

Representation Learning of Genomic Sequence Motifs with Convolutional Neural Networks

Peter K. Koo¹, Sean R. Eddy^{1,2}

July 6, 2018

1. Howard Hughes Medical Institute, Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA
2. John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA

Abstract

Although convolutional neural networks (CNNs) have been applied to a variety of computational genomics problems, there remains a large gap in our understanding of how they work. Here we perform systematic experiments on synthetic sequences to reveal principles of how CNN architecture influences the internal representations of genomic sequence motifs that are learned. We focus our study on representations learned by first convolutional layer filters. We find that deep CNNs tend to learn *distributed* representations of partial sequence motifs. However, we demonstrate that the architecture of a CNN can be modified to predictively learn more interpretable *localist* representations, *i.e.* whole motifs. We then validate that the representation learning principles established from synthetic sequences generalize to *in vivo* sequences.

Introduction

Deep convolutional neural networks (CNNs) have recently been applied to predict transcription factor binding motifs from genomic sequences (Zhou and Troyanskaya, 2015; Quang and Xie, 2016; Kelley *et al*, 2016; Hiranuma *et al*, 2017). Despite the promise that CNNs bring in replacing methods that rely on k -mers and position weight matrices (PWMS) (Ghandi *et al*, 2016; Foat *et al*, 2006), there remains a large gap in our understanding of why CNNs perform well.

A convolutional layer of a CNN is comprised of a set of filters, each of which can be thought of as a PWM. Each filter scans across the inputs, and outputs a non-linear similarity score at each position, a so-called feature map. The filters are parameters of the CNN that are trained to detect relevant patterns in the data. Deep CNNs are constructed by feeding the feature maps of a convolutional layer as input to another convolutional layer. This can be repeated to create a network with any depth. CNNs typically employ *max-pooling* after each convolutional layer, which down-sample the feature maps by setting non-overlapping windows with a single maximum score, separately for each filter. Max-pooling enables deeper layers to detect features hierarchically across a larger spatial scale of the input. CNN predictions are typically made by feeding the feature map of the final convolutional layer through a fully-connected hidden layer followed by a linear classifier.

In genomics, it is commonly thought that the first convolutional layer filters learn PWMS of sequence motifs, while deeper layers learn combinations of these motifs, so-called regulatory grammars (Alipanahi *et al*, 2015; Angermueller *et al*, 2016; Zeng *et al*, 2016; Quang and Xie, 2016; Kelley *et al*, 2016). Thus, first convolutional layer filter sizes are usually chosen to be larger than the expected motifs, *i.e.* 6-19 nucleotides (nts) (Mathelier *et al*, 2016). A quantitative motif comparison search typically reveals that less than 50% of

the first layer filters find any statistical match against a database of motifs (Kelley *et al*, 2016; Quang and Xie, 2016). Unmatched filters have been suggested to be either partial representations of known motifs or novel motifs, *i.e.* motifs not included in the database. It remains unclear to what extent we should expect first layer filters to learn whole-motif representations in the first convolutional layer.

Learning whole-motif representations by first layer filters is not indicative of a deep CNN's classification performance. For instance, a deep CNN that employs a small first layer filter, *i.e.* 8 nts (Zhou and Troyanskaya, 2015), which is shorter than many common motifs found *in vivo*, has demonstrated comparable performance as CNNs that employ larger filters, *i.e.* ≥ 19 nts (Quang and Xie, 2016; Kelley *et al*, 2016). In principle, smaller filters that capture partial motif representations can be combined in deeper layers to assemble whole-motif representations, thereby allowing the CNN to make accurate predictions.

Here we perform systematic experiments to demonstrate that a CNN's architecture, specifically max-pooling and filter size, is indicative of how internal representations of motifs are learned. We focus our study on the representations learned by first convolutional layer filters using a synthetic dataset with a known ground truth. By systematically modifying the architecture, we learn general principles that are predictive of the extent that first layer filters learn representations that resemble whole-motifs versus partial-motifs. We then demonstrate that the same principles learned from synthetic sequences generalize to *in vivo* sequences.

Results and Discussion

Internal representations of motifs depend on architecture

We conjecture that motif representations learned in first layer filters are largely influenced by a CNN's ability to assemble whole-motif representations in deeper layers, which is determined by architectural constraints set by: 1. the convolutional filter size, 2. the stride of the filter, which is the number of steps the filter takes (usually set to 1), 3. the max-pool size, and 4. the max-pool stride, which is the number of steps for each max-pool window (usually set to the max-pool size to create non-overlapping max-pool windows).

