

I

INTEGRATIVE MODELLING AS A TOOL TO EXPLORE BIOLOGICAL SYSTEMS

1.1 INTRODUCTION

Large-scale chromatin data has recently been produced by multiple consortia, most notably ENCODE^[1] (Section ??) but also the NIH Roadmap Epigenomics project.^[2] The breadth and depth of this new data offers unprecedented opportunities to advance our understanding of the complex biology of the chromatin landscape. To this end, studies have already enjoyed success in integrating these data through modelling techniques, with the subsequent dissection of these models revealing novel insights into complex biological phenomena.

Recent studies have shown convincingly that local chromatin state measurements can accurately predict expression levels of genes on a genome-wide basis. Tippmann *et al.*^[3] designed a linear model to predict steady-state mRNA levels in mouse embryonic stem cells based on just four predictors: 3 histone modifications (H3K36me3, H3K4me2 and H3K27me3) and Pol-II occupancy. Remarkably, the linear model was found to explain 84.6% of an estimated 91% maximal variance that could be explained (as calculated through a detailed determination of noise).^[3] An additional finding of this study was that mRNA half-life and microRNA mediated transcript degradation both had relatively minor influence on steady-state mRNA levels, with the authors concluding that “the lion’s share of regulatory contribution is at the level of mRNA synthesis and predictable from chromatin alone.”^[3] An independent study used a similar regression modelling approach to chromatin and transcription factor data and again concluded that models built with histone modifications and chromatin accessibility data were almost as accurate as those which also included binding data for 12 transcription factors.^[4]

A key study from the ENCODE consortium, that of Dong *et al.*^[5], used ChIP-seq datasets to predict gene expression in a range of cell types as measured by a variety of experimental techniques. The authors here developed a two-stage model which first attempts to classify each transcription start site (TSS) into an ‘on’ or ‘off’ state using a powerful ensemble classifier technique called Random Forests (RF). The second stage of the model used the same range of histone modifications as regressors in a simple linear modelling framework to quantify predicted expression. This approach proved very successful, producing a median Pearson correlation coefficient (PCC) between predicted and empirical expression levels using 10-fold cross-validation of 0.83 across all cell lines and expression level technologies.^[5] Additionally, this study highlighted

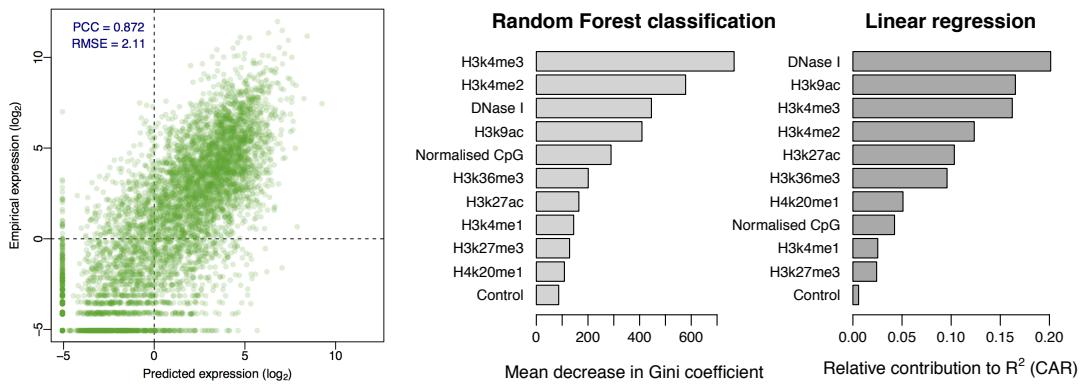


Figure 1: Highly accurate models of expression were built following Dong *et al.* A scatterplot shows the strong correlation between predicted and observed expression levels per transcript (*left*). Variable importance metrics are shown for the on/off RF classification step and continuous linear regression step (*right*).

cap analysis of gene expression (CAGE) as the technology, relative to RNA-seq and RNA-PET, which produced the most predictable expression response. CAGE uses 5' capped transcripts to generate short, specific tags which precisely identify TSS positions as well as quantifying the abundance of a given transcript.^[6,7]

These recent publications highlight the importance and relevance of advancing our understanding of chromatin biology through a model-based approach. We can extend this idea to the related domain of nuclear architecture, in the hope of understanding the relationships between chromatin and higher order structure in the same way that chromatin features have been related to transcriptional output.

1.2 EXTENDING DONG *et al.*

We reimplemented the published modelling framework of Dong *et al.*^[5] to replicate their results and analyse the strengths and caveats of their approach.

We were able to reproduce the reported results and generate highly accurate models of transcriptional output. An example is shown for a model of CAGE expression in the H1 hESC cell type (Fig. 1). Note that not all variables used in Dong *et al.*^[5] were made available for this particular modelling scenario, however the Pearson correlation between predicted and observed expression (0.87) is above the study's median value (0.83), and in-line with other models predicting CAGE data (median $PCC \approx 0.87$).^[5]

We also re-calculated measures of variable importance for this model of transcription (Fig. 1). We note some small variations from the example model shown in Dong *et al.*^[5] which was predicting cytosolic CAGE data recorded in K562 cells. This hints at some degree of some variability between cell type models, though broad

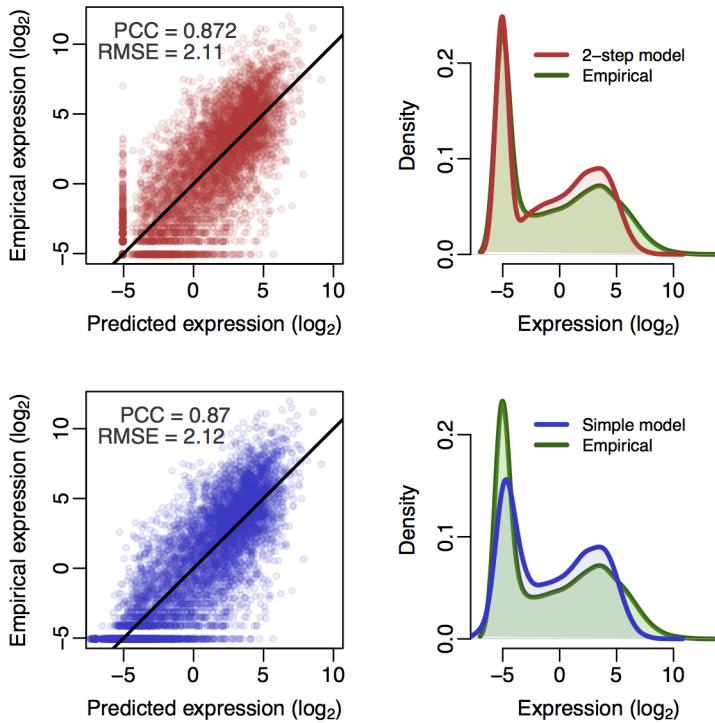


Figure 2: Comparison of a published two-step classification-regression model of transcription with a simple linear regression model. Scatterplots of predicted against empirical \log_2 reads per million (RPM) expression values for the two-step model of Dong *et al.* [5] and simple multiple linear regression are shown (*left*) along with frequency distributions of predicted and observed expression levels (*right*). Scatterplots are annotated with Pearson’s correlation coefficient (PCC) and the root mean squared error (RMSE); black trendlines describe $y = x$. Overall correlations calculated with 10-fold cross-validation.

similarities also exist such as both models ranking DNase I hypersensitivity as a relatively informative variable (Fig. 1).

When replicating the modelling approach of Dong *et al.* [5], we were surprised to find that the two-step classification then regression (firstly assessing a gene as ‘on’ or ‘off’ and then predicting its expression level) added little additional accuracy relative to a simple linear regression model (Fig. 2). Indeed, it appears the “best bin” technique (Section 1.2.1) had much greater impact on overall predictive power than the addition of this classification step.

