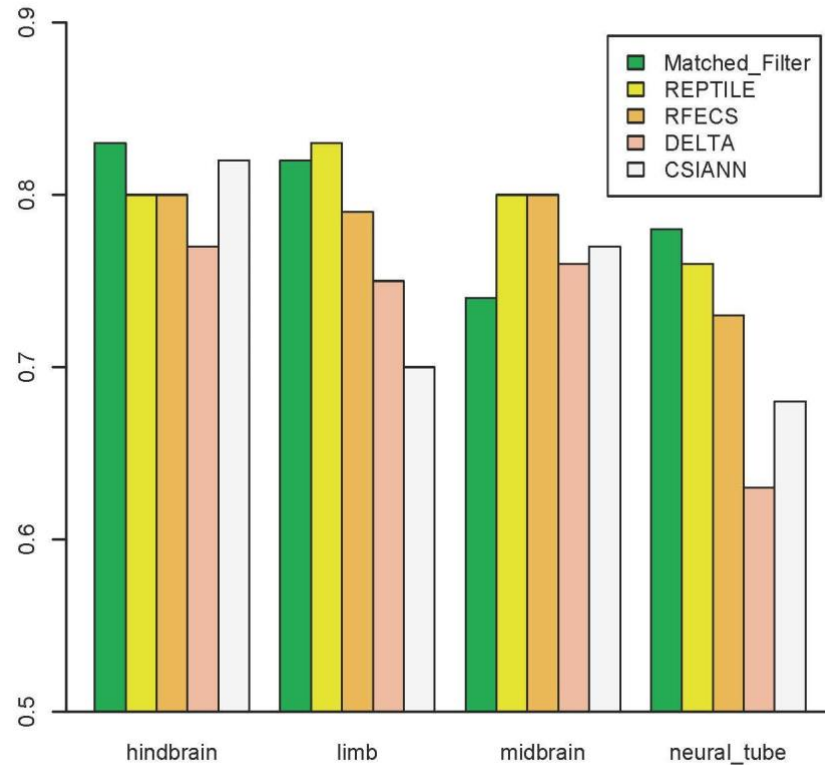
	<p>redundant, as indicated by the model.</p> <p>H3K27ac has the highest performance even when we compare all histone marks independently. So it's not surprising that the model selects H3K27ac as the highest predictive value.</p> <p>As for the redundancy between DHS and H3K27ac, we agree that the dip in between the two maxima is usually where the DHS peak would occur, which provides good explanation for the redundancy. We have added this discussion in the manuscript as shown below.</p>
Excerpt From Revised Manuscript	<p>According to the model, the acetylations (H3K27ac and H3K9ac) are the most important feature for predicting active regulatory regions. The DHS matched filter performed well as an individual feature (AUPR in Figure 2) to predict enhancers, but had a lower weight among the six features likely due to the fact that the information in DHS is redundant with the information contained within the histone mark, eg. the DHS peaks usually occur at the trough region between two maxima in the histone signal. Despite the redundancy, combination of the DHS and histone signals is more predictive of regulatory activity as the complementary signals are strengthened compared to the uncorrelated noise in each signal.</p>

**-- Ref1.4 – Comparison with previous methods --**

Reviewer Comment	<p>The authors compare their approach to chromHMM and SegWay, which are both not built for enhancer prediction but rather to segment the genome into different types of regions. A more relevant comparison to a supervised machine learning approach (Capra, ref 64) is presented only superficially in the methods section and without any (supplementary) figure.</p>
Author Response	<p>ChromHMM and SegWay are initially built for segmentation of the genome and provide annotations for different genome regions. The ChromHMM and SegWay enhancer annotations of the Roadmap Epigenetics samples has been used in many publications as a way to define enhancer regions. We want to compare with them to show that our framework provides a better set of enhancers readily available for related studies. Based on the referee's suggestion, we also did more comparison with other published methods, and we have included the results in our manuscript as shown below.</p>
Excerpt From Revised Manuscript	<p>In addition to the comparison with unsupervised segmentation based methods, we also compared with other published enhancer prediction tools, including CSIANN, a neural network</p>



based approach; DELTA, an ensemble model integrating different histone modifications; RF ECS, a random forest model based on histone modifications, and REPTILE, a more recent published method that integrates histone modifications and whole genome bisulfite sequencing data. We show that our method also outperforms these previous approaches assessed by in vivo transgenic experiments.



