

1 Motif Interactions Affect Post-Hoc Interpretability of 2 Genomic Convolutional Neural Networks

3 **Marta S. Lemanczyk^{1,*}, Jakub M. Bartoszewicz^{1,2}, and Bernhard Y. Renard^{1,*}**

4 ¹Hasso Plattner Institute for Digital Engineering, Digital Engineering Faculty, University of Potsdam

5 ²Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology

6 *corresponding author(s): Bernhard Y. Renard (Bernhard.Renard@hpi.de)

7 **ABSTRACT**

Post-hoc interpretability methods are commonly used to understand decisions of genomic deep learning models and reveal new biological insights. However, interactions between sequence regions (e.g. regulatory elements) impact the learning process as well as interpretability methods that are sensitive to dependencies between features. Since deep learning models learn correlations between the data and output that do not necessarily represent a causal relationship, it is difficult to say how well interacting motif sets are fully captured. Here, we investigate
8 how genomic motif interactions influence model learning and interpretability methods by formalizing possible scenarios where interaction effects appear. This includes the choice of negative data and non-additive effects on the outcome. We generate synthetic data containing interactions for those scenarios and evaluate how they affect the performance of motif detection. We show that post-hoc interpretability methods can miss motifs if interactions are present depending on how negative data is defined. Furthermore, we observe differences in interpretability between additive and non-additive effects as well as between post-hoc interpretability methods.

9 **Background**

10 Convolutional neural networks (CNN) excel at various sequence-based tasks due to their capability to learn patterns
11 and complex interactions making these models an efficient method for many predictive tasks in the field of genomics
12 [1]. However, for many biological applications, predictions alone are insufficient for understanding the underlying
13 mechanisms for a given problem [2]. Besides verifying that a model learned meaningful predictions, interpreting
14 CNNs can lead to new insights for genomic questions [3, 4]. With the help of such interpretations, the models'
15 decisions can be verified, or new insights can be obtained [5]. To explain the outcome of a CNN, post-hoc
16 interpretation methods are a commonly used approach. Instead of training intrinsically interpretable models, post-
17 hoc interpretation methods are applied after the training process on a fully trained model. Known methods include
18 feature permutation [6], Integrated Gradients [7], and DeepLIFT [8]. When applied to sequence data, scores are
19 assigned to each position in an input sequence based on their contribution to an individual prediction. By aggregating

20 attribution scores from multiple input sequences, it is possible to extract meaningful motifs [9]. Identified motifs can
21 be compared with known motifs in task-specific databases to determine their biological relevance by using methods
22 for motif comparison like TOMTOM [10, 11].

23 The quality of contribution scores can be affected by multiple factors besides model complexity so the evaluation of
24 interpretability performance for machine learning models is crucial [12, 13, 14]. This also applies to interpretability
25 for biological neural networks. Some differences in contribution scores can be attributed to the architectural choices
26 for the model. In [15] it was shown that exponential activation functions can lead to more interpretable motif
27 representations in first-layer filters than for other functions, like ReLU, sigmoid or tangent activation function, as
28 well as specific choices for filter size, max pooling width and model depth [16]. Interpretability can be also improved
29 by introducing robustness with the help of regularization, random noise injection, and adversarial training [17].

30 While those approaches improve interpretability in general, one has to keep in mind that post-hoc interpretations
31 represent what a model learned. The learned correlative features do not necessarily represent causal effects.
32 Dependencies between features complicate learning causal effects with machine learning models. Not only are many
33 features in biological sequences in relationship with others, but groups of such locally dependent features can also be
34 part of a higher interaction representing a regulatory logic hidden in deeper layers [18]. Such biological interaction
35 can be, for instance, cooperative binding of transcription factors to DNA [19]. This can result in misleading
36 explanations if the complete underlying biological mechanism is not uncovered by the model. Furthermore, some of
37 the attribution methods are based on the assumption that input features are independent. Dependencies between
38 motifs can influence attribution methods so that subsets of interacting motifs can produce incomplete or noisy
39 interpretations. It is crucial to analyze if post-hoc methods are capable of capturing interacting motifs and, therefore,
40 the underlying causal effects in an understandable manner.

41 To tackle those challenges, we design suitable data for different sources of interactions and evaluate the interpretability
42 performance of models trained with that data. For that, we first investigate how interactions affect model training
43 and interpretability in general by using different negative data sets, forcing the model to learn interactions explicitly
44 or not (Fig. 1, left). The selection of negative data for machine learning is an ongoing problem for various biological
45 tasks [20, 21] since it can influence the discriminative power of models. One mistake is not to include data where
46 there is some uncertainty for the data label, for example, due to the similarity between positive and negative instances.
47 Using easily distinguishable data results in a simpler training task and can still lead to good predictions if the new
48 data is similarly structured to training data. However, the predictive performance decreases when the model is
49 confronted with uncertain data like novel populations. In the case of genomic sequences, this can happen when
50 using random sequences for the negative training data set instead of carefully curated samples. Models that are
51 trained on non-random negative data have to learn more complex relationships between motifs in the positive class
52 in contrast to random negative sequences, which allow the model only to learn subsets or even individual motifs to

53 distinguish between classes. While the influence of negative data on prediction tasks is a known problem, it is not
54 well explored how it influences interpretability.
55 Secondly, we look at the interaction effects between motifs (Fig. 1, right). In regulatory genomics, different
56 experiments allow us to determine the functions and characteristics of non-coding DNA regions [22]. This includes,
57 among other methods ChIP-seq, DNase-seq, and ATAC-seq where the enrichment of sequence fragments is measured
58 for various functions like protein binding locations or chromatin accessibility. The output is then mapped to the
59 genome resulting in per-position counts, which can then be binarized based on signal peaks representing the sites of
60 interest. Previous deep learning models focused on the prediction of binarized peaks [23] and perform, therefore, a
61 classification problem. Currently, many state-of-the-art models are used to directly predict the signal of an assay
62 rather than just the presence of a peak so that the output gives a more precise prediction of the signal of interest [24].
63 However, the regression task is more complex since the model needs to learn a function between the input sequence
64 and the numerical output instead of just distinguishing between classes. Since multiple regulatory elements can be
65 involved in a regulatory mechanism, interactions between motifs complicate the prediction task. Motif interactions
66 can occur in multiple forms, including additive effects as well as multiplicative interactions [25]. Here, we explore if
the complexity of interaction effects influences model interpretability.

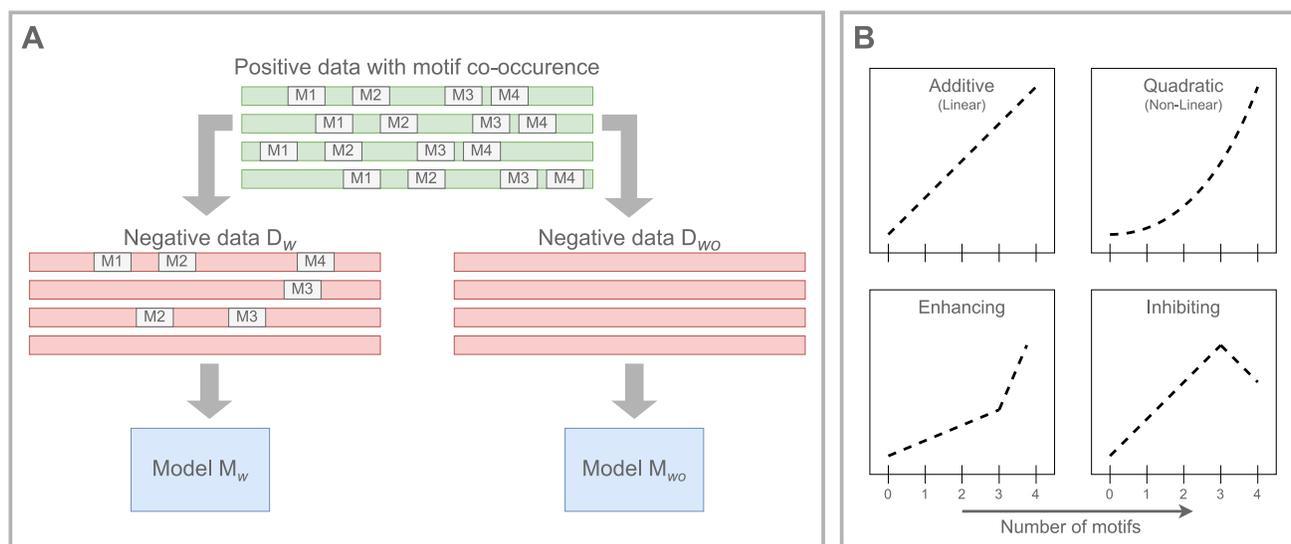


Figure 1. Data design for motif interactions. **(Classification)** Capturing a motif interaction in the form of co-occurrence does not only depend on the positive dataset containing all motifs. Depending on the nature of the negative data, models differ in how the class boundaries and the interactions are learned. Here, we also analyze the impact on interpretability. **(Regression)** Deep learning models are capable of learning interactions with varying complexities. We investigate if complexity also influences post-hoc interpretability. Detailed experiment and data generation settings can be found in Fig. 4.