Assuming that accurate classification can only be made if the correct motifs are detected, a CNN that learns partial-motif representations in the first layer must assemble whole-motif representations at some point in deeper layers. To help explain how architecture can influence representation learning in a given layer, we introduce the concept of a receptive field, which is the sensory space of the data that affects a given neuron's activity. For the first convolutional layer, each neuron's receptive field has a size that is equal to the filter size at a particular region of the data. Each neuron's receptive field is only activated to an extent that depends on the similarity of a sequence and a given filter. Since there are typically many filters in a convolutional layer, there are many neurons which have a receptive field that share the same spatial region. However, each neuron's activation is determined by a different filter. Max-pooling combines multiple neurons of a given filter within a specified window size to a single max-pooled neuron, thereby augmenting the size of its receptive field. In doing so, max-pooling obfuscates the exact positioning of the max-activation within each window. Thus the location of the max-activation has spatial invariance within its receptive field with an amount equal to the max-pool size.

Although max-pooling creates spatial uncertainty of the max-activation within a max-pooled neuron's receptive field, we surmise that neighboring max-pooled neurons of different filters, which share significantly overlapping receptive fields, can help to resolve spatial positioning of an activation. To illustrate, Figure 1A shows a toy example of two convolutional filters, each 7 nts long, which have learned partial-motifs: 'GTG' and 'CAC'. An example sequence contains three embedded patterns (highlighted in green): 'CACGTG', 'GTG CAC', and 'CACNNNGTG', where 'N' represents any nucleotide with equal probability. The resultant max-pooled, activated convolutional scans for each filter are shown above the sequence with a blue shaded region highlighting the receptive field of select max-pooled neurons. Even though the first convolutional layer filters have learned partial-motifs, the second convolutional layer filters can still resolve each of the three embedded patterns by employing filters of length 3. Of course situations may arise where the three second convolutional layer filters are unable to fully resolve the embedded patterns with fidelity. For instance, 'CACNGTG' could be activated by the same filter for 'CACGTG'. A CNN can circumvent these ambiguous situations by either learning more information about each pattern within each filter or by dedicating additional filters to help discriminate the ambiguous patterns.

It follows that by creating a situation where partial-motif representations cannot be assembled into whole-motifs in deeper layers, learning whole-motifs by first layer filters becomes obligatory for accurate classification. One method to limit the information flow through a CNN is by employing large max-pool sizes relative to the filter size. The max-pooled neurons then have large receptive fields with a large spatial uncertainty and only a small overlap in receptive fields with neighboring neurons of different filters. A deeper layer would be unable to resolve the spatial ordering of partial motifs to assemble whole-motifs with fidelity. To exemplify, figure 1B shows a toy example of a CNN that employs a larger pool size of 20. Importantly, there are large spatial regions within a receptive field for which a neighboring neuron cannot help to resolve due to a lack of overlap in receptive fields. As a result, deeper convolutional layer filters which are dedicated to each pattern would yield the same signature, unable to resolve any of the three patterns.

More technically, the extent of motif information that each filter learns is guided by the gradients of the objective function, which serves as a measure of the classification error. Assuming accurate classification can only be achieved upon discriminating the underlying motifs of each class, once whole-motifs for each class are learned, then the objective function is minimized and the training gradients go to zero. If a CNN can build whole-motifs from partial-motifs in deeper layers, then there is no more incentive to learn additional