1.2.1 Bestbin method

An innovative element of the modelling approach used in Dong *et al.* [5] is the ‘bestbin’ method of matching chromatin measurements to the expression of a given TSS. This strategy first bins normalised signal intensities into 40×100 bp bins encompassing 4 kb around the TSS, and adds an additional bin representing the remaining gene

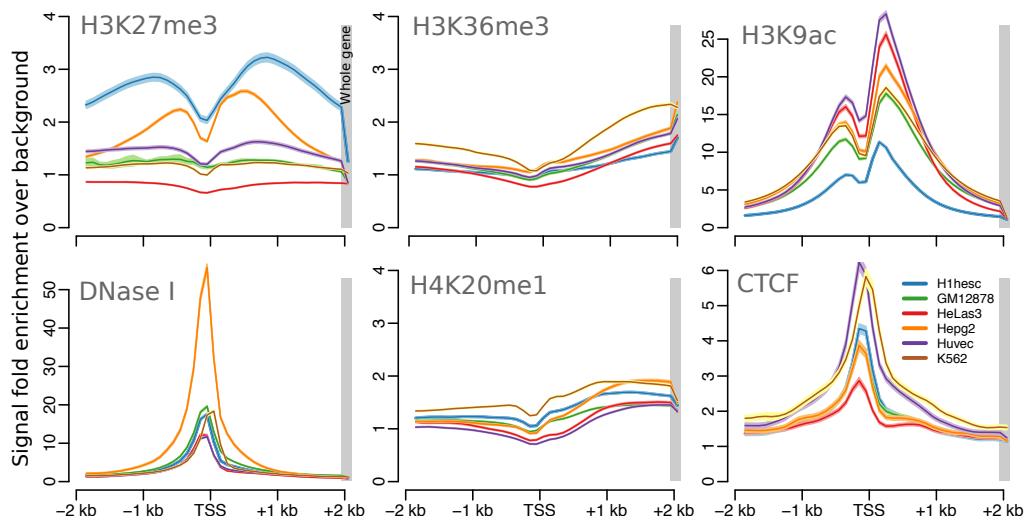


Figure 3: Average input feature profiles over transcription start sites. Mean ChIP-seq signal over input control is shown for 6 factors in 6 human cell types used in Dong *et al.* [5]. Each is averaged genome-wide over GENCODE v7 hg19 defined TSS ± 2 kb, and over whole genes (grey shading). Ribbons shown 99% confidence intervals of the mean.

body. Then the correlation between the signal of a given mark and the expression of a TSS across all genes is measured, then the bin producing the highest correlation is designated as the ‘bestbin’. That bin’s normalised ChIP-seq signal intensity is then used as an input feature for training the model of transcription. This strategy was shown to increase model performance, measured in terms of the Pearson correlation between predicted and observed expression, by 0.1 in the simple regression model, an increase of almost 13% relative to simply taking the average value across all bins. [5]

The justification for such an approach hinges on the idea that the multitude of input features (mostly histone modifications and DNA binding proteins) have a variety of biological functions, and so the bestbin method is one way of learning these functions in an automated and unbiased way. For example, the histone modification H3K36me3 is understood to be painted across exons that are being actively-transcribed, [3,8,9] thus the genome-wide summary statistic that best captures this function is likely the whole-gene measurement, rather than the level of H3k36me3 at a gene’s TSS or upstream promoter. A re-analysis of ENCODE data used in Dong *et al.* [5] highlights this kind of variability across input features (Fig. 3). Some marks clearly are enriched directly over TSS (CTCF, DNase; Fig. 3) while others show enrichments along the gene body (H3K36me3, H4K20me1; Fig. 3) as well as more-complex, asymmetrical “shoulder” patterns (H3K27me3, H3K9ac; Fig. 3). Bestbin will therefore, to some degree, capture these spatial relationships without *a priori* specification.

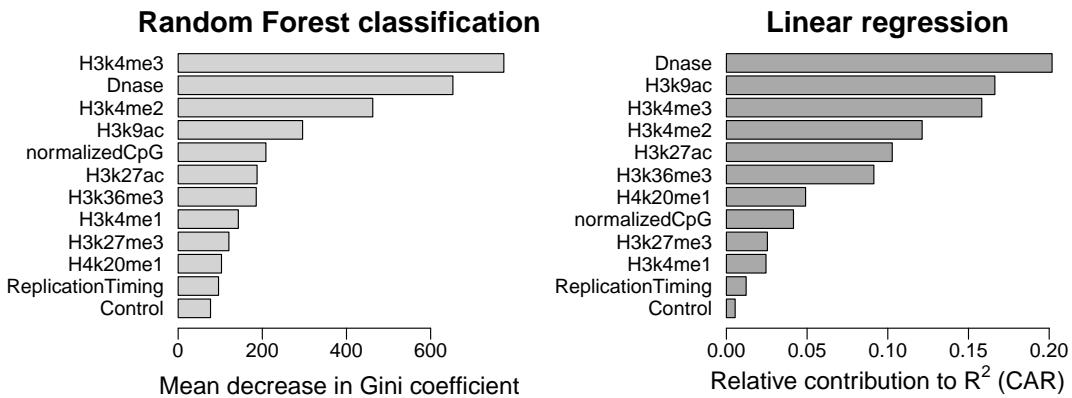


Figure 4: Relative importance metrics for variables in both stages of a reimplementation of a published model for predicting transcriptional output. Variable importance is measured by decrease in Gini coefficient for the RF classification step, and by CAR R^2 decomposition^[10] for the linear regression step.

1.2.2 Model exploration

We attempted to improve the accuracy of predicted expression values produced by Dong *et al.*^[5] through increasing the number of informative regressors. While Dong *et al.*^[5] included broad coverage of different histone modifications, they did not investigate the impact of higher order chromatin data. For this reason, we matched the TSS positions used in Dong *et al.*^[5] with previously-published genome-wide replication timing ratios measured in BGo2 ESCs.^[11] This data is of a different origin to the transcriptional data in this case (which was recorded in H1 hESC) but replication timing is thought to be largely conserved between cell types.^[12]

We then used these values as an additional regressor in both the two-step classification regression model and the simple linear model but saw no significant improvement in either model's accuracy (*data not shown*). The reasons for this are likely that the data were relatively low-resolution (1 Mb), from an imperfectly matched cell line and also that the existing model is already achieving such accurate results that they must already be accounting for most of the maximal explainable variance in gene expression given experimental and biological noise. With this in mind, additional regressors would be expected to yield diminishing returns. Even so, on closer examination, the replication timing data appeared only slightly more informative than the control ChIP-seq input measurements when evaluated with relative importance metrics (Fig. 4), implying that large-scale chromatin domains have relatively little influence on the expression of the genes resident within them.

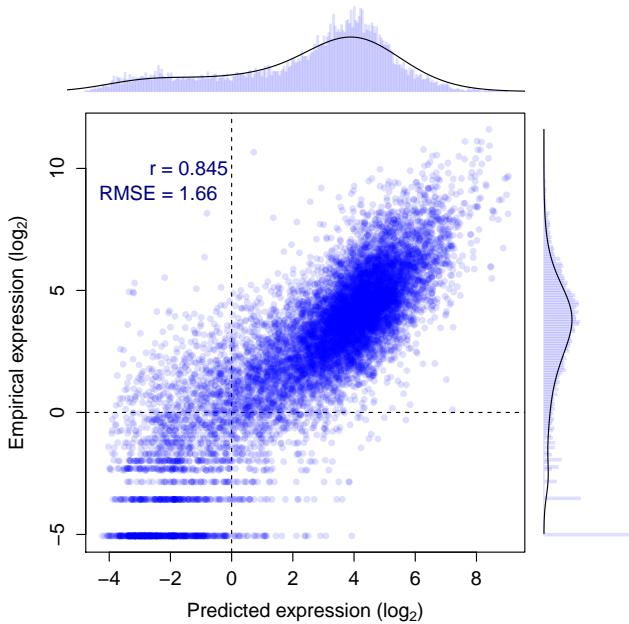


Figure 5: Random Forest predictions of FANTOM5 expression data. RF model predictions are plotted against their empirical values. The marginal distributions of predicted and empirical expression values are shown opposite their respective axes. Pearson’s correlation coefficient (r) and the root mean-squared error (RMSE) are shown (*inset*).

1.3 MODELLING FANTOM5 EXPRESSION DATA

Using FANTOM5 CAGE data^[13] and the approach established above (Section 1.2), we next attempted to model gene expression at timepoint zero (t_0) of a differentiation timecourse of human H1 embryonic stem cells (H1 hESC) to CD34+ hematopoietic stem cells. Applying this modelling strategy to a novel dataset will allow us to assess the portability of the model design, as well as enabling further analysis of model components such as the bestbin strategy.

We retrieved a number of genome-wide ChIP-seq datasets measured in H1 hESC cells and produced by the ENCODE consortium^[1] (Methods ??). These were matched to transcript annotations to build an input feature set for use in building a predictive model of transcriptional output.

Due to the finding that a two-step (classification–regression) approach added little additional modelling accuracy (Fig. 2), we employed a single-step design using a Random Forest (RF) regression model.^[14,15] With a total of 14 predictors (10 histone modifications, HDAC6, H2A.Z, DNase I and an input control, listed in Methods ??), we were able to build a highly accurate predictive model of transcriptional output of around 11,000 TSS (Fig. 5).

Model predictions evaluated with 10-fold cross validation show a highly significant correlation with measured CAGE levels ($PCC = 0.845 \pm 1 \times 10^{-4}$, $p < 2 \times 10^{-15}$), and

the model is able to explain around 71% of the variance in the expression response (for comparison a linear model resulted in $PCC = 0.825 \pm 3.2 \times 10^{-5}$, $p < 2 \times 10^{-15}$). This result is less impressive than that of Dong *et al.*^[5] who achieved cross-validated correlation coefficients of up to 0.9, but it is roughly equal to their median test set correlation of 0.83.^[5] Though the RMSEs of our predictions, when normalised by the range of observed values, compare more favourably (0.11, compared with Dong *et al.*'s: 0.14).

Our slightly lower predictive power could be explained by our streamlined model design. Dong *et al.*^[5] implemented a pseudocount optimisation step whereby an additional count added to each binned signal intensity prior to log transformation to maximise expression correlation. In the model presented above, a fixed pseudocount of 1 was used to avoid introducing an unwarranted positive bias towards higher correlation. We confirmed that a two-step classification–regression design did not improve our model performance metrics; indeed, the PCC and RMSE of a classification–regression framework with this data showed a slight decrease in prediction accuracy ($PCC = 0.834 \pm 0.007$, $RMSE = 1.77$ when applied to the same test and training data used in Fig. 5).