68 Results

69 We evaluated the influence of motif interactions on motif detection by first defining various [interactions](#) and then
70 [simulating data](#) containing those interactions. Model architecture and training setup are described in the [CNN](#)
71 section. We evaluated the model performance with regard to prediction and motif detection capability as well as the
72 post-hoc attribution methods with metrics introduced in the [evaluation](#) section.

73 The motif set used for the following evaluation consists of motifs in [table 1](#) obtained from the JASPAR database for
74 transcription factor binding sites [[26](#)].

75 **Negative Sequences** To evaluate the influence of negative data on model interpretability, we simulate the scenario
76 in a classification problem predicting an outcome based on the co-occurrence of a set of motifs. For that, we
77 create two data sets that contain the same sequences for the positive class but differ in their negative sequences.
78 Specifications can be found in the respective [method section](#). The positive data set is described by the co-occurrence
79 of all n motifs that are here set to $n = 4$. Regarding the negative data, we distinguish between the data set containing
80 0 to $n - 1$ motifs and a data set containing only random sequences without any motifs inserted. We call the model
81 trained on data with motif subsets in the negative data set \mathcal{M}_w and the model without motifs \mathcal{M}_{wo} .

82 Based on the negative data set used for training, models learn different ways to predict the positive class. To
83 investigate the underlying learning mechanisms, we use negative test data similar to the training data set with motif
84 subsets in the negative data set. Each possible combination of motifs is represented by the same number of sequences
85 in the negative data set, as can be seen in [Fig. 2A](#). The respective accuracy can be seen in [Fig. 2B](#). While both
86 models have a decent accuracy for the positive data ($acc(\mathcal{M}_w) = 0.9485$, $acc(\mathcal{M}_{wo}) = 0.9992$), the models differ in
87 the accuracy for the negative data set ($acc(\mathcal{M}_w) = 0.9112$, $acc(\mathcal{M}_{wo}) = 0.5916$, more detailed accuracies can be
88 found in [Table 2](#)).

89 We trained models with a minimal number of filters so that the number of filters equals the number of motifs (#filters
90 = 4), as well as CNNs with 32 filters. Both minimal filter models captured nearly identical motif weights in the
91 convolutional filter ([Fig. 2C](#)), showing that both models learned a similar representation of the inserted motifs.

92 We calculated the contribution scores for test sequences using DeepLIFT (DL), Integrated Gradients (IG), and
93 Feature Permutation (FP). Since the ground truth is known for our simulated data, we can quantify how well the
94 contribution scores were assigned to the motifs. The AUPRC value (see methods) should be high if an attribution
95 method assigned high absolute scores to the motif positions compared to random positions. By using the absolute
96 scores, the negative influence of the motifs is also captured.

97 In [Fig. 2D](#), AUPRC scores are shown for the dataset *distinct_1* for the NHLH1 Motif (ID: MA0048.1). The test
98 sequences contained subsets of motifs, of which at least one was the motif of interest. For the positive data with all
99 motifs present ([2D, iv](#)), differences can be observed between the AUPRC medians of both models with \mathcal{M}_w having

100 lower AUPRC ($\Delta IG : 0.173931, \Delta DL : 0.194141, \Delta FP : 0.133069$, see Table 5). The motif detection performance
101 for M_w drops for the negative sequences containing 1 to 3 motifs, while the performance for M_{wo} remains more
102 stable for all sequences (2D, i-iii). There are no major differences between the attribution methods for motif NHLH1
103 (see 5). Similar observations can be made for the other motifs (see Supplement Fig. 5). Besides the data sets
104 containing heterologous motifs, data sets with homologous motifs were investigated. Similar to the models based on
105 heterologous motif sets, motif detection performance decreases for M_w the fewer motifs are present in the sequence.
106 However, for all motifs (MEF2A), an increase in the AUPRC scores can be observed for Feature Permutation when 1
107 or 2 motifs are present while the performance for M_{wo} drops.

108 **Interaction effects** As described in the method section, various interaction effects based on the co-occurrence
109 of an interacting motif set are represented by a function that generates labels for the input sequences. For that,
110 we assigned each sequence a numerical label based on the number of motifs contained in the input sequence. We
111 distinguish between an additive, enhancing, quadratic, and inhibiting effect. Further details can be seen in Fig.1
112 (right). To make the prediction accuracy comparable, we bin the outcomes around the possible labels and calculate if
113 the predicted value falls in the interval. We averaged the accuracy and the recall as well as the AUPRC across 5
114 models with different seeds to decrease potential noise. We compare the AUPRC values for all effects individually
115 for each motif. We also differentiate between a model with 4 and 32 filters.

116 We calculate the recall separately for each subset depending on the number of containing motifs (see Table 6).
117 The total accuracy for each effect lies between 0.8599 and 0.8729 for the model with 32 filters. However, there
118 are differences between the effects when it comes to recall. The additive interaction effect has a mostly stable
119 performance between subsets, while the other effects vary in recall. Additionally, the performance for negative
120 sequences without any motifs is low for all effects ranging between 0.3588 to 0.6568.

121 AUPRC scores are compared between the interaction effects as well as interpretability methods and motifs. For
122 IG, The motif-wise AUPRC for the additive effect is high for all motifs with median values between 0.86 and
123 1, similarly for the inhibiting effect (0.85-0.98). For the enhancing and quadratic effect, the AUPRC values do
124 not show a large decrease for motifs NHLH1 and ETS1, whereas for motifs CREB1 and TEAD1 the values drop,
125 especially for the quadratic interaction (AUPRC median: 0.75 for CREB1 and 0.68 for TEAD1) including an
126 increased variance in AUPRC scores. Feature Permutation performs in total worse than Integrated Gradients and
127 DeepLIFT having a large variance with median values around 0.5 to 0.85. Here, the highest values were assigned
128 to the inhibiting effect while Feature Permutation also performed worse for other non-additive effects. DeepLIFT
129 shows the most robust AUPRC along the motifs and effects. However, the assigned values to motifs CREB1 and
130 TEAD1 are slightly worse for the non-additive effects, including the inhibiting effect.