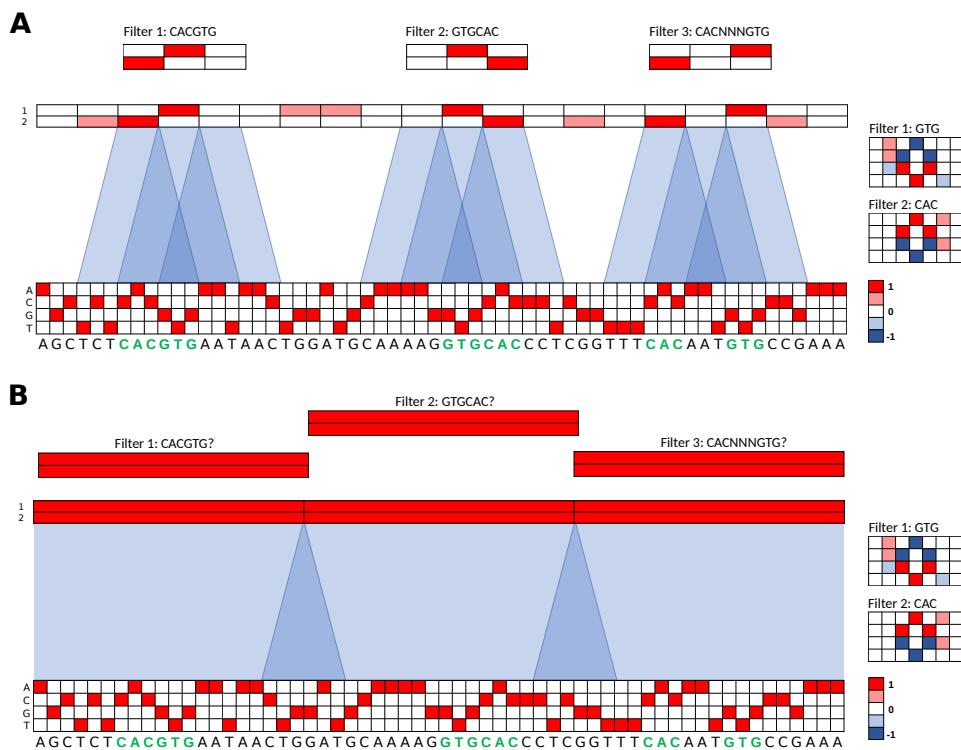


Figure 1: Toy model for representation learning of sequence motifs. (A,B) An example 60 nt one-hot encoded sequence contains 3 patterns (shown in green): CACGTG, GTGCAC, and CACNNNTG. Two filters, each of length 7 (7 columns and 4 rows, one for each nucleotide), are shown to the right. A partial-motif representation has been captured by each filter: GTG for filter 1 (Top) and CAC for filter 2 (Bottom). The max-pooled feature maps are shown above the sequence. The feature maps have the same size as the sequence by adding 3 zero-padding units to each end of the sequence prior to convolution (not shown in diagram). (A) Shows the feature maps when employing a small max-pooling size of 3, which creates overlapping receptive fields, highlighted in blue. 3 second layer convolutional filters, shown above, demonstrate a feature map pattern that can resolve each embedded sequence pattern. (B) Shows the feature maps when employing a larger pooling size of 20 using the same filters as (A). The larger receptive fields have a large spatial uncertainty along with a small overlap in receptive fields from neighboring neurons. Each of the 3 second layer convolutional filters, shown above, is unable to find a unique feature map pattern that can resolve any embedded sequence pattern.

information to build upon the partial-motif representations already learned. As a result, the first layer filters will maintain a *distributed* representations of motifs (Hinton *et al*, 1986). On the other hand, if architectural constraints limit the ability to build whole-motifs from partial-motifs in deeper layers, then accurate predictions cannot be made. Hence, gradients will persist because the objective function is not yet minimized, encouraging first layer filters to learn whole-motifs, also known as a *localist* representation of motifs (Hinton *et al*, 1986). Once the first layer filters have learned sufficient information of whole-motifs to discriminate each class, then the objective function can be minimized, signaling the end of training.

Max-pooling influences ability to build hierarchical motif representations

To test this idea, we created a synthetic dataset for the task of predicting which transcription factors bind to a given sequence. Briefly, synthetic sequences, each 200 nts long, were implanted with 1 to 5 known motifs, randomly selected with replacement from a pool of 12 transcription factor motifs embedded in random DNA (see Methods for details). The motifs were manually selected from the JASPAR database to represent a diverse, non-redundant set. The goal of this computational task is to simultaneously make 12 binary predictions for the presence or absence of each transcription factor motif in the sequence. Since we have ground truth for all of the relevant TF motifs and where they are embedded in each sequence, we can test the efficacy of the representations learned by a trained CNN model. We note that the ground truth is only from embedded motifs and not from motifs that occasionally arise by chance; the latter effectively creates false negative labels in this dataset.