1.3.1 Bestbin location

We again implemented the previously-described ‘bestbin’ strategy^[5] (Section 1.2.1) to objectively select the most-correlated binned signal for each chromatin H1 hESC mark. To explore the implications of this approach, we analysed the stability of chosen bestbins by calculating them on 200 sets of 1000 randomly selected TSS samples, with each sample representing approximately 8% of the complete dataset. Distributions of chosen bestbins across these 200 sub-samples are shown as boxplots (Fig. 6).

We find that bestbin selections are often highly consistent across sub-samples, indicating there are fairly static informative regions relative to a TSS for each chromatin feature. Furthermore, the selected bestbins match known biological mechanisms; for example the H3K36me3 mark’s bestbin is consistently the whole gene measurement (Fig. 6) and this mark is known to be enriched in actively transcribed exons.^[3,8,9] The negative control variable (ChIP-seq input) shows no strong location bias, as expected (Fig. 6). Other distributions are less easily explained, such as those features showing a tight distribution of informative regions slightly downstream of the TSS (H3K9ac, H4Kme2/3 and H3K79me2; Fig. 6). In the case of H3K9ac, we note that the selected bestbin appears to coincide with the highest summit of its bimodal average profile over all TSS (shown in Fig. 3).

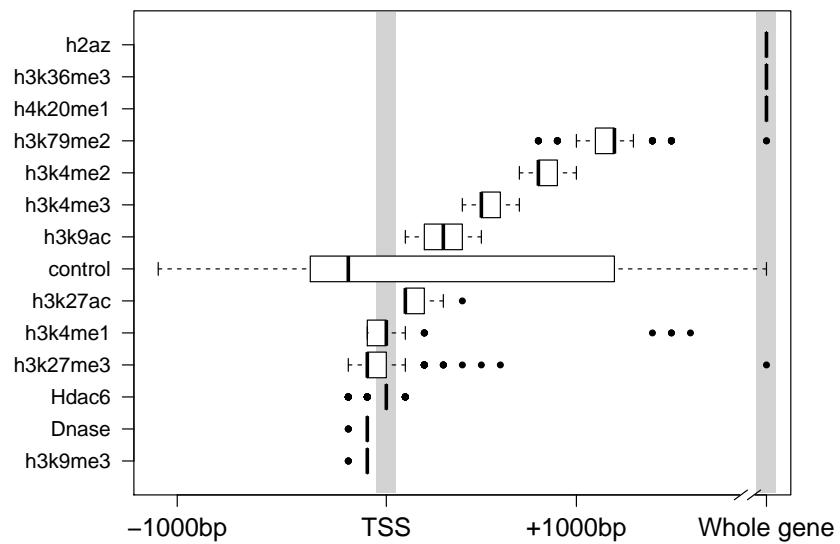


Figure 6: Distributions of bestbin locations relative to the TSS. Bestbins were selected for normalised ChIP-seq signal intensities for 10 histone marks, the H₂A.Z histone variant, Hdac6 histone deacetylase, DNase hypersensitivity and a ChIP-seq input chromatin control. Bins analysed extended 2 kb flanking the TSS, but more distal bins were never selected and hence are not shown. 'Whole gene' represents the averaged signal intensity from TSS to transcript end site, as defined by Ensembl Genes v69.

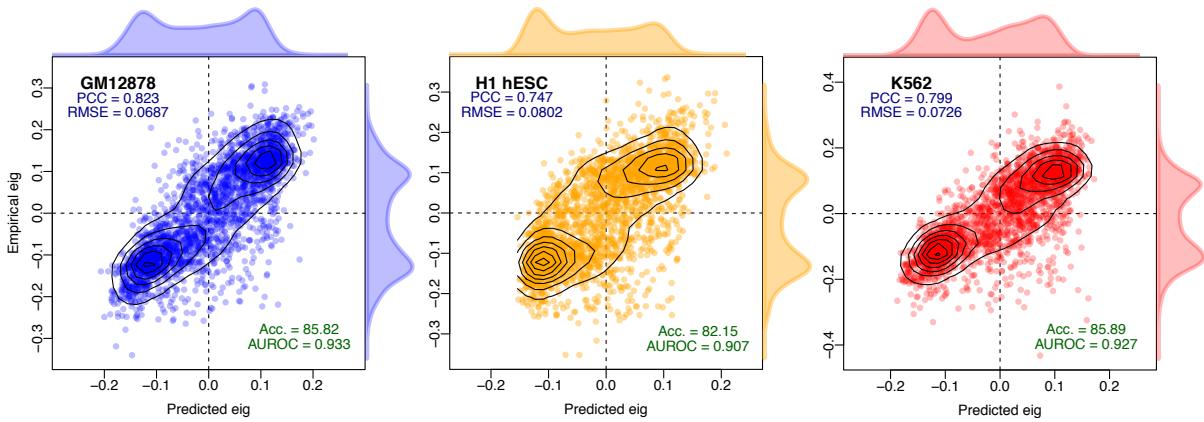


Figure 7: Compartment eigenvector model predictions are highly correlated with observed values. Pearson correlation coefficient (PCC) and root mean-squared error (RMSE) report the degree of success of the regression model, whereas accuracy (Acc.) and area under the receiver operating characteristic (AUROC) give the classification accuracy of binarized outcomes.

1.4 MODELLING HIGHER ORDER CHROMATIN

Accurate predictive modelling of transcription in a variety of cell types offered several novel insights into the interrelationships between locus-level chromatin features and transcriptional machinery, as well as advancing a quantitative explanation of the degree to which correlated features are informative. It is of interest then, to test whether this approach can be applied to other data, such as the reprocessed higher order chromatin organisation data assembled in this work (Chapter ??).

Previous publications have identified several correlates which track compartment eigenvector profiles to varying degrees,^[16,17] yet to date these relationships have not been quantitatively investigated. The above-described modelling framework offers a statistical approach towards understanding the drivers of these observed correlations.

1.4.1 Predictive model

We built Random Forest (RF) regression models (Methods ??) to predict compartment eigenvector profiles genome-wide in three human cell types. Models were found to have high predictive accuracy, with Pearson correlation between predicted and observed compartment eigenvectors in the range of 0.75–0.82 (Fig. 7), comparable to that achieved by Dong *et al.*^[5] in the prediction of transcription.

Our predictive models were also assessed in terms of classification performance, i.e. did the model correctly assign each block to an A or B compartment. Instead of training a classifier, thereby constructing a second model, we threshold our regression predictions at 0 (Methods ??). We found our RF models achieved high classification

accuracy with $\geq 82\%$ of all 1 Mb genomic bins correctly assigned in each cell type (Fig. 7).

This predictive performance underlines the strong connection between locus-level features and higher order chromatin structure previously noted by Lieberman-Aiden *et al.*^[16] Given such highly-predictive models can be generated, it is then of interest to dissect said models in an attempt to understand the nature of this captured relationship.

1.4.2 Cross-application

High predictive accuracy on cell type specific models could be the result of “overfitting”. In machine-learning, overfitting refers to the point at which parameters are being optimised to capture the random errors of a particular sample, on top of any underlying relationship between inputs and predictions, thereby giving an overoptimistic model performance which would not generalise to another featureset with an independent noise profile.^[18,19]

To test if overfitting was responsible for our high modelling accuracy, we cross-applied models learnt in one cell type to unseen input data from each of the other two cell types under study. If predictive accuracy is a lot lower on unseen data, this lends evidence to the idea that our models may be overfit to their respective cell types. Conversely, it could be the case that biologically-distinct mechanisms are in place that differ between cell types, preventing a simple cross-application.

We found cross-application between cell types was possible and resulted in similarly-high levels of accuracy to within cell-type cross-validation (Fig. 8). This gives good evidence not only that models are not overfitting to cell-type specific noise, but also that there exist broad rules linking chromatin conformation and locus-level features. Model performance under cross-application suggests there exists enough commonalities for compartment profile predictions to transcend the cell-type specific biology inherent to an embryonic stem cell or differentiated lymphoblast.

1.4.3 Between-cell variability

Given much of the higher order chromatin organisation is conserved between the three cell types used in this work (Fig. ??), a testable hypothesis is that these conserved regions are drivers of cross-applicability between cell types. Under this hypothesis we might also expect those genomic regions which vary most across cell types to be more difficult to predict.