131

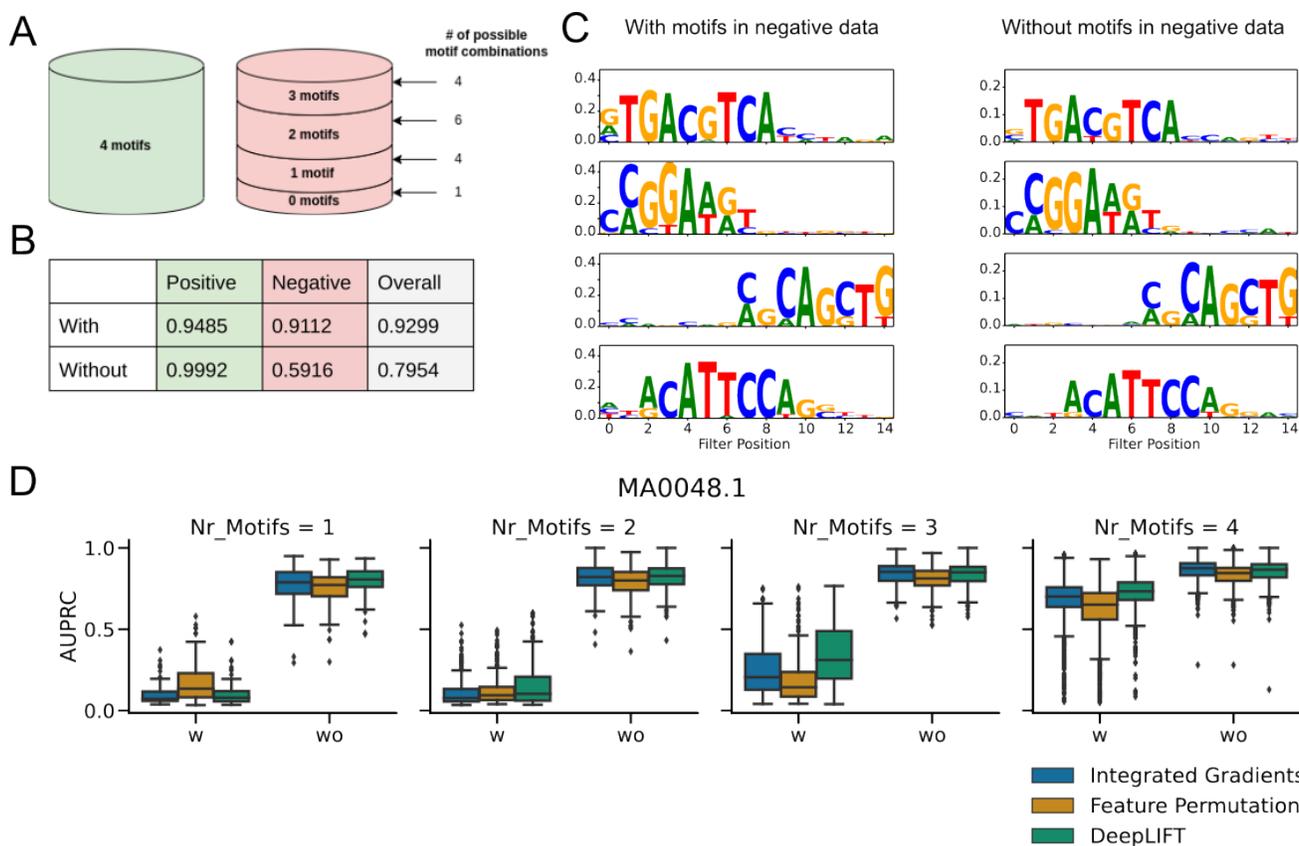


Figure 2. Results for models trained on two different negative sequence data sets to predict the co-occurrence of 4 motifs. **A** Both model types are trained on the same positive data containing sequences with 4 motifs inserted. One model type (M_w) additionally includes motif subsets with max. 3 motifs in the negative data set as depicted on the right side, while the second only includes random sequences without any motifs (M_{wo}). **B** Predictive accuracy for both models. While both models have good accuracy for the positive class containing 4 motifs, the model M_{wo} performs poorly on the negative data, as expected. **C** Weights of convolutional layer filters. Both models learned similar representations of the motifs within the layer. **D** AUPRC values for contribution scores for positive class. High values indicate high contribution scores for motif positions compared to random sequences and, therefore, better motif detection. For M_w , the AUPRC scores decrease significantly for models containing only subsets of motifs which indicates a lower motif detection capability.

132 **Subsets of interactive motifs in negative data can lead to different decision boundaries and**
 133 **therefore to varying motif detection performances**

134 While the performance for positive sequences with all motifs from the interactive set is similar for both models, there
 135 are strong differences for the negative sequences when it comes to accuracy and motif detection performance. As
 136 expected, M_{wo} is not capable of classifying all negative sequences containing subsets of motifs correctly since there
 137 were no sequences with subsets in the training data. Accuracy drops with increasing size of the motif subset. The

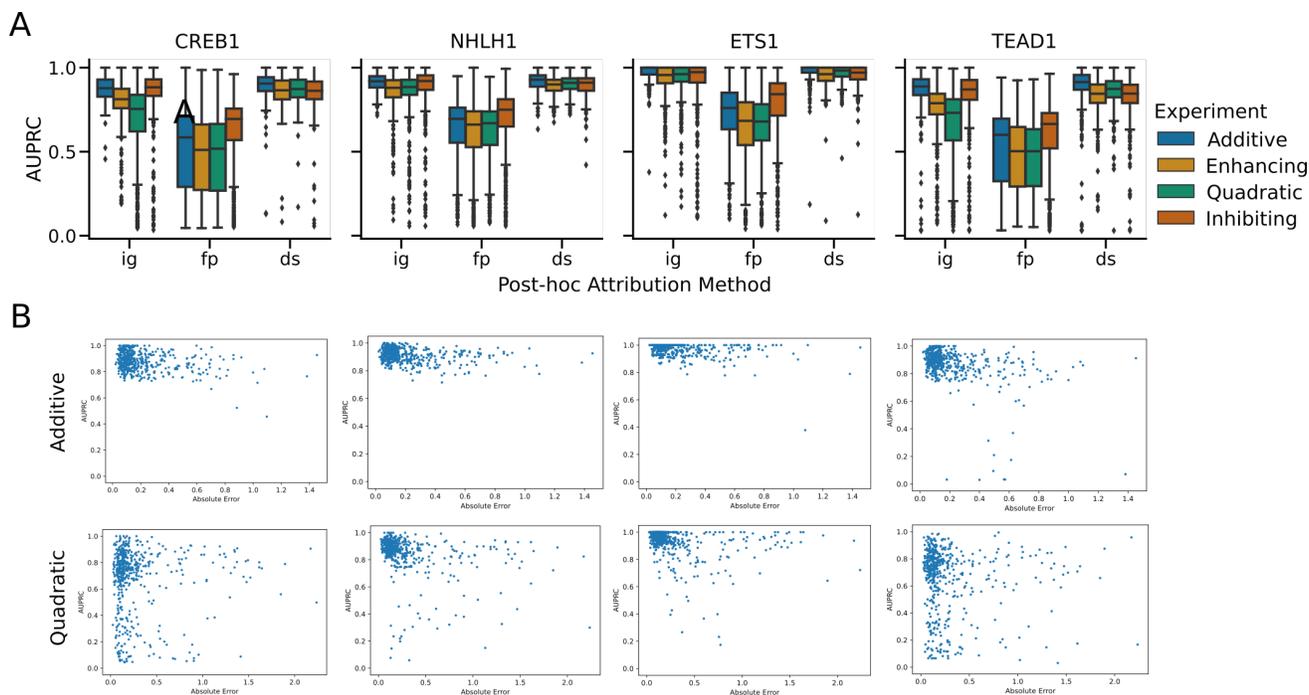


Figure 3. Results for interaction effects for sequences containing all 4 motifs.

A Motif-wise AUPRC values for heterologous motif set. For Integrated Gradients a drop in AUPRC values can be observed for the enhancing and quadratic effect for motifs for CREB1 and TEAD1 compared to motifs NHLH1 and ETS1. Only a small decrease can be observed for the non-additive interactions for DeepLIFT. Feature Permutation performs overall worse for the regression task having high variance in the scores. **B** Comparison between Accuracy and AUPRC values for the additive and quadratic effect for Integrated Gradients. Low AUPRC values for the quadratic effect do not show increased error values.