A CNN model that employs at least two convolutional layers is required to test our hypotheses of representation learning. We constructed a CNN with 3 hidden layers: two convolutional layers, each followed by max-pooling, and a fully-connected hidden layer. Specifically, our CNN takes as input one-hot encoded sequences, processes them with the hidden layers, and outputs a prediction for the binding probability for each of the 12 classes. The number of filters in each convolutional layer, the number of units in the fully-connected hidden layer, and the dropout probabilities are fixed (see Methods). The filter sizes, the max-pool window sizes, and the max-pool strides are the hyperparameters that can be varied. For a given hyperparameter setting, we trained the CNN as a multi-class logistic regression (see Methods for training details). All reported metrics are strictly drawn from the held-out test set using the model parameters that yielded the best performance on the validation set.

To explore how spatial uncertainty within receptive fields set by max-pooling influences the representations learned by first layer filters, we systematically altered the max-pool sizes while keeping all other hyperparameters fixed, including a first and second layer filter size of 19 and 5, respectively. To minimize the influence of architecture on classification performance, we coupled the max-pool sizes between the first and second layer, such that their products are equal, which makes the inputs into the fully-connected hidden layer the same size across all CNNs. The max-pool sizes we employed are (first layer, second layer): (2, 50), (4, 25), (10, 10), (25, 4), (50, 2), and (100,1). For brevity, we denote each CNN with only the first max-pool window size, *e.g.* CNN-2 for (2, 50).

We first verified that the performance of each model is similar as measured by the average area-under the receiver-operator-characteristic (AU-ROC) curve across the 12 classes (Table 1), which is in the range of previously reported values for a similar task using experimental ChIP-seq data (Zhou and Troyanskaya, 2015; Quang and Xie, 2016). Next, we converted each filter to a sequence logo to visually compare the motif representations learned by the first layer filters across the different architectures (Fig. 2). As expected, we found CNNs that employ large max-pool sizes (≥ 10) learn representations that qualitatively resemble the ground truth motifs. On the other hand, CNNs that employ a small max-pool size (≤ 4) do not seem to qualitatively capture any ground truth motif in its entirety, perhaps learning, at best, parts of a motif.

To quantify the number of filters that have learned motifs, we employed the Tomtom motif comparison search tool (Gupta *et al*, 2007) to compare the similarity of each filter against all motifs in the JASPAR 2016 vertebrate database (Mathelier *et al*, 2016) using an *E*-value cutoff of 0.1. In agreement with our qualitative observation, we found CNNs that employ a small max-pool size (≤ 4) have, at best, 33% of their filters match any known motifs. Of these, only 1 filter in CNN-4 matches a ground truth motif. In contrast, CNNs that employ a large max-pool size yield, at worst, a 90% match to ground truth motifs.