**-- Ref1.5 – Critique to main text and referencing --**

Reviewer Comment	The main text needs to be substantially revised to improve clarity and avoid repetitiveness. While some parts explain fundamental basics in great detail, such as the difference between ROC and PR statistics (pages 5-6), other more important details are missing. For example, it only becomes obvious in the methods but not in the main text (page 5) that only STARR-seq enhancers with a H3K27ac and DHS peaks are considered (page 3 in the supplement).
Author	ANS

Response	As STARR-seq quantifies enhancer activity in an episomal fashion, they mentioned in their paper that “the complementary DHS-seq and ChIP-seq determine enhancer-associated characteristics in the endogenous genomic context”. We took the overlap of the STARR-seq enhancers with H3K27ac/ DHS peaks to get a high confident set of enhancers that are active in vivo.
Excerpt From Revised Manuscript	Rewrite main text

### -- Ref1.6 – Negative control regions --

Reviewer Comment	The restriction of the STARR-seq enhancers to those that intersect with H3K27ac and DHS peaks (supplement page 3, see also my last point) and the selection of negatives as “randomly chosen regions in the genome with H3K27ac signal that had the same width distribution of the distance between double peaks near STARR-seq peaks (supplement pages 3-4) makes me wonder how H3K27ac can be the most predictive feature: if the negatives controls are chosen to match the positives in H3K27ac signals (which is a very powerful control), the predictive value of H3K27ac should be minimal or even zero. In this respect, the results are strange and the authors need to investigate the reasons for this outcome.
Author Response	Thanks the referee for the comment. For negative regions we match the width distribution which is essentially selecting regions that has similar lengths to the enhancers. These regions does not have the same H3K27ac signals in terms of the signal strength and pattern, but mostly have some background H3K27ac signals that the model would learn to distinguish from. We didn't choose non-STARR-seq peaks with no H3K27ac signal as they wouldn't provide enough information for training. Based on the comment, we have made it more clear how we select the negatives in this section of supplement as reproduced below.
Excerpt From Revised Manuscript	The negatives are randomly chosen non-STARR-seq-peak regions in the genome that had the same lengths distribution as the enhancers from the STARR-seq. We require most of the regions contain some H3K27ac signals, since negatives with no H3K27ac signal at all wouldn't provide enough information for training.

### -- Ref1.7.1 – Minor comments: Title and Abstract --

Reviewer Comment	The message that the authors' approach is trained on Drosophila enhancers und functions successfully across
------------------	---

	different species does not come across very clearly in the title and abstract, which could be improved.
Author Response	<b>To discuss</b> Current: A framework for supervised enhancer prediction with epigenetic pattern recognition and targeted validation across organisms
Excerpt From Revised Manuscript	

**-- Ref1.7.2 – Minor comments: Reference --**

Reviewer Comment	The referencing of manuscripts is broken and needs to be fixed: several references seem to not be correctly formatted (e.g. "cite 31, 50" on page 5, "linear SVM [54]" on page 7 points to the wrong paper, "(see Supplement)" on page 12 is an unclear reference).
Author Response	We thank the referee for pointing out the formatting issue and we've fixed the citations accordingly.
Excerpt From Revised Manuscript	The STARR-seq studies on <i>Drosophila</i> cell-lines provide the most comprehensive MPRA datasets as the whole genome was tested for regulatory activity within these assays and these assays were performed with multiple core promoters [31, 49].  We built an integrated model with combined matched filter scores of the most informative epigenetics marks (H3K27ac, H3K4me1, H3K4me2, H3K4me3, H3K9ac, and DHS) associated with active regulatory regions using a linear SVM [59].