138 accuracy for sequences containing only 1 motif is still high ($acc_{1 \text{ motif}} = 0.9795$), while for sequences with 2 and
 139 3 motifs, the model performs poorly ($acc_{2 \text{ motifs}} = 0.6$, $acc_{3 \text{ motifs}} = 0.089$) compared to M_w ($acc_{2 \text{ motifs}} = 0.9647$,
 140 $acc_{3 \text{ motifs}} = 0.7215$). This indicates that the model M_{w_0} does not generalize well and learns only subsets of the
 141 interactive motif sets to classify a sequence as positive instead of the full motif set like M_w . The accuracy scores for
 142 sequences containing 2 motifs differ depending on the present motif subset. For instance, the averaged accuracy for
 143 the motif set containing motifs CREB1 and TEAD1 reaches 0.874 while for the set containing NHLH1 and ETS1
 144 only 0.144 (see all accuracies in 3). Since the negative sequences containing 3 motifs are mostly classified as
 145 positive, and most sequences containing 1 motif are correctly classified as negative, those observations suggest that
 146 the decision boundary between classes is based on specific motif sets containing between 2 and 3 motifs.
 147 As we can see in the accuracy values, M_w has to learn a decision boundary between sequences containing 3 and 4
 148 motifs. If we permute one motif (which equals a motif removal) in the positive class sequence, the outcome should

149 change to negative. Here, the permuted motif has therefore a high contribution to the outcome. However, if we
150 look at a negative sequence containing 3 motifs and permute one motif, the resulting sequence with 2 motifs still
151 belongs to the negative class. The permuted motif does not contribute to an outcome change which can result in
152 low AUPRC scores. For M_{wo} , it is unclear how the decision boundary is learned and which individual motifs or
153 motif subsets must be present to influence the model's decision. Based on the observations in the accuracies, the
154 presence or absence of individual motifs could already impact the outcome in sequences with fewer motifs. That
155 impact could be reflected in the higher AUPRC scores also for the negative sequences.

156 **Non-additive interactions can influence interpretability independent of accuracy**

157 We observe a decrease in AUPRC values for non-additive interaction effects for the models with 32 convolutional
158 filters (see Fig. 3 A). DeepLIFT performed the most stable when it comes to interpretability, with only a small
159 decrease between the additive effect and the non-additive ones. On the other hand, IG performed worse on motifs
160 CREB1 and TEAD1 for the enhancing and the quadratic effect, while there is no large difference for motifs NHLH1
161 and ETS1 between the different interaction effects. Feature Permutation has overall low AUPRC values suggesting
162 that it is not suitable for the regression task. We also included the absolute errors of the predictions to validate if
163 lower AUPRC values result from bad predictions. The absolute errors of the predictions for the sequences with
164 low AUPRC values for the quadratic effect do not show an increase (see Fig. 3 B). Therefore, we can assume that
165 the worse performance in detecting the motifs can result from the more complex interaction independent of the
166 accuracy.

167 **Discussion**

168 Interpretation has become a crucial part of deep learning applied in the field of genomics. While the validation of
169 identified motifs can confirm that a model learned meaningful patterns, the lack of complete biological data makes it
170 difficult to prove the completeness of all motifs and therefore derive causal insights. Here we investigated from a
171 data-centric point of view how interactions in genomic datasets can result in missing or noisy interpretations.

172 We concentrate in this work on the effects of the co-occurrence of motifs and their effects on the outcome. However,
173 there are further aspects that can influence motif interactions. For example, the cooperativity of transcription factor
174 motifs can additionally depend on order, orientation, and distance (including 3-D genome distances) between
175 regulatory elements [27] as well as temporal causes [28]. Motifs in our generated data sets have a fixed order and
176 orientation as well as the same distance between each other for all sequences in one data set. In this way, we focus
177 solely on a fixed grammar to reduce complexity and other factors that could affect interpretability. It is also important
178 to point out that the model might see the tasks more as a classification counting the number of motifs instead of
179 learning a function since the labels are discrete based on the number of motifs in the sequence. In this case, the kind

180 of interaction function might not be relevant to the model since it is not learning the function itself. Including other
181 information like distances between motifs could improve the simulation by making a long function continuous and,
182 therefore, additive and non-additive interaction effects could be better explored.

183 In the negative sequence experiment, we observe a trade-off between accuracy and motif detection, especially
184 for negative data. Here, multiple data augmentation strategies can be evaluated to obtain better motif detection
185 performance or desired interpretability outcomes on interactive data while preserving biological functionality (e.g.
186 [29]). Motif detection can also be improved by accounting for interactions like in [18] where stochastic masks are
187 used to find sets of motif features that preserve or change the outcome and therefore avoid saturation effects.

188 In this work, we focused on simple CNNs to break down the problem of interpretability to interactions. Using
189 more complex architectures would result in additional sources affecting interpretability so it would be more difficult
190 to separate the interaction effects from the other sources. However, currently, more complex models with more
191 sophisticated modules such as the attention mechanism are applied to genomic problems to capture interactions
192 within genomic sequence data (see an overview on genomic large language models (LLM) [30]). So far, many
193 approaches to interpreting genomic LLM models focus on the analysis of the attention scores or the output with
194 post-hoc methods that mostly offer interpretations on the input token level. One ongoing challenge is to uncover
195 the grammar between interacting motifs so that interpreting genomic LLMs beyond those approaches could give
196 better explanations of underlying biological processes. Also, pre-training of genomic LLMs should be explored
197 in the context of interactions. Especially, if downstream tasks are missing relevant data, like in the negative data
198 experiment, it is necessary to analyze how the missing information is imputed.

199 As we could also see in our results, machine learning models do not necessarily learn the underlying causal effect
200 of biological mechanisms. Thus interpreting models after training is not always suitable for knowledge extraction.
201 Therefore, designing interpretable architectures that capture the interactions explicitly instead of only relying on
202 post-hoc model interpretation could be a better approach for motif identification as well as interaction detection.

203 **Conclusion**

204 We analyze the influence of motif interactions on post-hoc interpretability methods. First, we investigate how
205 motif co-operativity can affect model learning depending on how interacting motifs are present in the negative
206 data set. We observe that interpretability performance can decrease when interactions are learned more explicitly
207 by the model. Especially for negative sequences, evidence for the positive class can be missed. Secondly, we
208 formalize different interaction effects (additive as non-additive) and compare those with regard to interpretability.
209 We discovered differences between the effects as well as the interpretability methods, from which we deduce that
210 post-hoc interpretability is affected by complex interactions.

211 **Methods**

212 **Motif interactions**

213 We define each prediction task as a function $F : X \rightarrow Y$. The input $x \in X = \{0, 1\}^n$ describes the presence or absence
214 for all motifs $i \in M = \{1, \dots, n\}$ in the input sequence. If a motif i is present in a given sequence, then $x_i = 1$, if
215 absent then $x_i = 0$. The outcome Y depends on the task. For regression problems, we define $Y = \mathbb{R}$, while $Y = \{0, 1\}$
216 applies for binary classification tasks.

217 **Interaction effect on outcome**

218 Motif interactions, eg. co-occurrence, can be expressed as logical constructs using AND, OR, and NOT. The nature
219 and magnitude of the effects of these relationships on the outcome (eg. non-linearity, inhibition, activation) can be
220 encoded in the target values of a regression task. Since we assume that there are no other features that can influence
221 the outcome except the given motifs, we derive the following definitions from the definitions in [31].

222 Let $F(x)$ be the sum of the effects of all possible subsets of motifs, where each motif combination has its influence
223 on the outcome:

$$F(x) = \sum_{S \subset M} f_S(S) \quad (1)$$

224 The independent effect of a motif m_i on $F(x)$, which does not rely on the presence or absence of other motifs, is
225 called the main effect and is defined here as a subfunction $f_i(x_i)$. If $F(x)$ is only affected by the main effects of
226 the motifs and therefore does not contain any interactions between them, the function is described as an additive
227 interaction:

$$F(x) = \sum_i f_i(x_i) \quad (2)$$

228 The task contains at least one non-additive interaction if there is a subset $S \subset M, |S| \neq 1$ so that $f_S(S) \neq 0$ and
229 therefore $F(x) \neq \sum_i f_i(x_i)$ [32]. In that way, we can define different interactions as functions. We show two examples
230 that we use for our evaluation.

231 *Example 1: Enhancement and Inhibition*

232 Besides the main effects of individual motifs, we introduce an enhancement/inhibition term so that

$$F(x) = \sum_i f_i(x_i) + c \prod_{i \in M} x_i \quad (3)$$

233 for some constant $c \in \mathbb{R} \setminus \{0\}$. Here, a non-zero value is added to the outcome if and only if the input sequence
234 contains all motifs from the interacting motif set. The co-occurrence of all motifs in that set enhances the individual
235 main effects on the outcome.