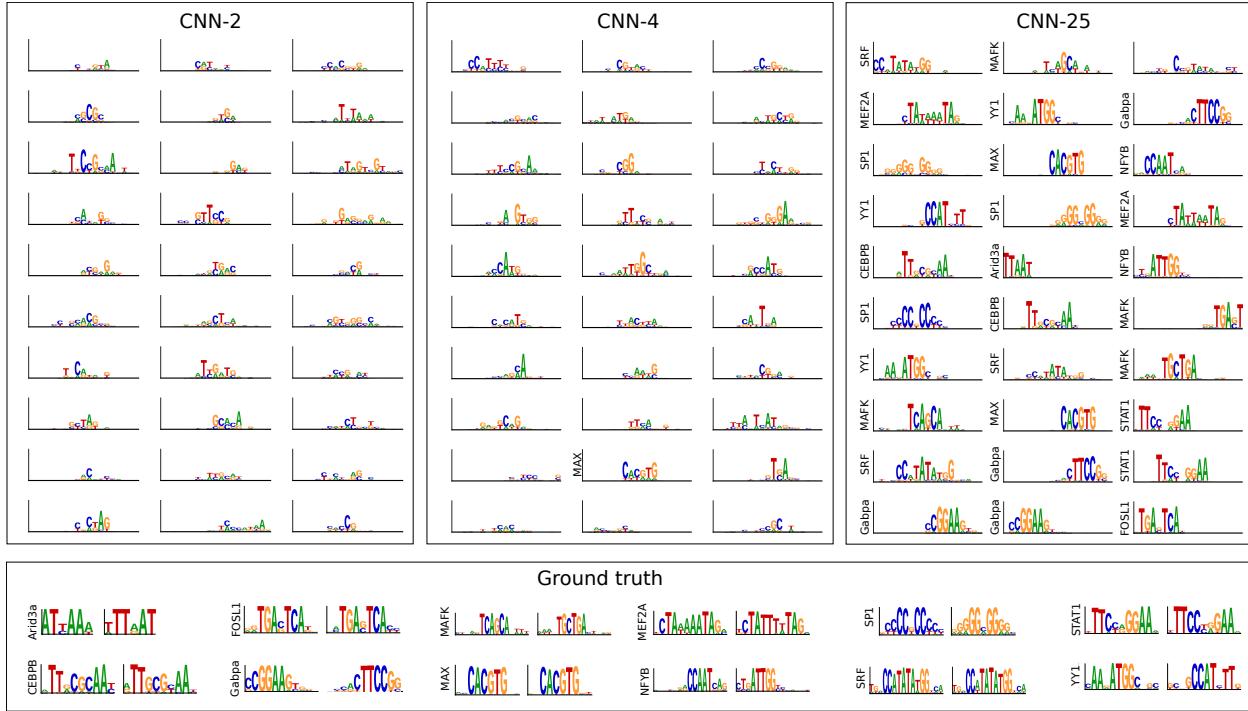


Figure 2: Comparison of first layer filters for CNNs with a different max-pool size. Sequence logos for normalized first convolutional layer filters are shown for CNN-2 (Left), CNN-4 (Middle), and CNN-25 (Right). The sequence logo of the ground truth motifs and its reverse complement for each transcription factor is shown at the bottom. The y-axis label on select filters represents a Tomtom match to a ground truth motif.

Table 1: Performance on the synthetic dataset. The table shows the average area under the receiver-operator-characteristic curve (AU-ROC) across the 12 TF classes, percentage of matches between the 30 first convolutional layer filters and the entire JASPAR vertebrates database (JASPAR), and the percentage of filters that match to any ground truth TF motif (Relevant) for different CNN models. Motif matches were determined by the Tomtom motif comparison search tool using an *E*-value cutoff of 0.1. The error bar in the Average AU-ROC represents the standard deviation of the AU-ROC across the 12 classes.

Model	Average AU-ROC	% Motif match (JASPAR)	% Motif match (Relevant)
CNN-2	0.968±0.028	26	0
CNN-4	0.952±0.056	33	3
CNN-10	0.970±0.038	97	97
CNN-25	0.954±0.094	97	97
CNN-50	0.965±0.040	97	93
CNN-100	0.961±0.042	90	90
CNN ₉ -4	0.958±0.039	10	3
CNN ₉ -25	0.961±0.038	87	80
CNN ₃ -2	0.968±0.027	3	0
CNN ₃ -50	0.652±0.060	13	0
CNN-50-2	0.917±0.136	87	70

Motif representations are not very sensitive to filter size

Because it has been widely thought that first convolutional layer filters learn motifs, deep learning practitioners have traditionally employed CNN architectures with large first layer filters to capture motif patterns in

their entirety. However, we have shown that employing a large filter does not necessarily lead to whole-motif representations. To test the sensitivity of filter size to representation learning, we created two new CNN models that employ a first layer filter size of 9 (CNN_9), in contrast to a filter size of 19 which was previously used, with max-pool combinations of 4 and 25, *i.e.* $\text{CNN}_9\text{-}4$ and $\text{CNN}_9\text{-}25$. Since the combination of a filter size of 9 with a max-pool size of 4 creates overlapping receptive fields with a small spatial uncertainty, we expect that this architecture setting will lead to partial-motif representations. On the other hand, the filter size of 9 is insufficient to resolve spatial positions when employing a max-pool size of 25. Hence, we predict that this architecture setting will yield whole-motif representations. As expected, $\text{CNN}_9\text{-}25$ learns representations that qualitatively better reflect the ground truth motifs compared to $\text{CNN}_9\text{-}4$ (Fig. 3, A-B). Interestingly, $\text{CNN}_9\text{-}25$ also learns partial motif representations of larger motifs, *i.e.* MEF2A, SRF, STAT1, CEBPB, but in a more visually identifiable way compared to $\text{CNN}_9\text{-}4$. By quantifying the percentage of filters that statistically match ground truth motifs, $\text{CNN}_9\text{-}25$ yields an 80% match compared to $\text{CNN}_9\text{-}4$ which only yields a single match (Table 1).