**-- Ref1.7.3 – Minor comments: BG3 cells --**

Reviewer Comment	On page 7, it seems that the authors conclude from a good performance in BG3 cells that the SVM model 'is applicable across species'. Please note that BG3 cells are also <i>Drosophila</i> cells.
Author Response	Thanks for pointing this out. Indeed, the validation experiments described later in the paper shows that the model is applicable across species, but the BG3 cell line validation here is to show that our model is applicable across different cell lines.
Excerpt From Revised Manuscript	The model is highly accurate at predicting active enhancers and promoters in the S2-cell line (Figure S6), indicating our framework of combining epigenetic features with a linear SVM model to predict enhancers is applicable <b>across different cell lines</b> .

**-- Ref1.7.4 – Minor comments: Term correction --**

Reviewer Comment	"impute chromatin status" (page 12) should be "segment the genome based on chromatin features" or similar.
Author Response	We have rephrased the sentence as shown in the excerpt below.
Excerpt From Revised Manuscript	We first did the comparison with ChromHMM[63], a well known method to segment the genome based on chromatin features

**-- Ref1.7.5 – Enhancer-specific factors --**

Reviewer Comment	The differential distribution of factor binding between enhancers and promoters (page 12 and figure 6) shows many signals for promoters but only very few (and relatively weak ones) for enhancers. Are there no enhancer-specific factors?
Author Response	There are some TFs that preferentially bind to enhancers as compared to promoters. However, a few of the TATA-binding proteins bind to every promoter predicted to be active according to our model. In comparison, among the TFs with experimentally measured ChIP-seq experiment, there is no single TF that binds to a majority of predicted enhancers. The TFs that bind to active enhancers tend to bind to smaller subsets of enhancers. This could explain why, unlike promoters, it has been hard to find a strong sequence signature associated with enhancers as a diverse set of TFs tend to bind to these regions.
Excerpt From Revised Manuscript	

**-- Ref2.1a – Comparison with FANTOM5 and ENCODE --**

Reviewer Comment	<p>Page 3: "In addition to the small numbers, the validated enhancers were typically selected based on conserved noncoding regions [17] with particular patterns of chromatin [18], transcription-factor binding, [19] or noncoding transcription [20]."</p> <p>Since the FANTOM5 Atlas is the most comprehensive collection of transcribed enhancers across different primary cells and tissues, I would like to see a comparison of the model predictions in human to the enhancer dataset of the FANTOM5 Atlas dataset taking into account cell-type/tissue specificity. In a similar fashion, what is the overlap with the integrative ENCODE annotation proposed by Hoffman et al. NAR 2013. Assuming that the size of training datasets is the only limiting factor for achieving high discrimination performance, what is the minimum number of samples that guarantees good performance in the deployed method?</p>
------------------	---

Author Response	<p>Thanks to the referee for this point. The FANTOM5 Atlas contains a good set of transcribed enhancers, although there is only a relatively small number of transcribed enhancers detected in each cell. Based on the referee's suggestion, we've checked our predictions against the FANTOM5 enhancer set and compared our overlap with the annotation provided by Hoffman et al, NAR 2013. We included the result in the supplement as reproduced below:</p> <p><b>put in the supplement figure caption</b></p> <p>We overlapped the CAGE-defined enhancers from FANTOM5 to our predicted enhancers and also to the enhancers predicted by the integrative ENCODE annotation method proposed by Hoffman for cell lines including GM12878, K562 and HepG2. We found that around 40% of the CAGE-defined enhancers overlap with our predicted enhancers, and only 23% to 34% of the CAGE-defined enhancers overlap with the enhancers predicted by integrative ENCODE annotation method, although the latter provides much larger numbers of predictions (about four times for GM12878 and K562, and three times for HepG2).</p>
Excerpt From Revised Manuscript	

**-- Ref2.1b – Comparison with FANTOM5 and ENCODE --**

Reviewer Comment	<p>Assuming that the size of training datasets is the only limiting factor for achieving high discrimination performance, what is the minimum number of samples that guarantees good performance in the deployed method?</p>
Author Response	<p>We performed detailed saturation analysis and is shown under comment 2.3. Briefly, the method achieves good performance in terms of prediction accuracy with even a small number of training dataset, but the variance of the performance decreases with larger training dataset. The large training dataset also decreases the false discovery rate. We found that with half of the size of the current training dataset, we would get reasonably good prediction accuracy and low FDR, but the stable performance and the saturation of precision-recall is achieved with around 80% of the current training data size.</p>
Excerpt From	