Indeed we found the most variable regions across cell types were then most difficult to predict through our RF modelling framework (Fig. 9). In each cell type, the third of the genome with the most consistent compartment eigenvectors across

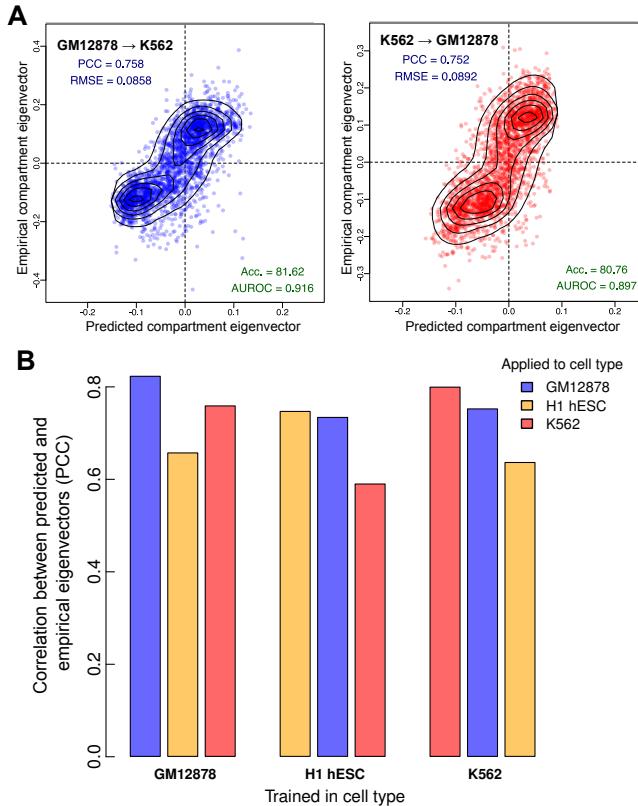


Figure 8: Models of higher order chromatin structure learned in one cell type can be cross-applied to others. Each model, trained in one cell type, was applied to the chromatin feature datasets from the other two cell types. (A) The GM12878 model achieved high accuracy when applied to K562 features ($PCC = 0.76$), as did the reciprocal cross ($PCC = 0.75$). (B) In each case, predictive accuracy decreased on cross-application but there remains significant agreement between predicted and empirical values. Acc., accuracy; AUROC, area under the receiver operating characteristic curve; PCC, Pearson correlation coefficient; RMSE, root mean-squared error.

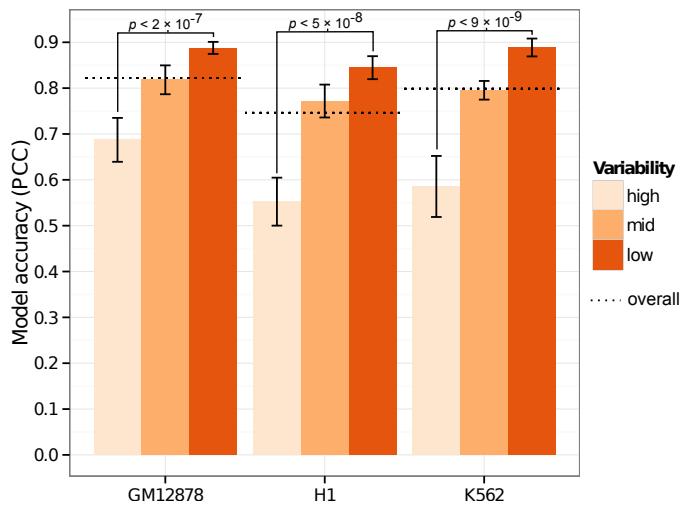


Figure 9: Genomic regions that vary across cell types are modelled less successfully than static regions. Genome-wide compartment eigenvectors were partitioned into thirds according to their median absolute deviation (MAD) across the three cell types under study. Models were fit independently to each third, and the modelling accuracy is compared.

cell types could then most accurately be modelled in each cell type, and conversely the most variable third shown significantly depleted predictability (Fig. 9). This result suggests these variable regions could either be those which are noisiest, where the eigenvector is least capturing compartment structure, or where cell-type specific biology is influencing compartment structure in ways not captured by our input feature set and low resolution modelling pipeline. Results presented in Section ??, showing that regions of variable structure are enriched for cell type specific enhancers and transcription, is suggestive of this latter explanation.

1.4.4 Variable importance

Having built accurate predictive models, we next dissect the relative variable contributions made from our range of input features and compare these across cell types. An overview of the top 10 most highly-ranked features in cell type specific models shows some agreement but also notable differences between cell types (Fig. 10)

Only one input feature, H3k9me3, is present in the top 10 most important variables of each model (Fig. 11). H3k9me3 is one of the few features to be negatively correlated with compartment eigenvectors, hence offers orthogonal information to the majority of other, positively-correlated input variables (Fig. 12; Section 1.4.5). Of those important variables shared between two cell type models, H3k27me3 is also a repressive mark and deposited by polycomb repressive complex 2 (PRC2)^[20] while H2A.Z is a histone variant again linked to polycomb-regulated genes and essential for embryonic development.^[21] Furthermore EZH2, the catalytic subunit of PRC2,^[22] is also included in the feature set and is highly ranked in the GM12878 cell type model

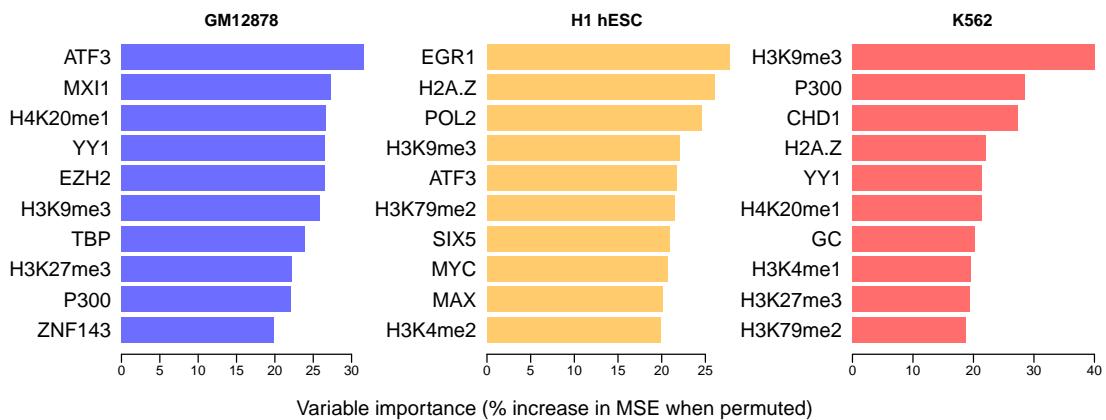


Figure 10: Variable importance per cell type specific model. Variable importance for each Random Forest model was measured in terms of percentage increase in mean squared error when permuted (Methods ??) and the top 10 ranking variables are shown for each model.

(Fig. 10). Other interrelated and important variables include MYC and MAX, which are found in the top 10 influential variables in H1 hESC, and MXI1, found to be an informative variable in GM12878. Recent results suggest MYC binds open chromatin as a transcriptional amplifier in embryonic stem cells, [23,24] with MAX and MXI1 acting as antagonistic co-regulators. [25] These biological relationships between variables may help explain the observed differences between models: different representatives of correlated clusters of input variables are likely being selected in each model (see Section 1.4.5).

To assess the significance of observed intersections (Fig. 11), the variable selection process could be modelled with, for example, a multivariate hypergeometric distribution or via simulation. Simulation was used here for simplicity: each intersection was calculated under 10,000 variables draws with uniform distribution and empirical p -values were then calculated accordingly. Under the assumption that variables are ranked independently in each cell type, drawing at least one variable in all three cell types would be expected by chance ($p = 0.6$). Similarly, the overlaps between pairs of cell types is within the range of expectation (probability of 7 or more variables appearing in exactly two sets: 0.39). Hence these data suggest the top 10 most influential variables are not significantly more alike across the three cell-type specific models than expected by chance, however ten is an arbitrary cutoff, and many of the rankings are based on small differences in variable importance, thus could be unstable between multiple generations of stochastic RF models.

In addition to rankings, raw variable importance metrics can also be compared between cell-type specific models (Fig. 13). Through this analysis we found that variables such as CTCF have a relatively small but highly consistent variable importance across the three cell type specific models, whereas other features like ATF3 are highly influential in one cell type but not the other two. Absolute differences in these figures

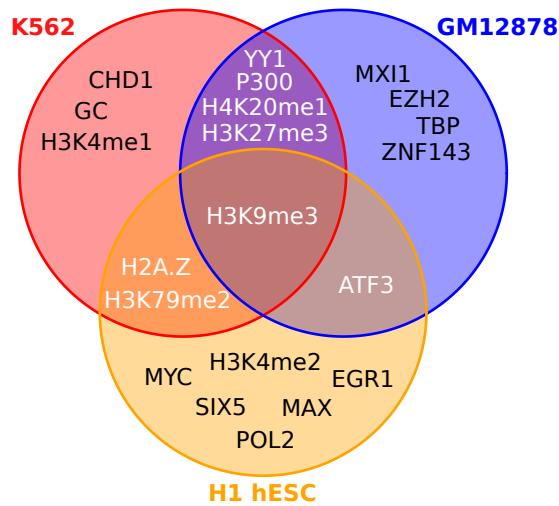


Figure 11: Intersections of the top 10 ranked variables in the cell type specific models. Venn diagram illustrating intersections between sets of ten most influential variables per cell type specific Random Forest regression model of compartment eigenvector (Fig. 10).