As a control, we created a CNN model with a filter size of 3 with max-pool size combinations of 2 and 50, *i.e.* $\text{CNN}_3\text{-}2$ and $\text{CNN}_3\text{-}50$. Since a max-pool size of 2 is smaller than the filter size, we expect that $\text{CNN}_3\text{-}2$ will still be able to assemble whole motifs to some extent in deeper layers, despite having a very small filter size. On the other hand, since $\text{CNN}_3\text{-}50$ has only one chance to learn whole motifs, we expect that the small filter size will lead to a poor classification performance. Indeed, $\text{CNN}_3\text{-}50$ yields a mean AU-ROC of 0.652 ± 0.060 across the 12 classes, compared to $\text{CNN}_3\text{-}2$ which yields 0.968 ± 0.039 (error is the standard deviation across the 12 classes).

Spatial uncertainty within receptive fields determines motif representations

One aspect of max-pooling that we did not consider in our toy model is the max-pool stride, which is typically set to the max-pool size. Employing a large max-pool size with a small max-pool stride can create

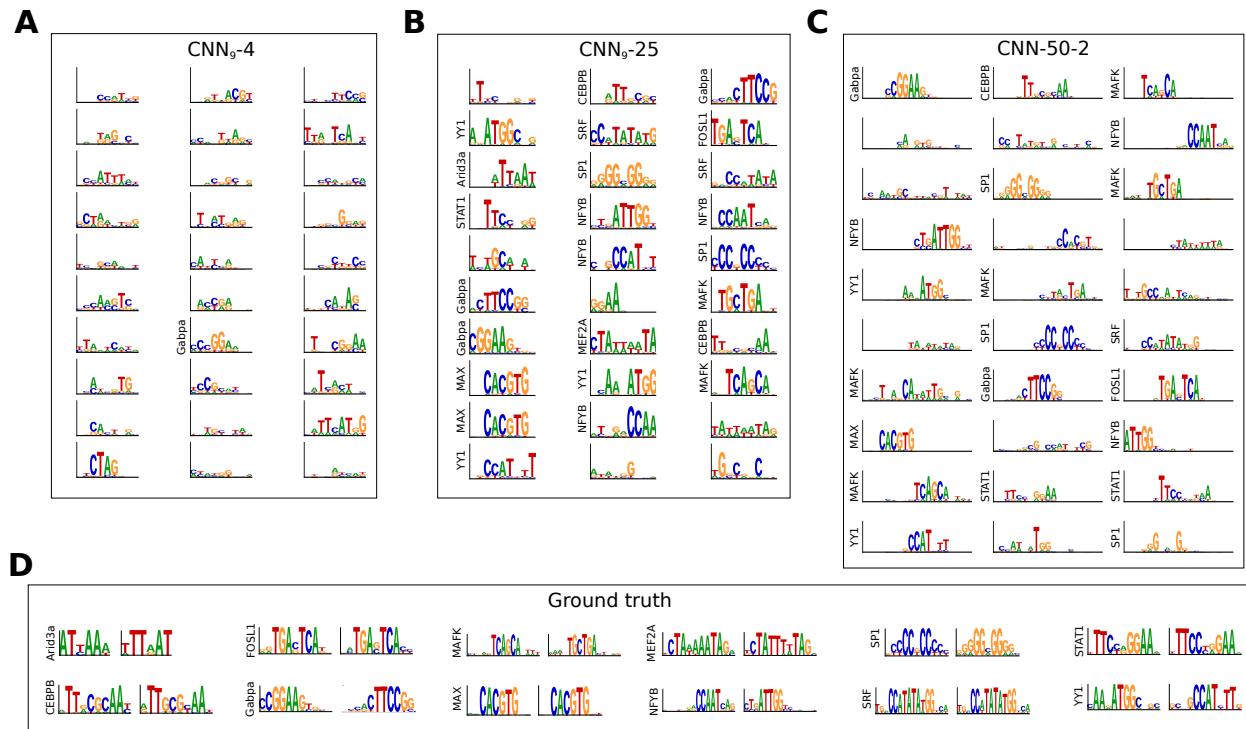


Figure 3: Representations learned by first layer filters for alternative CNN architectures. Sequence logos for normalized first convolutional layer filters are shown for (A) $\text{CNN}_9\text{-}4$ (B) $\text{CNN}_9\text{-}25$, and (C) $\text{CNN}\text{-}50\text{-}2$. (D) shows the ground truth motifs and its reverse complement for each transcription factor. The y-axis label on select filters represents a Tomtom match to a ground truth motif.