236

237 *Example 2: Non-linear relationship*

238 The relationship between motifs can also be expressed with a non-linear function depending on the subsets of the
239 interacting motif set. As an example of a nonlinear interaction, we use a quadratic relationship:

$$F(x) = \left(\sum_i f_i(x_i) \right)^2 \quad (4)$$

240 Combinations of different interactions are also possible and add complexity to the task. However, we use the
241 described interactions to investigate the differences between additive and non-additive interactions.

242 **Negative sequences**

243 The performance of a model strongly depends on the available data. The nature of the data set can have different
244 effects on how the model learns the interactions between motifs. Here, we simulate a binary classification problem
245 with positive and negative classes. While the data for the positive class plays a major role in binary classification,
246 the negative class can also impact the resulting model. In this case, the positive class represents the co-occurrence of
247 all motifs in the interacting motif set M so that $\forall i \in M : x_i = 1$. Different negative data sets can be chosen for the
248 same task while the model can still have a similar predictive performance in the end. We distinguish between two
249 different negative data sets. One data set includes individual motifs or subsets of the interacting motif set in the
250 negative data set so that $\exists i \in M : x_i = 0$. In contrast, sequences from the second negative data set do not contain any
251 motifs from the interacting motifs set and therefore $\forall i \in M : x_i = 0$.

252 **Sequence data**

253 Genomic sequences have to be transformed into numerical matrices so they can be processed by CNNs. Each column
254 of this matrix stands for one sequence position where the base at this position is represented by a one-hot-encoding
255 vector. We use sequences with the length of 250 base pairs resulting in matrices with the size of 4x250.

256 We obtain real transcription factor binding motifs from the JASPAR database [26] for the evaluation. We distinguish
257 here between subsets of homologous and heterologous motif subsets to investigate if motif similarity influences
258 interpretability. The similarity was measured by the Pearson correlation coefficient for motif similarity [10]. We
259 used the implementation from the biopython package [33]. Motifs can have different lengths, e.g. transcription
260 factor binding sites have a length of around 5-31 nt [34], which may also influence interpretability if cooperating
261 motifs differ in size. We picked for our experiment motifs with approximately similar lengths. The selected motifs
262 can be found in Table 1. For each task, a grammar is generated with a fixed order of motifs and distances between
263 the motifs. The grammar itself is inserted randomly within the sequence to ensure invariance regarding the position
264 of the motif grammars. The distance between the motifs is larger than the filter size, so the filters learn individual
265 motifs, not overlapping regions. Labels were generated by following the interaction definitions above. Training, test,

266 and validation sequences are the same for models that are compared, and the data sets only differ in the labels that
267 encode the interactions.

268 **Convolutional Neural Networks**

269 To ensure that differences in interpretability performance cannot be traced back to differences in predictive perfor-
270 mance, one requirement is that models that are compared have similar performance.

271 We use CNNs with one convolutional layer with filters approximating the length of the chosen motifs to learn localist
272 representations as described in [16]. 3 dense layers follow the convolutional layer to learn the interactions between
273 the motifs. We apply batch normalization on the inputs before passing them to a ReLU function. Additionally, we
274 apply max pooling in the convolutional layers.

275 **Interpretability Methods**

276 We use feature permutation (FP) [6], Integrated Gradients (IG) [7], and DeepLIFT (DL) [8] as post-hoc attribution
277 methods. The analyses are performed with the method implementations from the Captum library for PyTorch [35].
278 Since we obtained similar results for average and zero reference sequences, we use the zero reference sequence due
279 to the shorter computational time. Contribution scores are averaged over the 5 models to reduce noise.

280 **Evaluation**

281 Contribution scores were evaluated similarly to [15]. Each position in the sequence gets a label assigned depending
282 on if it belongs to a motif (1) or not (0). Contribution scores of motif positions are then compared to those of
283 random positions by calculating the AUPRC (see overview of the model evaluation in Supplement Fig. 4). AUPRC
284 scores are calculated for each motif separately. AUPRC scores were then visualized via boxplots. Interpretability
285 performance was then compared between the models of interest by analyzing the differences of the AUPRC on the
286 same sequence test set.

287 **Code availability**

288 Code for the evaluation is available under www.gitlab.com/dacs-hpi/interpret-interaction upon publication.

289 **Author contributions statement**

290 MSL designed the experiment. MSL, JMB, and BYR consulted on analytical decisions. MSL performed the analysis
291 and wrote the manuscript with input from all co-authors. The authors read and approved the final manuscript.

292 **Competing interests**

293 No competing interests.

294 References

- 295 1. Eraslan, G., Avsec, v., Gagneur, J. & Theis, F. J. Deep learning: new computational modelling techniques for
296 genomics. *Nat. Rev. Genet.* **20**, 389–403, [10.1038/s41576-019-0122-6](https://doi.org/10.1038/s41576-019-0122-6) (2019).
- 297 2. Talukder, A., Barham, C., Li, X. & Hu, H. Interpretation of deep learning in genomics and epigenomics.
298 *Briefings Bioinforma.* [10.1093/bib/bbaa177](https://doi.org/10.1093/bib/bbaa177) (2020).
- 299 3. Avsec, v. *et al.* Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat. Genet.* **53**,
300 354–366, [10.1038/s41588-021-00782-6](https://doi.org/10.1038/s41588-021-00782-6) (2021).
- 301 4. Bartoszewicz, J. M., Seidel, A. & Renard, B. Y. Interpretable detection of novel human viruses from genome
302 sequencing data. *NAR Genomics Bioinforma.* **3**, lqab004, [10.1093/nargab/lqab004](https://doi.org/10.1093/nargab/lqab004) (2021).
- 303 5. Novakovsky, G., Dexter, N., Libbrecht, M. W., Wasserman, W. W. & Mostafavi, S. Obtaining genetics
304 insights from deep learning via explainable artificial intelligence. *Nat. Rev. Genet.* **24**, 125–137, [10.1038/
305 s41576-022-00532-2](https://doi.org/10.1038/s41576-022-00532-2) (2023).
- 306 6. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32, [10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324) (2001).
- 307 7. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. In *Proceedings of the 34th*
308 *International Conference on Machine Learning*, vol. 70, 3319–3328 (2017).
- 309 8. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation
310 differences. *PMLR 70:3145-3153, 2017* (2017). [1704.02685](https://arxiv.org/abs/1704.02685).
- 311 9. Shrikumar, A. *et al.* Technical Note on Transcription Factor Motif Discovery from Importance Scores (TF-
312 MoDISco) version 0.5.6.5. *arXiv:1811.00416 [cs, q-bio, stat]* (2020). ArXiv: 1811.00416.
- 313 10. Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. S. Quantifying similarity between motifs.
314 *Genome Biol.* **8**, R24, [10.1186/gb-2007-8-2-r24](https://doi.org/10.1186/gb-2007-8-2-r24) (2007).
- 315 11. Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME Suite. *Nucleic Acids Res.* **43**, W39–W49,
316 [10.1093/nar/gkv416](https://doi.org/10.1093/nar/gkv416) (2015).
- 317 12. Adebayo, J. *et al.* Sanity Checks for Saliency Maps. In *Advances in Neural Information Processing Systems*,
318 vol. 31 (Curran Associates, Inc., 2018).
- 319 13. Sixt, L., Granz, M. & Landgraf, T. When explanations lie: Why many modified bp attributions fail. *Proc. 37th*
320 *Int. Conf. on Mach. Learn.* (2019). [1912.09818](https://arxiv.org/abs/1912.09818).
- 321 14. Zhou, J., Gandomi, A. H., Chen, F. & Holzinger, A. Evaluating the Quality of Machine Learning Explanations:
322 A Survey on Methods and Metrics. *Electronics* **10**, 593, [10.3390/electronics10050593](https://doi.org/10.3390/electronics10050593) (2021).
- 323 15. Koo, P. K. & Ploenzke, M. Improving representations of genomic sequence motifs in convolutional networks
324 with exponential activations. *Nat. Mach. Intell.* **3**, 258–266, [10.1038/s42256-020-00291-x](https://doi.org/10.1038/s42256-020-00291-x) (2021).
- 325 16. Koo, P. K. & Eddy, S. R. Representation learning of genomic sequence motifs with convolutional neural