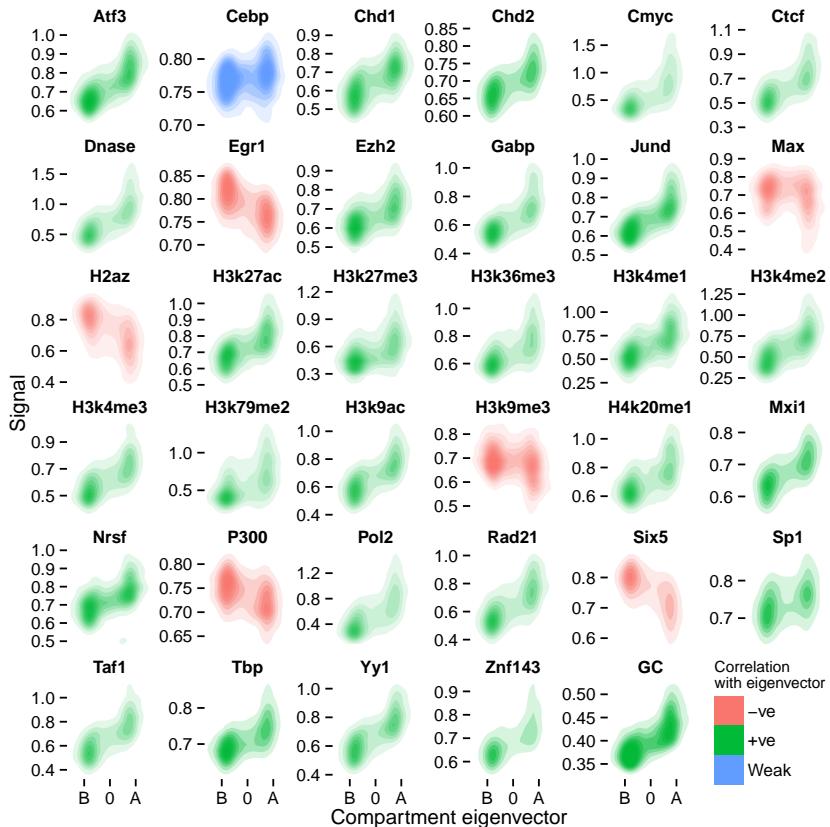


Figure 12: Correlations of individual features with compartment eigenvector in the H1 hESC cell type. Two-dimensional kernel density estimates show the density of points in a scatterplot of compartment eigenvector (x-axis) against each input feature individually (y-axes). Features with a PCC against eigenvector of above or below 0.1 are coloured as positive or negative, respectively.

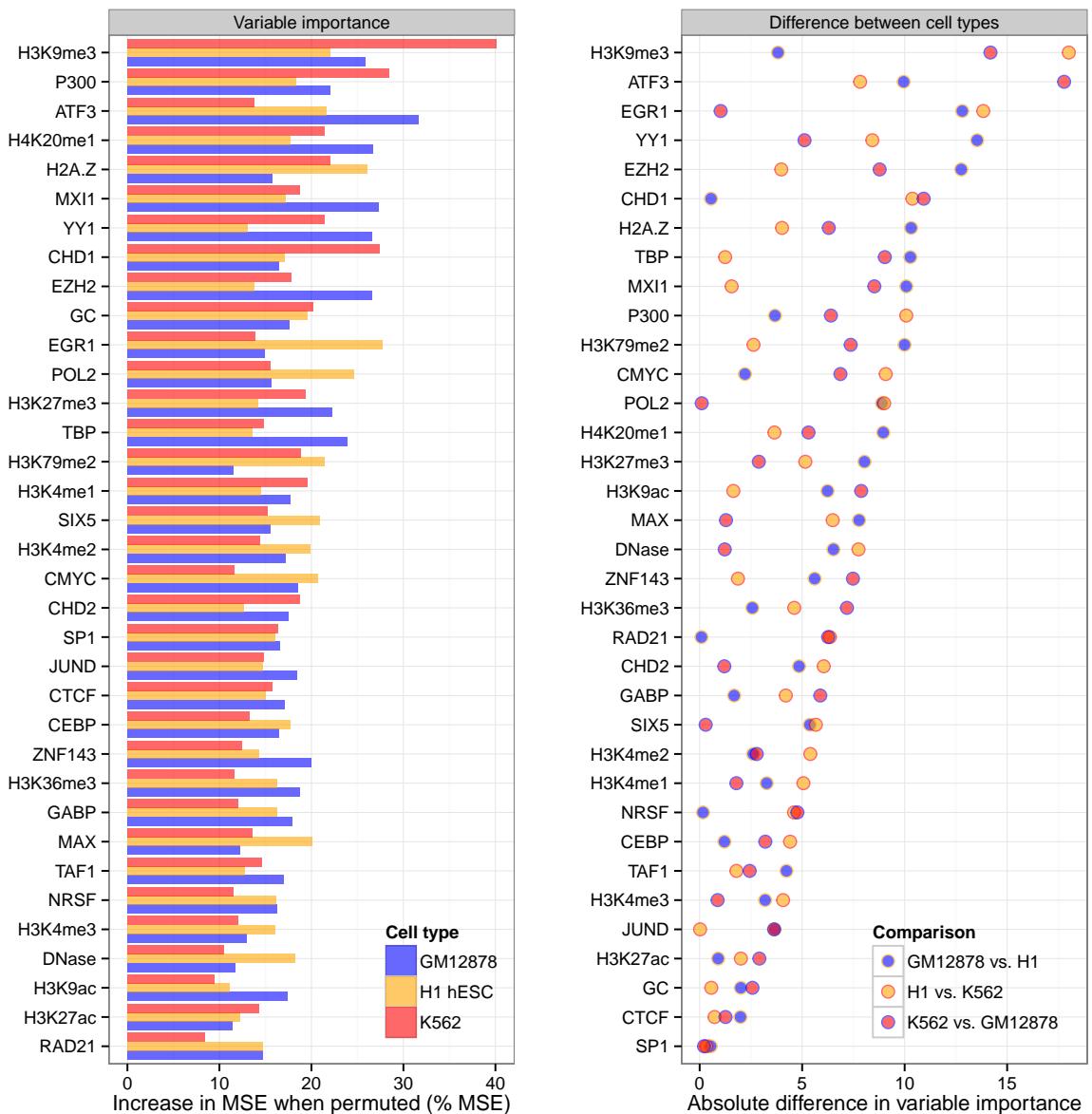


Figure 13: Comparison of variable importance between three cell type specific Random Forest models. Variable importance for each Random Forest model was measured in terms of percentage increase in mean squared error when permuted (Methods ??). Results are shown sorted by mean variable importance (*left*) and by largest absolute difference in pairwise comparisons (*right*).

should not be over interpreted and will be affected to some degree by data quality, eigenvector calculation and other sources of noise. Nevertheless there are observations which may reflect biological phenomena, such as the higher relative importance of P300 in both hematopoietic cell line models, potentially reflecting its activity as a histone acetyl transferase that regulates hematopoiesis,^[26] compared with the more consistent influence of CTCF in each model, a transcription factor widely known as a regulator of genome architecture (discussed in Section ??).

1.4.5 Correlating input features

We have an *a priori* expectation of multicollinearity in our feature set, for example between those that each broadly correlate with transcriptional activity (including POL2, H3K36me3 and sequence GC content). To explore these relationships, we performed unsupervised clustering of our feature sets in each cell type (Fig. 14).

We found pervasive multicollinearity across our feature sets, with the majority of input variables in each model falling into a persistent "active" cluster containing regions with high DNase hypersensitivity, POL2 binding and histone modifications H3K36me3 as well as GC content (Fig. 14).

Outliers are also present. H3K9me3, noted for high variable importance in each model (Fig. 10) and the only feature ranked within the top 10 in each model (Fig. 11) is a clear outgroup in the H1 hESC and GM12878 correlation heatmaps, and in K562 forms a stable cluster only with the P300 transcription factor (Fig. 14). This suggests H3K9me3 is providing orthogonal information to many of the other input variables, and likely explains its high variable importance.

1.5 TECHNICAL CONSIDERATIONS

1.5.1 Resolution

Thus far models were built at 1 Mb resolution, but if we are capturing true biological relationships we would expect these to hold at higher or lower resolutions. To test this, models leaned at 1 Mb resolution were applied to feature sets binned at 100 kb, an order of magnitude higher resolution.

Model accuracy when applied to higher resolution input features proved to be similarly high, with empirical PCC being 88 to 95% as high as that at 1 Mb native resolution (Fig. 15).