**-- Ref2.2 – Method justification --**

<p>Reviewer Comment</p>	<p>Page 3: “For example, two widely used methods for predicting enhancers were based on the fact that these elements are expected to contain a cluster of transcription factor binding sites [24] and their activity is often correlated with an enrichment of particular post-translational modifications on histone proteins [27, 30].”</p> <p>In a similar fashion one can argue that the authors use STARR-seq peaks that overlap with DHS or H3K27ac peaks to identify active regulatory regions in the genome. See comment below. This requires much better justification.</p>
<p>Author Response</p>	<p><b>ANS</b></p> <p>After applying STARR-seq to identify potential enhancers in the fly genome, Stark and co-authors showed that the STARR-seq peaks that overlap with enriched DHS or H3K27ac signal in the same cell-line were close to genes that had higher gene expression. In contrast, the STARR-seq peaks that overlapped with reduced DNase hypersensitivity overlapped with the repressive mark H3K27me3 and were situated close to genes with lower gene expression in the same cell-type. While STARR-seq identifies regions that could be potential enhancers or promoters, it does not guarantee that the region will be active or repressed in that cell-type. In machine learning models, the training data should be as well annotated as possible. As our attempt is to use the cleanest set of experimentally verified enhancers that could be active in a cell-type specific fashion, we used the experimentally active STARR-seq peaks that overlapped with DHS or H3K27ac peaks as our training data as these are more correlated with active regions in the genome as per the STARR-seq study. We have clarified this in the supporting information of the manuscript.</p>
<p>Excerpt From Revised Manuscript</p>	<p>STARR-seq identifies regions that could be potential enhancers or promoters, it does not guarantee that the region will be active or repressed in that cell-type. In machine learning models, the training data should be as well annotated as possible. As our attempt is to use the cleanest set of experimentally verified enhancers that could be active in a cell-type specific fashion, we used the experimentally active STARR-seq peaks that overlapped with DHS or H3K27ac peaks as our training data as these are more correlated with active regions in the genome as per the STARR-seq study.</p>