- 326 networks. *PLOS Comput. Biol.* **15**, e1007560, [10.1371/journal.pcbi.1007560](https://doi.org/10.1371/journal.pcbi.1007560) (2019).
- 327 **17.** Koo, P. K., Qian, S., Kaplun, G., Volf, V. & Kalimeris, D. Robust Neural Networks are More Interpretable for
328 Genomics. Tech. Rep. (2019). [10.1101/657437](https://doi.org/10.1101/657437). Type: article.
- 329 **18.** Linder, J. *et al.* Interpreting neural networks for biological sequences by learning stochastic masks. *Nat. Mach.*
330 *Intell.* **4**, 41–54, [10.1038/s42256-021-00428-6](https://doi.org/10.1038/s42256-021-00428-6) (2022).
- 331 **19.** Morgunova, E. & Taipale, J. Structural perspective of cooperative transcription factor binding. *Curr. Opin.*
332 *Struct. Biol.* **47**, 1–8, [10.1016/j.sbi.2017.03.006](https://doi.org/10.1016/j.sbi.2017.03.006) (2017).
- 333 **20.** Rentzsch, R., Deneke, C., Nitsche, A. & Renard, B. Y. Predicting bacterial virulence factors – evaluation
334 of machine learning and negative data strategies. *Briefings Bioinforma.* **21**, 1596–1608, [10.1093/bib/bbz076](https://doi.org/10.1093/bib/bbz076)
335 (2020).
- 336 **21.** Sidorczuk, K. *et al.* Benchmarks in antimicrobial peptide prediction are biased due to the selection of negative
337 data. *Briefings Bioinforma.* **23**, bbac343, [10.1093/bib/bbac343](https://doi.org/10.1093/bib/bbac343) (2022).
- 338 **22.** Monti, R. & Ohler, U. Toward Identification of Functional Sequences and Variants in Noncoding DNA. *Annu.*
339 *Rev. Biomed. Data Sci.* **6**, null, [10.1146/annurev-biodatasci-122120-110102](https://doi.org/10.1146/annurev-biodatasci-122120-110102) (2023).
- 340 **23.** Kelley, D. R., Snoek, J. & Rinn, J. Basset: Learning the regulatory code of the accessible genome with deep
341 convolutional neural networks. *Genome Res.* gr.200535.115, [10.1101/gr.200535.115](https://doi.org/10.1101/gr.200535.115) (2016).
- 342 **24.** Avsec, v. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat.*
343 *Methods* **18**, 1196–1203, [10.1038/s41592-021-01252-x](https://doi.org/10.1038/s41592-021-01252-x) (2021).
- 344 **25.** Zhou, J., Guruvayurappan, K., Chen, H. V., Chen, A. R. & McVicker, G. Genome-wide analysis of CRISPR
345 perturbations indicates that enhancers act multiplicatively and without epistatic-like interactions. *bioRxiv*
346 [10.1101/2023.04.26.538501](https://doi.org/10.1101/2023.04.26.538501) (2023).
- 347 **26.** Castro-Mondragon, J. A. *et al.* JASPAR 2022: the 9th release of the open-access database of transcription factor
348 binding profiles. *Nucleic Acids Res.* gkab1113, [10.1093/nar/gkab1113](https://doi.org/10.1093/nar/gkab1113) (2021).
- 349 **27.** Zeitlinger, J. Seven myths of how transcription factors read the cis-regulatory code. *Curr. Opin. Syst. Biol.* **23**,
350 22–31, [10.1016/j.coisb.2020.08.002](https://doi.org/10.1016/j.coisb.2020.08.002) (2020).
- 351 **28.** Zhou, J. *et al.* Accurate genome-wide predictions of spatio-temporal gene expression during embryonic
352 development. *PLOS Genet.* **15**, e1008382, [10.1371/journal.pgen.1008382](https://doi.org/10.1371/journal.pgen.1008382) (2019).
- 353 **29.** Lee, N. K., Tang, Z., Toneyan, S. & Koo, P. K. EvoAug: improving generalization and interpretability of
354 genomic deep neural networks with evolution-inspired data augmentations. *Genome Biol.* **24**, 105, [10.1186/
355 s13059-023-02941-w](https://doi.org/10.1186/s13059-023-02941-w) (2023).
- 356 **30.** Consens, M. E. *et al.* To Transformers and Beyond: Large Language Models for the Genome. Tech. Rep. (2023).
357 ArXiv:2311.07621 [cs, q-bio] type: article.
- 358 **31.** Friedman, J. H. & Popescu, B. E. Predictive Learning via Rule Ensembles. *The Annals Appl. Stat.* **2**, 916–954

- 359 (2008).
- 360 **32.** Tsang, M., Enouen, J. & Liu, Y. Interpretable Artificial Intelligence through the Lens of Feature Interaction.
361 *arXiv:2103.03103 [cs]* (2021). ArXiv: 2103.03103.
- 362 **33.** Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular biology and
363 bioinformatics. *Bioinformatics* **25**, 1422–1423, [10.1093/bioinformatics/btp163](https://doi.org/10.1093/bioinformatics/btp163) (2009).
- 364 **34.** Stewart, A. J., Hannenhalli, S. & Plotkin, J. B. Why Transcription Factor Binding Sites Are Ten Nucleotides
365 Long. *Genetics* **192**, 973–985, [10.1534/genetics.112.143370](https://doi.org/10.1534/genetics.112.143370) (2012).
- 366 **35.** Kokhlikyan, N. *et al.* Captum: A unified and generic model interpretability library for PyTorch.
367 *arXiv:2009.07896 [cs, stat]* (2020). ArXiv: 2009.07896.

368 **Appendix**

Table 1. Motif data sets. The data set '*distinct_1*' was used for the evaluation of heterologous motifs while the dataset '*MEF2A*' represents homologous data sets.

distinct_1		MEF2A	
Motif name	Motif ID	Motif name	Motif ID
CREB1	MA0018.3		MA0052.1
NHLH1	MA0048.1	MEF2A	MA0052.2
ETS1	MA0098.3		MA0052.3
TEAD1	MA0090.3		MA0052.4

Table 2. Recall for CNNs based on data with negative sequences containing motif subsets (M_w) and without motif subsets (M_{wo}). While M_w is capable of distinguishing between positive sequences (4 motifs) and negative sequences (0-3), M_{wo} has low accuracies for sequences with 2 or 3 motifs since sequences with subsets of the interacting motif set were not present during training. The number of filters has no large influence on the accuracy.

(a) MEF2A					(b) Distinct 1				
MEF2A	large (32 filter)		small (4 filter)		distinct_1	large (32 filter)		small (4 filter)	
#Motifs	M_w	M_{wo}	M_w	M_{wo}	#Motifs	M_w	M_{wo}	M_w	M_{wo}
4	0.9832	0.9992	0.9772	0.9999	4	0.9307	0.9994	0.9485	0.9992
3	0.91	0.034	0.8965	0.04	3	0.6815	0.0615	0.7215	0.089
2	0.9973	0.3336	0.9963	0.4534	2	0.9893	0.6177	0.9647	0.6
1	1	0.91	1	0.961	1	1	0.988	0.9985	0.9795
0	1	1	1	1	0	1	0.998	1	1

Table 3. Subset recall for M_{wo} trained on the *distinct_1* data set. Each test sequence contains 2 motifs. The accuracies are calculated for each pairwise motif combination to see if the subsets meet the overall accuracy (Large model: 0.6177 and Small: 0.6) or if there are preferences in motifs. Low accuracy means that many sequences were predicted as positive and, therefore, that motif subset is evidence for the positive class for the model.