Note however, there is some indirect leakage between test and training set when 100 kb bins have been used in aggregate in learning the 1 Mb models. Nevertheless,

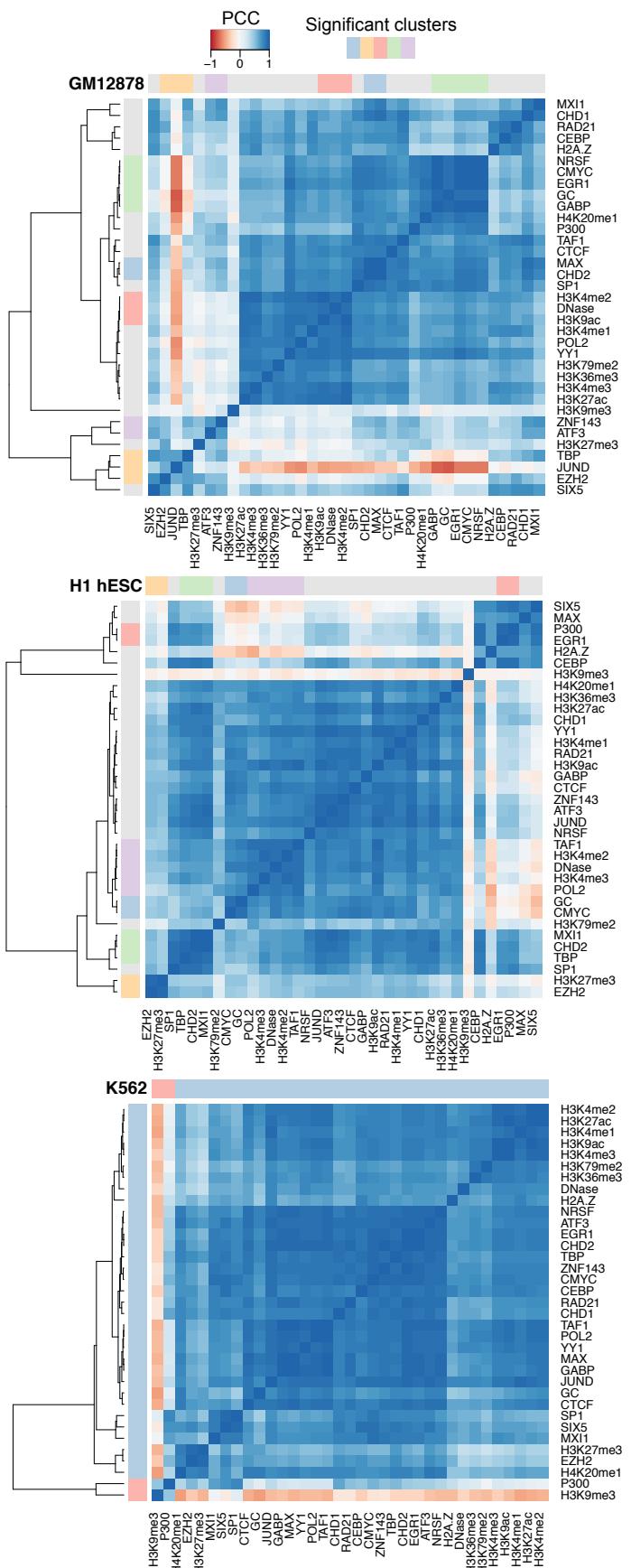


Figure 14: Correlation heatmaps of the 35 features used to model compartment eigenvectors. The Pearson correlation coefficient (PCC) of genome-wide 1 Mb bins of each feature were pairwise correlated with each other. The features were also clustered using hierarchical clustering. The significance of these clusters was determined through multi-scale bootstrap resampling, with those clusters that were stable across different sizes of resampling deemed significant, as implemented in the *pvclust* R package. [27]

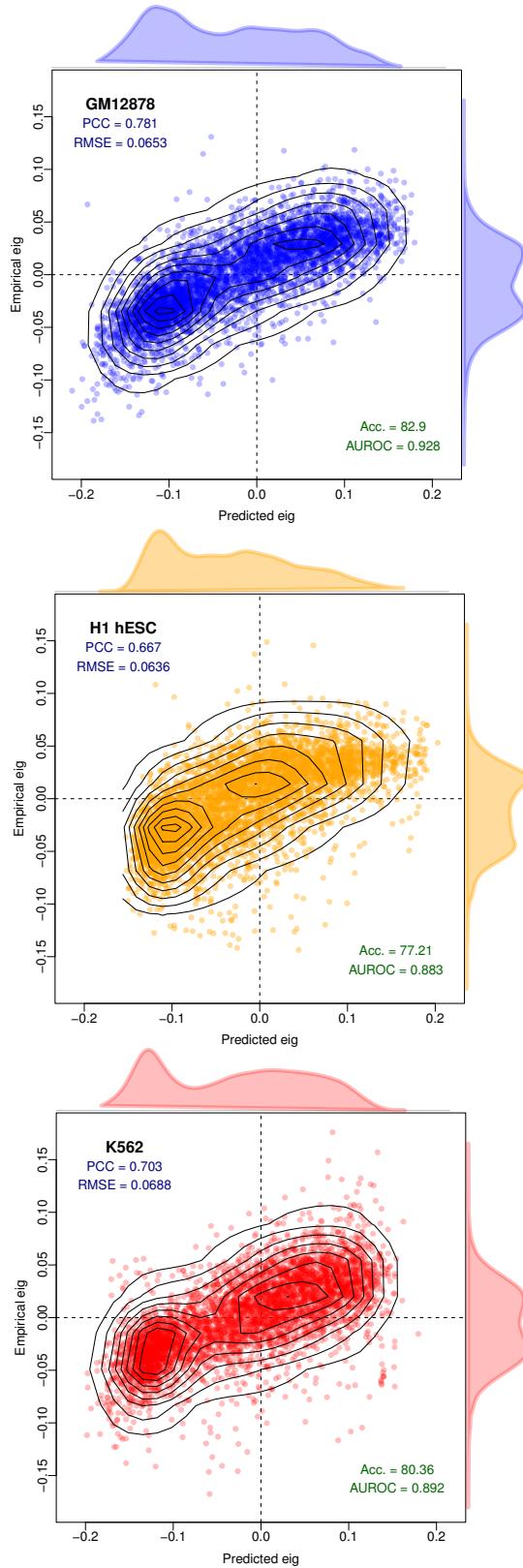


Figure 15: Models learned at 1 Mb resolution can be applied to higher resolution datasets. Despite having been trained on low resolution training sets, the Random Forest models generated can successfully predict compartment eigenvectors at higher resolution (100 kb, a 10 \times zoom). Eigenvectors at a higher resolution than this do not necessarily reflect A/B compartmentalisation.

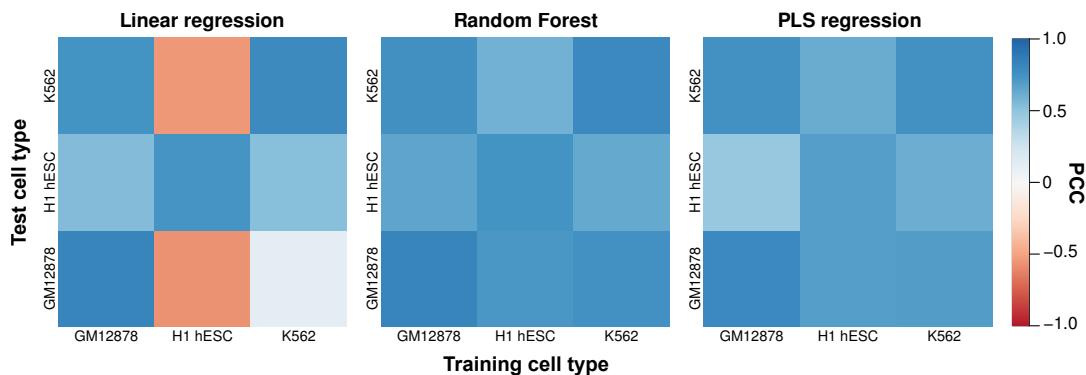


Figure 16: Comparison of Random Forest performance with other modelling approaches. Heatmaps show the Pearson correlation coefficient between predicted and observed compartment eigenvectors genome-wide for three regression techniques: multiple linear regression (LM), Random Forest (RF) and partial least squares (PLS). Results are summarised in Table 1.

sustained accuracy is evidence that our models are not resolution-sensitive, and could likely be applied to higher resolutions than the 1 Mb predominantly used in this work.

1.5.2 Other modelling approaches

Random Forest (RF) was *a priori* chosen as an appropriate and powerful modelling tool for this work. Other methods could have been used and should be compared. Here we compare our RF approach with two other options: multiple linear regression and partial least squares regression (Methods ??).

Our results confirm RF as a suitable and powerful approach for modelling our relationships of interest in this work (Fig. 16), with both the highest cell-type specific performance (PCC between predicted and observed = 0.790) and on cross-applications (mean PCC = 0.689).

Multiple linear regression assumes linear relationships between model parameters and input features and allows for simple, normally-distributed errors. Surprisingly, this simple approach is capable of accurate cell-type specific predictions (mean PCC = 0.787; Table 1), likely due to the high raw correlation between the inputs and dependent variable. However this simple approach fails to cross-apply between cell types (mean PCC = 0.139; Table 1) indicating a problems with overfitting. This can be remedied through variable selection procedures, however a strength of the RF approach is that this step is not necessary, and pre-selection of model variables may result in a sub-optimal end result.^[15]

Partial least squares regression is a technique that uses dimensionality reduction to engineer a lower-dimension and orthogonal feature set. Hence this method is well-suited to collinear inputs, such as the set of variables used in this work (e.g. Fig. 14). As expected, PLS regression provides highly accurate cell type specific predictions

Table 1: Performance comparison of different modelling techniques. Comparison of mean Pearson correlation coefficient between predicted and observed compartment eigenvectors for three different modelling approaches: LM: linear regression; RF: Random Forest regression; PLS: partial least squares regression. Correlations were averaged per cell type over three cell types (cell type specific) and in the six possible crosses (cross-application) shown in Fig. 16.