### -- Ref2.3 – Training and test data --

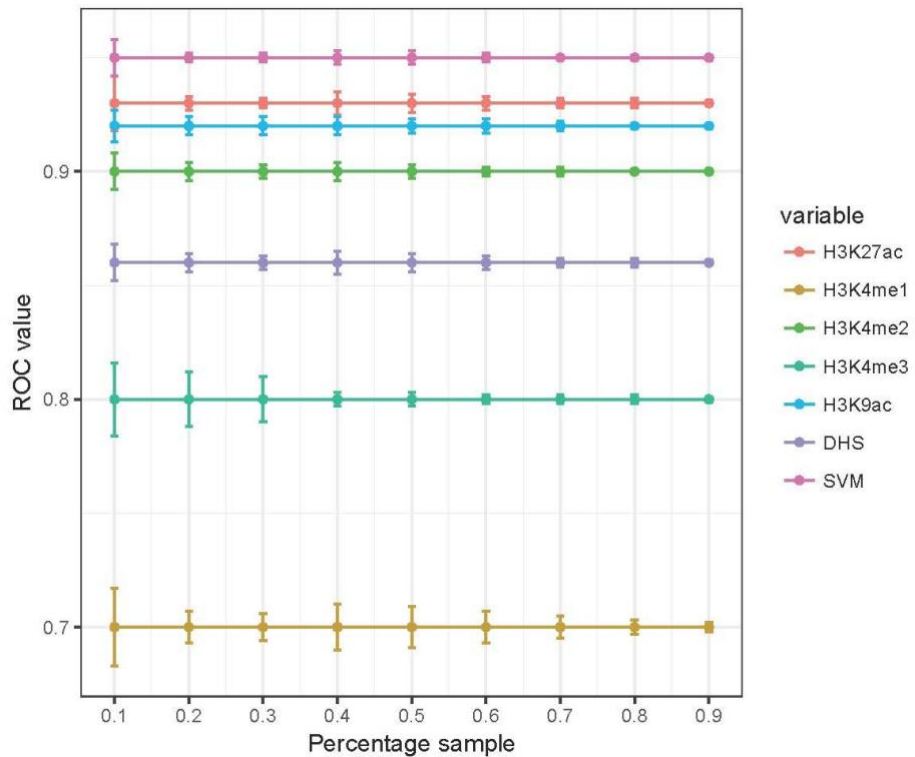
<p>Reviewer Comment</p>	<p>Page 3: "However, the optimal method to combine information from multiple epigenetic marks to make cell-type specific regulatory predictions remains unknown. For the first time, using data from several MPRAs, we have the ability to properly train our models based on a large number of experimentally validated enhancers and test the performance of different models for enhancer prediction using cross validation"</p> <p>By no means this is an optimal method. This may only be considered optimized but under very specific constraints. Most of the existing methods for the prediction of regulatory regions based on epigenetic markers such as RFECS, ChromaGenSVM, DEEP, CSI-ANN, Chromia, DELTA and others including the proposed method apply heuristic techniques to identify solutions that are close to the best possible answer. So, they are optimized. The sub-optimality of the achieved solutions using epigenetic markers is not due to the training procedure of the methods, but mainly due to the variability of the epigenetic profiles across different cells or developmental stages. However, the problem-solving technique (e.g., heuristic or analytic) is not related by any means to the proper training of the method, meaning that a method is properly trained as long the training data are completely independent from the testing. <b>Following, the previous points, the authors need to provide more evidence about the effect of the number of training samples on the performance maximization and make clear in their manuscript that the testing data are completely independent from the training.</b></p>
<p>Author Response</p>	<p>Thanks for the comment. In our original text, we didn't mean to claim that our method is the optimal method. Here, our goal is to build a framework with small number of inputs requirement to ensure that we had a widely applicable method that could be used across species. Our advantage was to use large scale STARR-seq experimental data to train the model, which was not used in previous methods.</p> <p>To demonstrate the effect of the training sample size on model performance, as suggested by the referee, we did a saturation analysis where we down-sampled the training data to different sizes. We added the result of this analysis in the supplement as reproduced below.</p> <p>For each cross-validation performed in this paper, the test dataset is completely separated from the training dataset. We have made that clear in the main manuscript and supplement as well. In addition, the many independent sources of validation</p>

performed in this paper shows that the model has good ability to generalize and has wide applications.