Motif IDs	CREB1	CREB1	CREB1	NHLH1	NHLH1	ETS1
	NHLH1	ETS1	TEAD1	ETS1	TEAD1	TEAD1
Large	0.628	0.822	0.92	0.15	0.458	0.728
Small	0.488	0.714	0.874	0.144	0.62	0.76

Table 4. Motif-wise AUPRC values for contribution scores for models trained on MEF2A homologous data set. The scores are shown for positive sequences containing all 4 interactive motifs

		MA0052.1	MA0052.2	MA0052.3	MA0052.4
DeepLIFT	w	1.0	0.7445	0.9006	0.7073
	wo	1.0	0.8085	0.9743	0.7429
Feature Permutation	w	0.9513	0.6713	0.7903	0.5942
	wo	0.9463	0.7787	0.9049	0.7165
Integrated Gradients	w	1.0	0.7268	0.8819	0.6698
	wo	1.0	0.8317	0.9933	0.7741

Table 5. AUPRC values for contribution scores for models trained on *distinct_1* heterologous data set for one motif (NHLH1). M_w has lower values, especially for negative data containing 1-3 motifs.

#Motifs	Integrated Gradients		Feature Permutation		DeepLIFT	
	M_w	M_{wo}	M_w	M_{wo}	M_w	M_{wo}
4	0.701	0.8749	0.6513	0.8454	0.7331	0.8661
3	0.2053	0.8529	0.1432	0.8124	0.3106	0.85
2	0.0765	0.8213	0.094	0.7998	0.1032	0.8279
1	0.0720	0.7881	0.1346	0.772	0.0789	0.8061

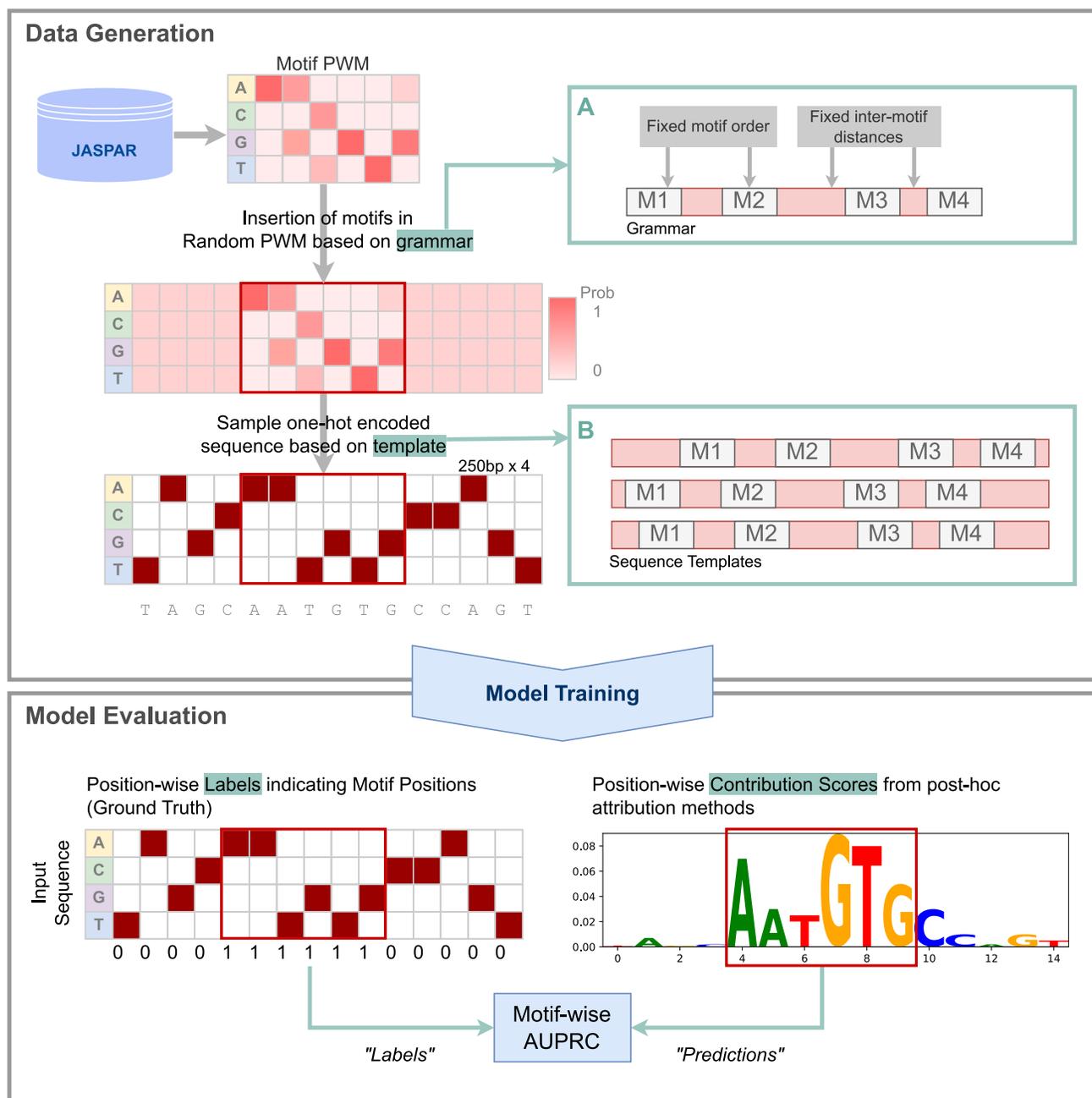


Figure 4. Data Generation We obtained the PWMs for 5 motif sets (see Table 1) from the JASPAR database. We create grammars for each set which consist of a specific motif order and distances between motifs (A). The presence of a motif depends on the investigated interaction (see methods section). One-hot-encoded sequence templates in the form of PWMs are generated for each input sequence from the grammar (B) from which the input sequence is then sampled. **Model Evaluation** A motif can be identified if the contribution scores higher than for random positions. To quantify how well a model captured a motif, we used an approach similar to [15]. Motif positions in the input sequence are labeled as positive, while random positions are labeled as negative. AUPRC values are then calculated based on those labels and the contribution scores similarly to prediction probabilities.