	LM	RF	PLS
Cell type specific	0.787	0.790	0.750
Cross-application	0.139	0.689	0.641

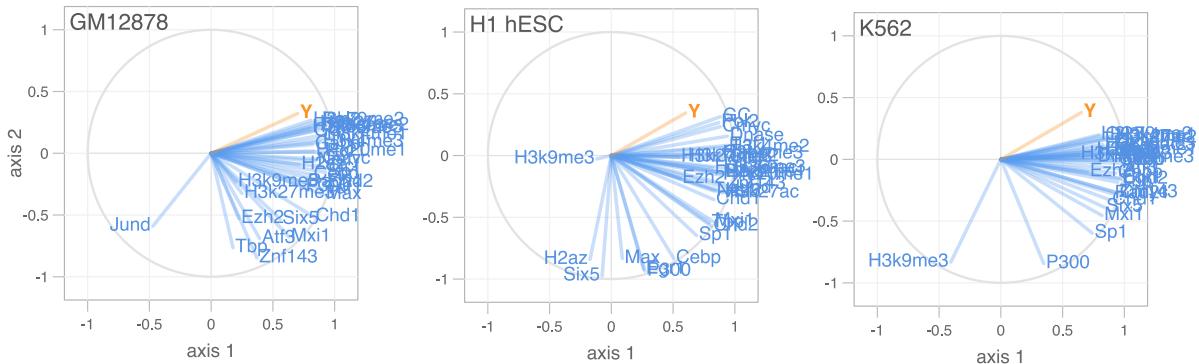


Figure 17: Circle of correlations of variables compared with PLS axes. Model variables are plotted against the first two axes used in PLS regression models per cell type. Y represents our compartment eigenvector.

(mean PCC = 0.750; Table 1) and performs well during cross-application (mean PCC = 0.641; Table 1), though in both cases produces slightly inferior results to RF models (Fig. 16).

PLS uses a type of dimensionality reduction, which offers another way to explore the inter-relationships between our feature set. Plotting input features against these lower-dimension components can give a revealing insight beyond simple correlations (e.g. Fig. 14). Figure 17 shows a “circle of correlations”, where features are plotted onto polar co-ordinates against the first two PLS components (Methods ??). Interpretation of this figure is that nearby variables in the scatterplot are positively correlated, and the vector length from the circle centre is proportional to said variable’s representation in the model. Negatively correlated variables point in opposite directions while uncorrelated variables are orthogonal to each other.^[28] We therefore see the known multicollinearity represented as groupings of overlapping variables in each cell type, with a smaller number of orthogonal and negatively correlated variables in each cell type (Fig. 17).

1.5.3 Non-independence

As recognised through our use of Hidden Markov Models (Methods ??), consecutive bins along a chromosome are non-independent yet thus far predictive models have not considered this inter-dependence.

This is for two reasons: firstly non independence could be thought of as an artefact of bin-sizing (we have elected to use regular, fixed binning beneath the scale of compartments themselves whereas another approach could use variable bin sizes, for example per compartment, TAD or restriction fragment); secondly using information of a bin's surroundings may obscure by proxy the chromatin features which would otherwise prove predictive. As an example, knowing that bin x_{i-1} and bin x_{i+1} are in compartment state A would allow us with high confidence to predict the state of bin x_i , but without learning anything of the region's relationship with its encompassed histone modifications and bound factors.

1.6 PARSIMONIOUS MODELS FROM EXPANDED FEATURE SETS

Strongly predictive models can be useful tools to reason about a complex system, however from a researcher's perspective there also exists a trade-off between predictive power and parsimony. Namely simpler models with fewer inputs may be more interpretable and of wider utility, for example they could be applied to cell types with less ChIP-seq data available than those used in this work. For this reason we explore parsimonious models with reduced feature sets, with an aim to build simpler models of chromatin state while retaining, if possible, similar levels of predictive accuracy.

On the other hand, the 35 variables used thus far as model inputs are not the complete set available in each cell type, but only the subset of those assayed in all three cell types under study. The ENCODE consortium has produced a significantly greater number of datasets^[1,29] in each cell type which have thus far gone unused. Here we explore models of higher order chromatin structure, in some cases built from over 100 variables, and then generate parsimonious models using optimal subsets guided by statistical techniques that penalise model complexity.

1.6.1 Stepwise regression

Multiple linear regression is a simple and analytically well-described modelling framework which is amenable to regularisation through a variety of methods. A simple approach is to start with a complete model and serially remove and/or add variables, then calculate a metric (here we use the Bayesian information criterion, BIC) which weighs the the model likelihood against model complexity. This process is

Table 2: Performance comparison of full and optimised RF and ML models. PCC between predicted and empirical compartment eigenvectors is shown for a range of modelling scenarios, including multiple linear regression (LM) and Random Forest (RF) approaches. For model selection, two methods are used: stepwise BIC-regularised linear models and LASSO regression; in each case those same features were then also used in building a separate RF for comparison.

	GM12878			H1 hESC			K562		
	n	LM	RF	n	LM	RF	n	LM	RF
All features	115	.836	.828	71	.744	.755	187	.811	.813
Matched subset	35	.827	.823	35	.740	.747	35	.796	.799
LASSO ℓ_1	23	.823	.836	23	.734	.750	39	.779	.811
Stepwise BIC	21	.840	.831	13	.746	.738	27	.819	.810

iterated until the metric reaches a (local) minimum, thus creating a more parsimonious model which retains predictive accuracy and should be less prone to overfitting. Stepwise regression also aids interpretation by selecting representative features from collinear clusters.^[30] A detailed explanation of this feature selection procedure can be found in Methods ???. It should be noted that despite its continued widespread usage, several statistical issues have been identified with the stepwise procedure for model selection.^[31,32]

In terms of model performance alone, stepwise regression gives the highest predictive accuracy on a held-out validation set in each cell type specific model of compartment eigenvector (Table 2), however it must be said that differences in model performance across all comparisons are modest. These results do show that even expanded feature sets of up to 187 input features add little explanatory power beyond that of much less complex models with 20 or fewer input variables (Table 2).

1.6.2 LASSO regression

A more modern technique for regularisation of linear models is the least absolute shrinkage and selection operator (LASSO). In brief, the LASSO is a form of ℓ_1 regularisation that penalises the sum of absolute values of standardised regression coefficients. By penalising total absolute values, rather than squared values as in ℓ_2 regularisation, coefficients can be shrunk to 0 thereby removing terms from the model.^[33,34] Thus LASSO combines the coefficient shrinkage of techniques like Ridge regression with a type of feature selection as seen in stepwise regression. A more rigorous description of this method is given in Methods ??.

Again we can perform a simplistic comparison of model performance using LASSO regression and other techniques (Table 2). LASSO retrieves comparable numbers of informative variables to the stepwise regression technique in each cell type, and again removes the majority of input features from expanded sets as redundant or relatively uninformative.

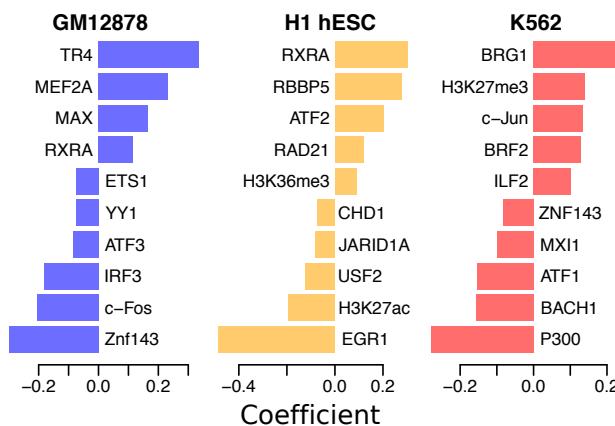


Figure 18: Ten largest LASSO model coefficients from expanded feature sets. Those coefficients with the largest absolute value are plotted for each cell type specific LASSO model.

Of those variables with a non-zero coefficient at the optimally-selected tuning parameter λ (Methods ??), the ten largest in each cell type are shown (Fig. 18). Similarities can be observed with variable importance from previous 35 input models (Fig. 11), including the large (negative) coefficient for EGR1 in the H1 hESC model as well as that of P300 in K562 (Fig. 18).

Of interest is the appearance of group of factors known to collectively form the heterodimeric activator protein-1 (AP-1), these include c-Fos, c-Jun and ATF1–3; all are spread across the most highly-ranked variables in each model of chromatin organisation (Fig. 18). The AP-1 complex has been shown to have DNA bending properties,^[35] and recently FOS and JUN members were associated with long-range chromatin interactions^[36] suggesting an under-explored role for this complex in genome organisation.