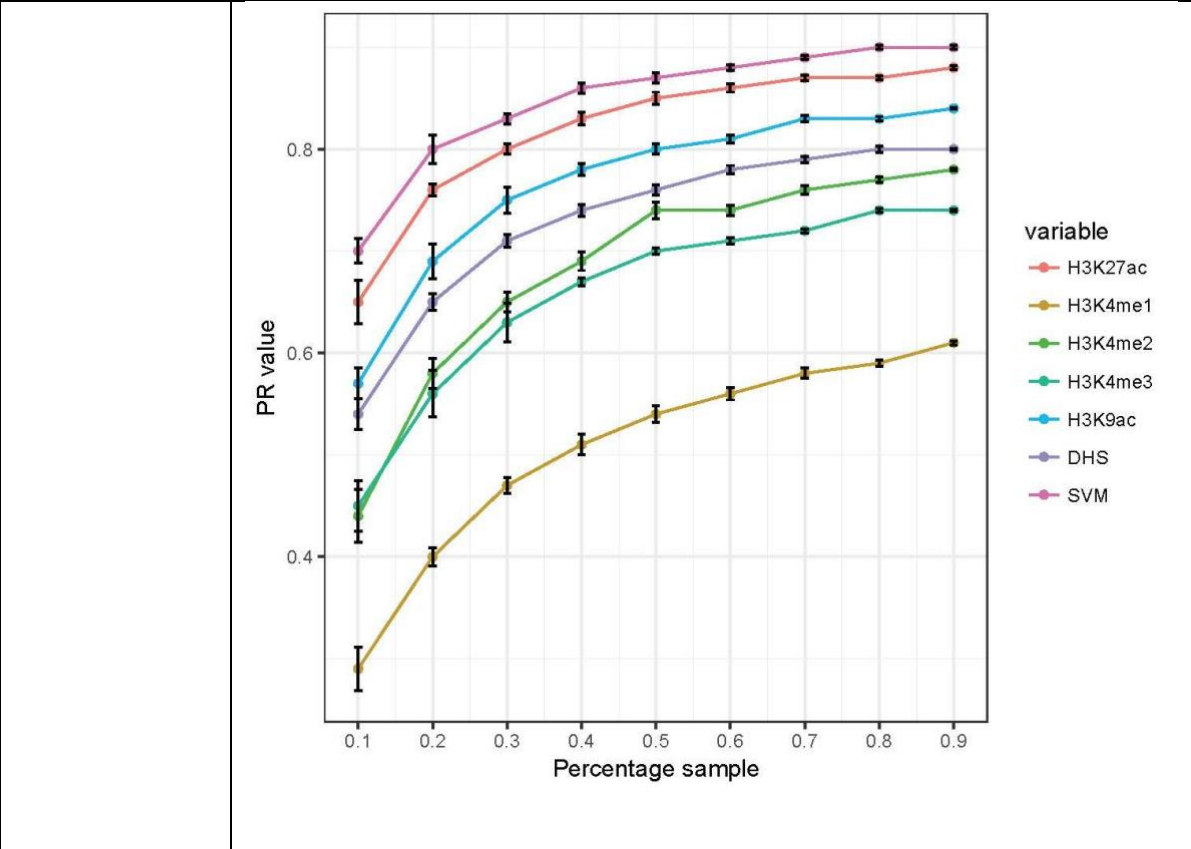
Excerpt From Revised Manuscript

To evaluate the impact of the training sample size on model performance, we did a saturation analysis where we down sampled the training data to different levels of fractions and evaluated the model performance on the remaining data. For each fraction level, we did a 10-fold cross-validation (see methods) and then took the average of the ten output result. The result shows that the average AUPR increases with increasing size of training data, and it starts to saturate for our SVM model with 80%-90% of the experimental data for training. In contrast, the average AUROC remain comparable with varying training size, but the performance variances decrease with increasing training data size.

[[ In methods section: The metaprofile and SVM models are trained on x% of samples and tested on the rest of the data, so the testing data is completely independent from the training.]]







-- Ref2.4.1 – Figure 2 --

<p>Reviewer Comment</p>	<p>Figure 2 requires more information: The authors assessed the performance of the deployed matched filter algorithm by predicting active STARR-seq peaks, and they concluded that H3K27ac is the most informative predictor. However, H3K27ac together with DHS has been used for the selection of the active STARR-seq peaks. <b>Thus, the authors should exclude those two markers and repeat the analysis without them.</b></p>
<p>Author Response</p>	<p>Thanks the reviewer for the comment.</p> <p>Thanks the reviewer for the comment. aS requested we did XXX &amp; YYY we generated a new suppl. figure it shows as expected that ... However, don't think the method we are using here is really that biased... argument blew</p> <p>As STARR-seq tested enhancer activity in an episomal fashion, in their original paper \cite{27831498} they noted that “the</p>

	complementary DHS-seq and ChIP-seq determine enhancer-associated characteristics in the endogenous genomics context". Here, we did use STARR-seq regions with H3K27ac or DHS signals which indicates these regions are active in endogenous genomic context, but we didn't require them to have any shape pattern. For our analysis in Figure 2, the 'peak-trough-peak' based the shape pattern matched filter score gives the highest predictive power. The ROC and PR curve for the other histone marks shown in Figure 2 are independently assessed just based on those histone modifications, thus H3K27ac and DHS are excluded.
Excerpt From Revised Manuscript	