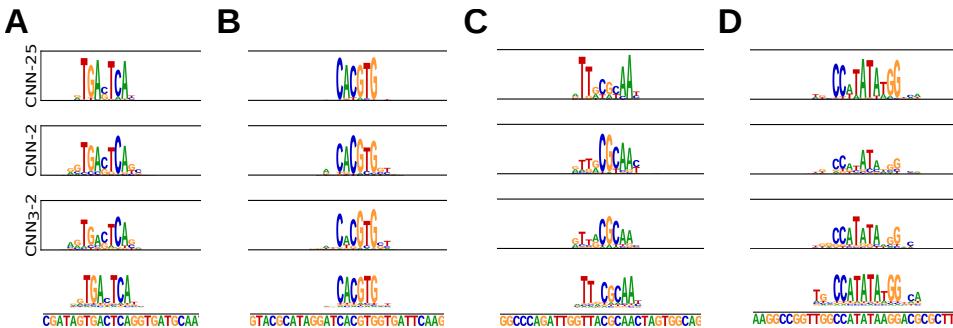


Figure 4: Representative sequence logos of the saliency maps generated by CNNs. Sequence logos of the saliency maps generated by CNN-25, CNN-2, and CNN₃-2 for a sequence that contains a label for (A) FOSL1, (B) MAX, (C) CEBPB, and (D) SRF. The underlying sequence and the sequence logo of the sequence model (ground truth) is shown below. A saliency map was generated by employing guided-backprop from the pre-activated output neuron with respect to the input layer. The saliency map was then normalized to a PWM and converted to a sequence logo. Each saliency map was further clipped about the embedded motif for brevity. The unclipped saliency plots are shown in Supplemental Fig. S1.

a situation where the receptive field of max-pooled neurons overlap significantly, which should improve the spatial resolution of partial-motifs. However, a deeper convolutional filter would still be unable to assemble whole motifs, because each receptive field has a large spatial uncertainty, making it challenging to discriminate between partial-motifs that assemble into whole motifs from partial-motifs that are spatially distant.

To test this, we created a new CNN model which employs a large max-pool size of 50 with a max-pool stride of 2 (CNN-50-2). Consequently, the length of the feature maps after the first convolutional layer are half of the input sequence, which is the same feature map shape as CNN-2, which employs a max-pool size of 2 with a stride of 2. Similar to CNN-2, CNN-50-2 employs a max-pool size and stride of 50 after the second convolutional layer. Strikingly, CNN-50-2 learns whole-motif representations with 70% of its filters matching ground truth motifs in the synthetic dataset (Table 1). Moreover, the motifs that are learned by CNN-50-2 qualitatively better resemble whole-motif representations (Fig. 3C) compared to CNN-2 (Fig. 2). Together, this result further supports that architecture, specifically the ability to assemble whole-motifs in deeper layers, plays a major role in how CNNs learn genomic representations in a given layer.

Distributed representations build whole-motif representations in deeper layers

The high overall classification performance of the CNNs suggests that they must have learned whole-motif representations of each embedded TF at some point. Thus, CNN-2, whose first layer filters did not match any relevant motifs, must be assembling whole-motif representations in deeper layers. To verify that CNN-2 eventually learns whole-motif representations, we visualize the representation learned throughout the entire network with saliency analysis, specifically guided-backpropagation (Springenberg *et al*, 2014), which is a technique to probe the independent importance of each nucleotide in a sequence towards a given prediction (see Methods). By visualizing a representative sequence logo of a saliency map generated by CNN-2 and CNN-25 for sequences associated with different TF classes, we confirm that the underlying motif representations are indeed learned irrespective of whether the first layer learns whole-motifs or partial-motifs (Fig. 4). We note that the quality of the saliency maps generated by CNN-2 can occasionally lead to noisier importance scores for nucleotide variants in the motif compared to CNN-25, which tends to better reflect the embedded motifs. Interestingly, we also show that CNN₃-2 is also able to largely learn representations of whole-motifs, despite employing a very small first layer filter size of 3 (Fig. 4).