1.6.3 Regularised Random Forest

Random Forest (RF) comparisons are included for comparison in Table 2 where RF models were built using those features selected by procedures based on linear regression. Thus the linear regression-based feature selection acts as a “filter” method for feature selection, fully independent of the RF learning algorithm. A more coherent approach might be an “embedded” method, where a regularisation procedure is integrated with the learning algorithm.^[37,38]

While RF is a much younger technique than linear models, a framework for Regularised Random Forests (RRF) has recently been described^[39] and implemented in the R package RRF.^[40] The RRF algorithm uses the idea that at each node in a tree, unused variable should only be included if they offer a significant information gain over those available variables which have already been used in the tree. This differs

from the standard RF algorithm where splitting decisions at each node are entirely independent of each other (Methods ??).

We found that this algorithm was unable to perform feature selection on our highly collinear feature set, instead leaving full or almost full feature sets in each case (*data not shown*) and so providing equal results to a standard RF model using expanded feature sets (Table 2).

REFERENCES

- [1] ENCODE (2012) An integrated encyclopedia of DNA elements in the human genome. *Nature*, **489**(7414): 57–74.
- [2] Bernstein BE, Stamatoyannopoulos Ja, Costello JF, Ren B, Milosavljevic A, Meissner A, Kellis M, Marra Ma, Beaudet AL, et al. (2010) The NIH Roadmap Epigenomics Mapping Consortium. *Nature biotechnology*, **28**(10): 1045–8.
- [3] Tippmann SC, Ivanek R, Gaidatzis D, Schöler A, Hoerner L, van Nimwegen E, Stadler PF, Stadler MB, Schübeler D (2012) Chromatin measurements reveal contributions of synthesis and decay to steady-state mRNA levels. *Molecular systems biology*, **8**(593): 593.
- [4] McLeay RC, Lesluyes T, Cuellar Partida G, Bailey TL (2012) Genome-wide in silico prediction of gene expression. *Bioinformatics*, **28**(21): 2789–96.
- [5] Dong X, Greven MC, Kundaje A, Djebali S, Brown JB, Cheng C, Gingeras TR, Gerstein M, Guigó R, et al. (2012) Modeling gene expression using chromatin features in various cellular contexts. *Genome biology*, **13**(9): R53.
- [6] Shiraki T, Kondo S, Katayama S, Waki K, Kasukawa T, Kawaji H, Kodzius R, Watahiki A, Nakamura M, et al. (2003) Cap analysis gene expression for high-throughput analysis of transcriptional starting point and identification of promoter usage. *Proceedings of the National Academy of Sciences of the United States of America*, **100**(26): 15776–81.
- [7] Kodzius R, Kojima M, Nishiyori H, Nakamura M, Fukuda S, Tagami M, Sasaki D, Imamura K, Kai C, et al. (2006) CAGE: cap analysis of gene expression. *Nature methods*, **3**(3): 211–22.
- [8] Kolasinska-Zwierz P, Down T, Latorre I, Liu T, Liu XS, Ahringer J (2009) Differential chromatin marking of introns and expressed exons by H3K36me3. *Nature genetics*, **41**(3): 376–81.
- [9] Schaft D (2003) The histone 3 lysine 36 methyltransferase, SET2, is involved in transcriptional elongation. *Nucleic Acids Research*, **31**(10): 2475–2482.
- [10] Zuber V, Strimmer K (2011) High-Dimensional Regression and Variable Selection Using CAR Scores. *Statistical Applications in Genetics and Molecular Biology*, **10**(1): 1–27.
- [11] Ryba T, Hiratani I, Lu J, Itoh M, Kulik M, Zhang J, Schulz TC, Robins AJ, Dalton S, Gilbert DM (2010) Evolutionarily conserved replication timing profiles predict long-range chromatin interactions and distinguish closely related cell types. *Genome research*, **20**(6): 761–70.
- [12] Chambers EV, Bickmore Wa, Semple CA (2013) Divergence of mammalian higher order chromatin structure is associated with developmental loci. *PLoS computational biology*, **9**(4): e1003017.

- [13] Consortium TF, Pmi R, Dgt C (2014) A promoter-level mammalian expression atlas. *Nature*, **507**(7493): 462–70.
- [14] Breiman L (2001) Random forests. *Machine learning*, **45**: 5–32.
- [15] Díaz-Uriarte R, Alvarez de Andrés S (2006) Gene selection and classification of microarray data using random forest. *BMC bioinformatics*, **7**: 3.
- [16] Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ, et al. (2009) Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*, **326**(5950): 289–93.
- [17] Imakaev M, Fudenberg G, McCord RP, Naumova N, Goloborodko A, Lajoie BR, Dekker J, Mirny LA (2012) Iterative correction of Hi-C data reveals hallmarks of chromosome organization. *Nature methods*, **9**(10): 999–1003.
- [18] Babyak Ma (2004) What You See May Not Be What You Get: A Brief, Nontechnical Introduction to Overfitting in Regression-Type Models. *Psychosomatic Medicine*, **66**(3): 411–421.
- [19] Hawkins D (2004) The Problem of Overfitting. *Journal of Chemical Information and Modeling*, **44**(1): 1–12.
- [20] Beringer M, Ballar C, Croce LD, Viz P (2015) Role of PRC2-associated factors in stem cells and disease. **282**: 1723–1735.
- [21] Creyghton MP, Markoulaki S, Levine SS, Hanna J, Lodato Ma, Sha K, Young Ra, Jaenisch R, Boyer La (2008) H2AZ Is Enriched at Polycomb Complex Target Genes in ES Cells and Is Necessary for Lineage Commitment. *Cell*, **135**(4): 649–661.
- [22] Deb G, Singh AK, Gupta S (2014) EZH2: Not EZHY (Easy) to Deal. *Molecular cancer research : MCR*, **12**(5): 639–53.
- [23] Nie Z, Hu G, Wei G, Cui K, Yamane A, Resch W, Wang R, Green DR, Tessarollo L, et al. (2012) c-Myc is a universal amplifier of expressed genes in lymphocytes and embryonic stem cells. *Cell*, **151**(1): 68–79.
- [24] Kieffer-Kwon KR, Tang Z, Mathe E, Qian J, Sung MH, Li G, Resch W, Baek S, Pruett N, et al. (2013) Interactome maps of mouse gene regulatory domains reveal basic principles of transcriptional regulation. *Cell*, **155**(7): 1507–20.
- [25] Zervos AS, Gyuris J, Brent R (1993) Mxi1, a protein that specifically interacts with Max to bind Myc-Max recognition sites. *Cell*, **72**(2): 223–232.
- [26] Sun XJ, Man N, Tan Y, Nimer SD, Wang L (2015) The Role of Histone Acetyltransferases in Normal and Malignant Hematopoiesis. *Frontiers in Oncology*, **5**(May): 1–11.
- [27] Suzuki R, Shimodaira H (2006) Pvclust: an R package for assessing the uncertainty in hierarchical clustering. *Bioinformatics (Oxford, England)*, **22**(12): 1540–2.
- [28] Abdi H, Williams LJ (2010) Principal component analysis. *Wiley Interdisciplinary Reviews: Computational Statistics*, **2**(4): 433–459.

- [29] Boyle AP, Araya CL, Brdlik C, Cayting P, Cheng C, Cheng Y, Gardner K, Hillier LW, Janette J, et al. (2014) Comparative analysis of regulatory information and circuits across distant species. *Nature*, **512**(7515): 453–456.
- [30] Mantel N (1970) Why Stepdown Procedures in Variable Selection. *Technometrics*, **12**(3): 621–625.
- [31] Hurvich CM, Tsai CI (1990) The Impact of Model Selection on Inference in Linear Regression. *The American Statistician*, **44**(3): 214.
- [32] Whittingham MJ, Stephens PA, Bradbury RB, Freckleton RP (2006) Why do we still use stepwise modelling in ecology and behaviour? *Journal of Animal Ecology*, **75**(5): 1182–1189.
- [33] Tibshirani R (1994) Regression Selection and Shrinkage via the Lasso.
- [34] Hastie T, Tibshirani R, Friedman J (2009) *Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2 edition. ISBN 978-0-387-84858-7.
- [35] Kerppola TK, Curran T (1993) Selective DNA bending by a variety of bZIP proteins. *Molecular and cellular biology*, **13**(9): 5479–5489.
- [36] Heidari N, Phanstiel DH, He C, Grubert F, Jahanbanian F, Kasowski M, Zhang MQ, Snyder MP (2014) Genome-wide map of regulatory interactions in the human genome. *Genome Research*.
- [37] Guyon I, Weston J, Barnhill S, Vapnik V (2003) An Introduction to Variable and Feature Selection. *Journal of Machine Learning Research*, **3**: 1157–1182.
- [38] Kohavi R, Kohavi R (1997) Wrappers for feature subset selection. *Artificial Intelligence*, **97**(1-2): 273–324.
- [39] Deng H, Runger G (2012) Feature selection via regularized trees. In: *Neural Networks (IJCNN), The 2012 International Joint Conference on*. IEEE, pp. 1–8.
- [40] Deng H, Runger G (2013) Gene selection with guided regularized random forest. *Pattern Recognition*, **46**(12): 3483–3489.