**--- Ref2.4.2 – Figure 2 --**

Reviewer Comment	Another more technical comment is about usage of 10-fold cross validation. If the number of training and testing sample is large enough 10-fold cross validation is not necessary. 5-fold cross validation is sufficient or even 2-fold cross validation assuming big numbers of training and testing data (e.g., more than few thousands). Finally, there is a minor comment about the quality of Figure 2 and some other figures. In my pdf many of them appear a bit blurry.
Author Response	We thank the referee for the comment. We agree to the referee that the 5-fold or even 2-fold cross validation might be sufficient. <b>We have included in the supplementary the results for 5-fold and 2-fold cross validation result.</b> We used the original PDF of figure 2 but we apologize it looks a bit blurry upon upload. We'll make sure it is upload in the full size and is in the clear form.
Excerpt From Revised Manuscript	

**-- Ref2.5 – Feature selection --**

Reviewer Comment	I need more justification about the selection of six predictors for the development of the integrated model. <b>I agree that the selected epigenetic marker datasets are widely available for many cell-lines from publicly available resources.</b> Without doubt, this way increase the utilization of the method in new cases. <b>My question is why</b>
------------------	---

	<p><b>six and not another combination out of the 30?</b> Continuing the previous comment about optimality of the heuristically identified solutions, is there any guarantee that the integration of the selected six predictors is optimized? For example, one can apply an exhaustive search algorithm and find the best combination. One also can argue that since the performance differentiation with Random Forests is small, the latter classifier is more effective since it integrates an "out-of-bag" feature selection technique. For example, this is the biggest advantage of RFECS method that pooled together multiple epigenetic markers and identifies the most informative. <b>Authors have to elaborate more on the available dimensionality reduction techniques to select the best combination of predictors.</b> To keep it as simple as possible, combining filtering techniques such as mRMR or Gini index with the linear SVM is quite powerful and provides interpretable results.</p>
<p>Author Response</p>	<p>Thanks to the referee for the question. In our model, we chose these 6 histone marks because we wanted to test the applicability of the model trained with fly data for predicting active enhancers and promoters in mouse and human tissues. The 30 histone marks we tested are from drosophila experiments, and most of them does not have available data even in top tier tissues and cell lines for mouse and human. We didn't seek to pursue an optimal combination of all histone marks. While optimality of marks could potentially be used to identify other histone marks that provide complementary information about activity of enhancers and promoters, it could potentially reduce the applicability of the model to mouse and human tissues and cell-lines. We select the features to which are both widely available and have good individual performance. Also, we allow our model to be flexible so even one of the histone mark is missing the model still works.</p> <p>Based upon GINI index for the random forest model (Supporting Information), H3K27ac and H3K9ac are two of the epigenetic marks whose matched filters provide the best performance among the thirty marks for identifying active enhancers and promoters. In addition, H3K4me1 and H3K4me3 marks provide the ability to distinguish between promoters and enhancers (and Figure 3). In addition, DHS and H3K4me2 are also widely used within the literature to identify enhancers and promoters. The set of histone marks selected in our model is in agreement with RFECS, where H3K4me1, H3K4me2, H3K4me3 are identified as the most predictive histone marks, with H3K9ac following as the commonly available highly predictive histone mark. They also adopted H3K27ac as it is the most commonly</p>

	available histone mark with prior knowledge of being predictive for enhancers, although H3K27ac is not among the top important histone marks in their importance analysis.
Excerpt From Revised Manuscript	