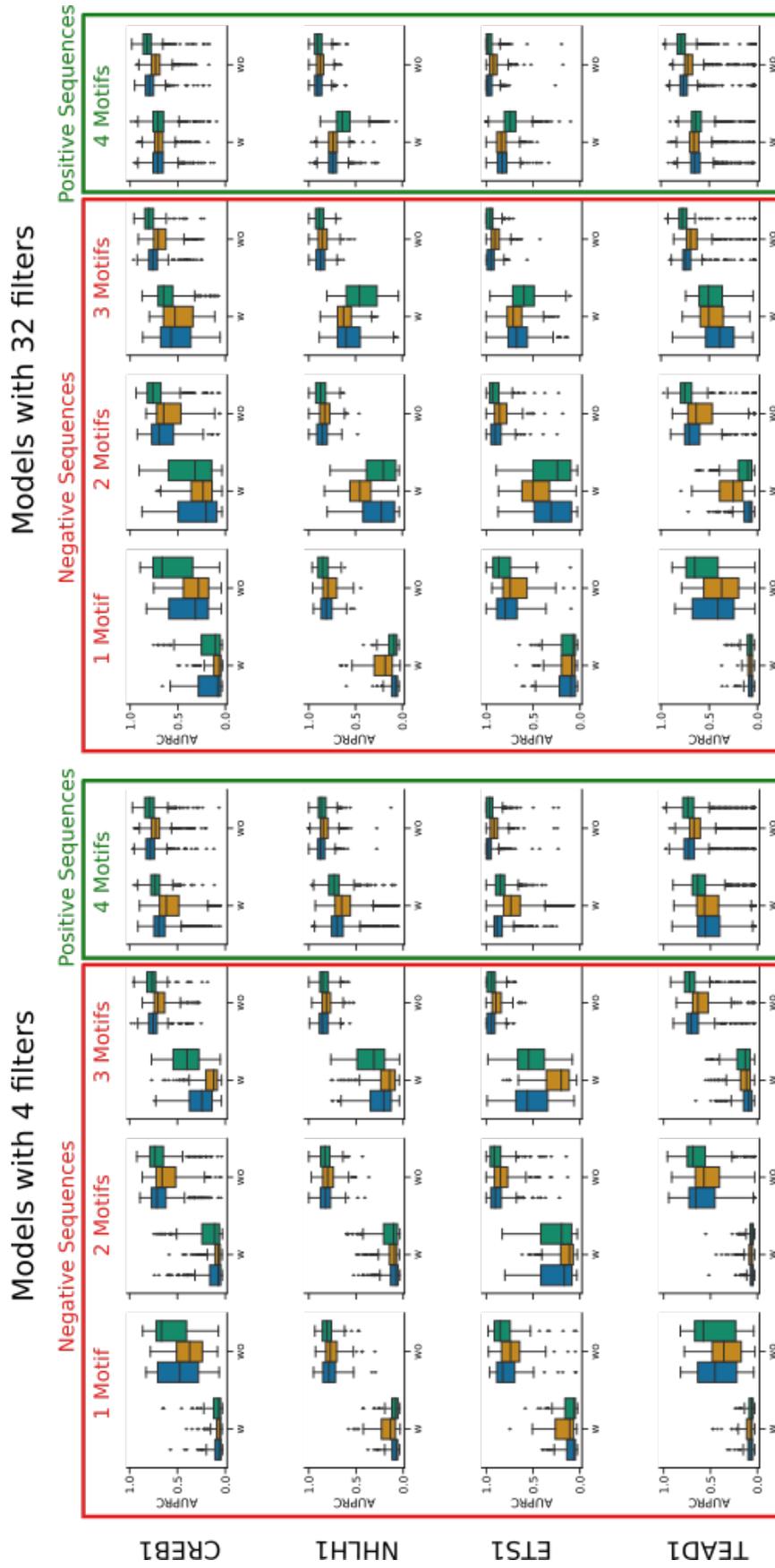


Figure 5. AUPRC scores for negative sequence experiment for distinct_1 motif set. AUPRC scores for the contribution scores of $M_{w/o}$ are lower for all motifs compared to $M_{w/o}$.

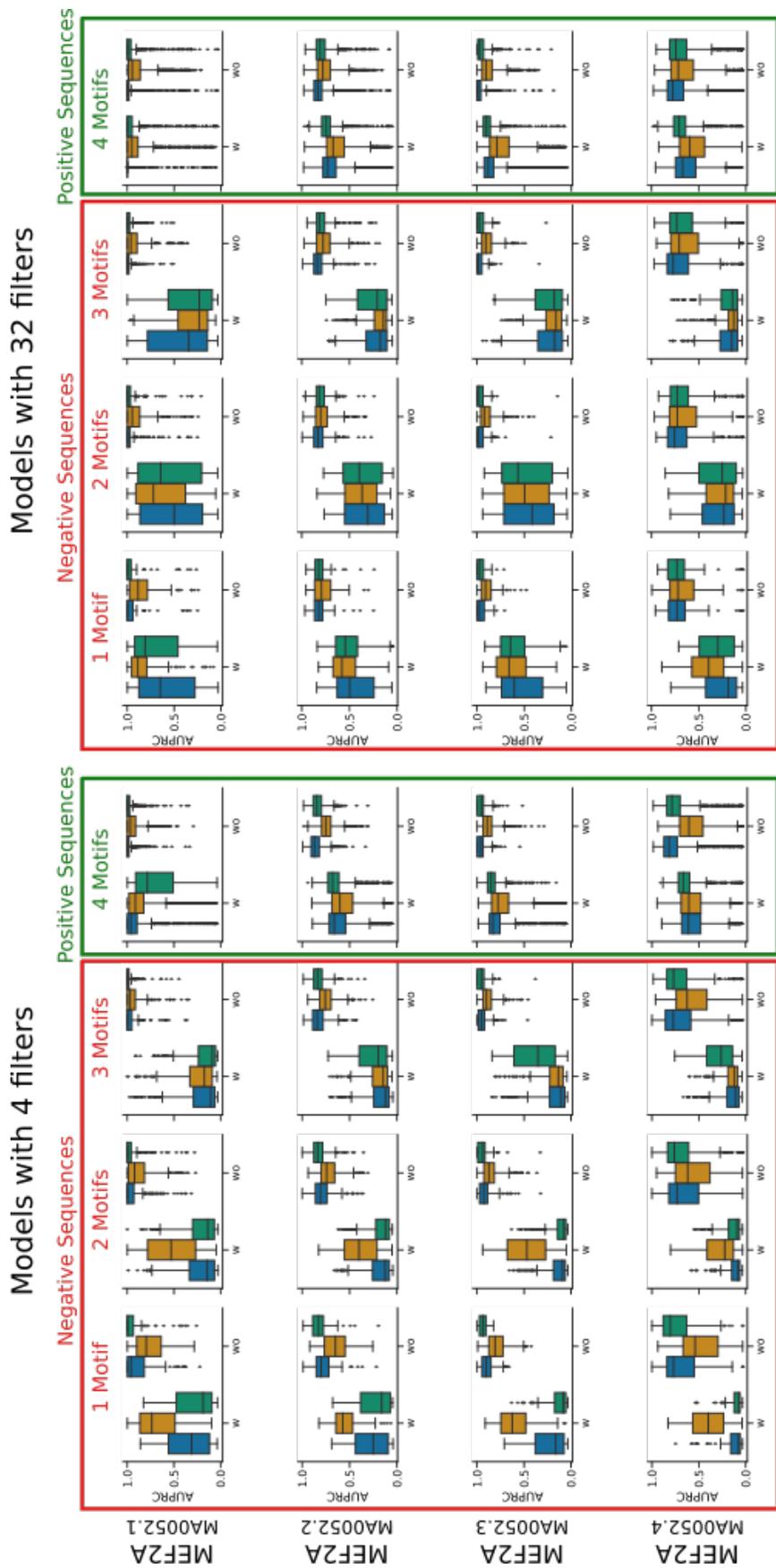


Figure 6. AUPRC scores for negative sequence experiment for MEF2A motif set. While the AUPRC values decrease with a lower number of motifs present in the case of heterologous motifs (see Supplement Fig. 5, AUPRC values are higher in the homologous case for the sequences containing 1 or 2 motifs compared to 3 since those sequences are closer to the decision boundary, especially for Feature Permutation).

Table 6. Subset recall and total accuracy for regression models trained on *distinct_1*.

distinct_1	large (32 filter)				small (4 filter)			
#Motifs	Additive	Enhancing	Inhibiting	Quadratic	Additive	Enhancing	Inhibiting	Quadratic
4	0.8876	0.9112	0.8452	0.9144	0.7032	0.7368	0.652	0.7944
3	0.8856	0.8871	0.8802	0.9041	0.8289	0.8207	0.6579	0.8709
2	0.8869	0.8777	0.8903	0.8731	0.8535	0.8113	0.8403	0.8567
1	0.8894	0.8895	0.8819	0.9073	0.8463	0.8321	0.7849	0.8483
0	0.6568	0.538	0.646	0.3588	0.4104	0.266	0.306	0.3228
Total	0.8729	0.8639	0.8599	0.8676	0.8085	0.7801	0.8209	0.7357

Table 7. Subset recall and total accuracy for regression models trained on *MEF2A*.

distinct_1	large (32 filter)				small (4 filter)			
#Motifs	Additive	Enhancing	Inhibiting	Quadratic	Additive	Enhancing	Inhibiting	Quadratic
4	0.9492	0.9576	0.9416	0.9576	0.9284	0.9328	0.9196	0.9256
3	0.9401	0.9457	0.9341	0.9458	0.9128	0.8956	0.9024	0.9036
2	0.9361	0.9245	0.9437	0.923	0.8992	0.877	0.9049	0.8832
1	0.9086	0.9105	0.909	0.9243	0.8795	0.8723	0.8688	0.8974
0	0.8612	0.8504	0.8588	0.74	0.8156	0.7656	0.8144	0.6484
Total	0.9190	0.9177	0.9174	0.8981	0.8871	0.8687	0.88201	0.8516