-- Ref2.6 – Definition of promoters and enhancers --

### Leave out this section, to be finished

Reviewer Comment	<p>Separation of active STARR-seq peaks to promoters and enhancer based on the distance from known TSSs is the adopted practice, however it is too "quick and dirty". The truth is that, it is very difficult to discriminate sharply enhancers from promoters based on the distance from TSSs since promoters have frequently function of enhancers and vice versa, and both of them share similar transcriptional architecture and have similar properties (ref. PMID: 26073855). <b>From a technical point of view and based on the existing results, I would like to see the performance of the deployed method by varying the distance from TSS for selecting enhancers and promoters for testing.</b> In the extreme case the binary classification problem is transformed to one-class classification problem that the method should handle. <b>An alternative way is to repeat the analysis, using appropriate CAGE-defined promoter and enhancer datasets that coincide with STARR-seq peaks.</b> There are also data from studies such as "Systematic dissection of regulatory motifs in 2000 predicted human enhancers using a massively parallel reporter assay" or "High-throughput functional testing of ENCODE segmentation predictions" that could be used as baseline for benchmarking the performance of the method in a more orthogonal way.</p>
Author Response	<p>The referee is making a reasonable point we've generated a suppl figure that shwos how sensitive our calc is to enhancer promotor def'n</p> <p>#1 is reasonable.</p> <ol style="list-style-type: none"> <li>1) Take all genes in H1-hESC that are active (&gt;1TPM).</li> <li>2) Take their promoters - closest active activatory region to that gene.</li> <li>3) Histogram of distance between promoter and TSS - will most promoters be &lt;2kb from TSS?</li> <li>4) Histogram of distance between the rest of the active regulatory regions and the gene - will most enhancers be &gt;</li> </ol>

	<p style="text-align: center;"><b>2kb from TSS?</b></p> <p style="text-align: center;">Distance to Gene start</p>
Excerpt From Revised Manuscript	

**-- Ref2.7 – Comparison analysis for human cell lines --**  
**### Leave out this section, to be finished**

Reviewer Comment	<p>Page 9: “Similarly, we did genome wide prediction of regulatory regions in ENCODE top tier human cell lines, including H1-hESC, GM12878, K562, HepG2 and MCF-7 (all available through our website)”.</p> <p>Following my previous comment, I would like to see the comparison analysis with CAGE-defined enhancers and promoters for some cell-specific cases, comparison with the integrative ENCODE annotation proposed by Hoffman for all top-tier cell-lines as well as comparison with other studies (see previous papers) that validated the regulatory activity of different segments in K562, HepG2 or H1-hESC cell-lines.</p>
Author Response	<p>Thanks for the suggestion.  MTG redo</p> <p>move to suppl {{ We overlapped the CAGE-defined enhancers from FANTOM5 to our predicted enhancers and also to the</p>

	<p>enhancers predicted by the integrative ENCODE annotation method proposed by Hoffman for cell lines including GM12878, K562 and HepG2. We found that around 40% of the CAGE-defined enhancers overlap with our predicted enhancers, and only 23% to 34% of the CAGE-defined enhancers overlap with the enhancers predicted by integrative ENCODE annotation method, although the latter provides much larger numbers of predictions (about four times for GM12878 and K562, and three times for HepG2).</p> <p>For the promoter (XXXXX) (I can not find the CAGE-defined promoters for these cell lines from FANTOM5.) }}</p>
Excerpt From Revised Manuscript	MTG to send CY the promoters and compare

**-- Ref2.8 – Comparison with previous methods --**

Reviewer Comment	<p>The comparison analysis is limited to ChromHMM and Segway. However, there are more methods available such as RFECS, DEEP, CSI-ANN that provide predictions for top tier ENCODE cell-lines. I would like to see a comparison analysis similar to the one presented in Figure 5 of the RFECS paper. Are the predictions of the competitor methods supported by same TF-binding sites? This might reveal that STARR-seq peaks that overlap with specific TFs such as p300 or CBP provide a better training dataset. Related to the comparison with ChromHMM and Segway. Both ChromHMM and Segway are based on probabilistic graphical models (HMM and Bayes). They should include a method of different type for example using SVM or Random Forest that is more close to what they have been developed.</p>
Author Response	<p>We compared with ChromHMM and SegWay as their enhancer annotation has been used in many publications as a way to define enhancer regions. Based on the referee's suggestion, we also did more comparison with other published methods, and we have included the results in our manuscript as shown below.</p>
Excerpt From Revised Manuscript	<p>In addition to the comparison with unsupervised segmentation based methods, we also compared with other published enhancer prediction tools, including CSIANN, a neural network based approach; DELTA, an ensemble model integrating different histone modifications; RFECS, a random forest model based on histone modifications, and REPTILE, a more recent published method that integrates histone modifications and whole genome bisulfite sequencing data. We show that our method also outperforms these previous approaches assessed by in vivo transgenic experiments.</p>