Generalization to *in vivo* sequences

To test whether the same representation learning principles generalize to *in vivo* sequences, we modified the DeepSea dataset (Zhou and Troyanskaya, 2015) to include only *in vivo* sequences that have a peak called

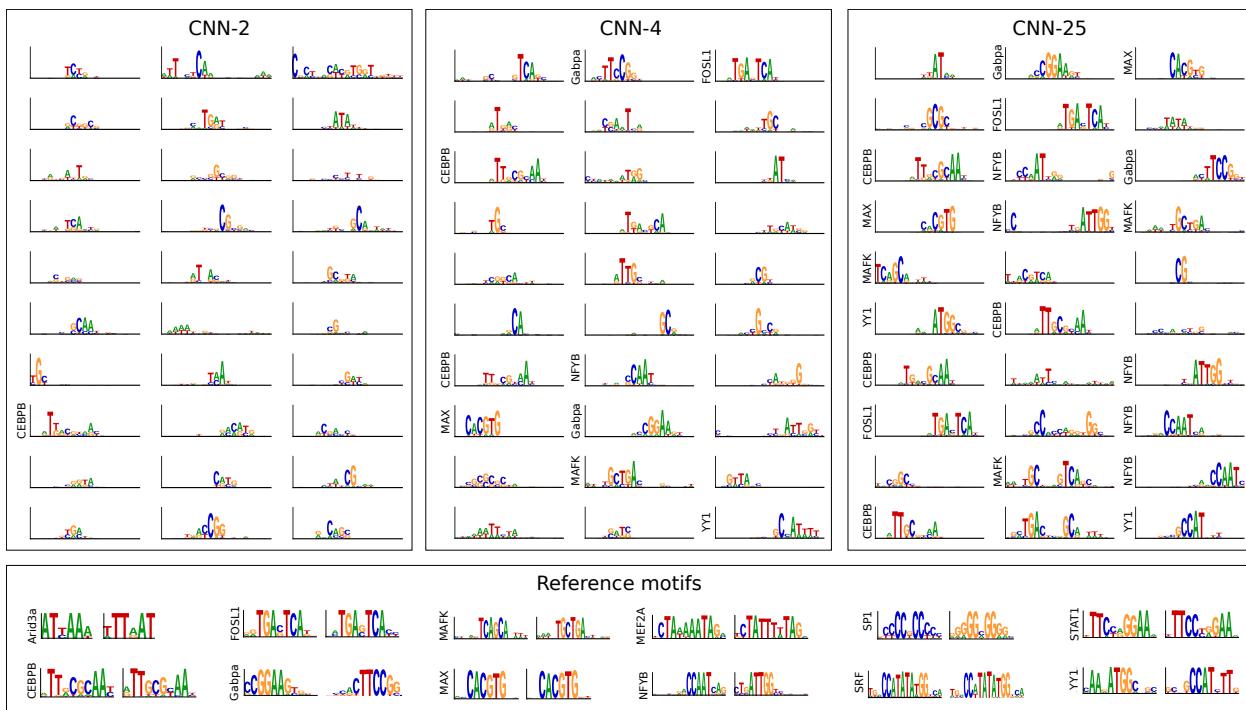


Figure 5: Comparison of the first layer filters for CNN models trained on *in vivo* sequences. Sequence logos for normalized first convolutional layer filters are shown for CNN-2 (Left), CNN-4 (Middle), and CNN-25 (Right). The sequence logos of reference motifs and their reverse complements for each transcription factor from the JASPAR database is shown at the bottom. The y-axis label on select filters represents a Tomtom match to a reference motif.

across 12 ChIP-seq experiments, each of which correspond to a TF in the synthetic dataset (see Supplemental Table S1). Thus, the truncated-DeepSea dataset is similar to the synthetic dataset with sequences that have a corresponding label for the presence or absence of a peak across the 12 ChIP-seq experiments. The truncated-DeepSea dataset consists of 1,000 nt sequences in contrast to the 200 nt sequences in the synthetic dataset.

We trained each CNN model on the *in vivo* dataset following the same protocol as the synthetic dataset. Similar to CNN models trained on the synthetic dataset, a qualitative comparison of the first layer filters of different CNN models shows that employing a larger pool size yield representations that better reflect whole-motifs (Fig. 5). By employing the Tomtom motif comparison search tool, we quantified the percentage of significant hits between the first layer filters against the JASPAR database (see Table 2). Similar to the synthetic dataset, CNNs that employ a smaller max-pool size (≤ 4) yield a percent match that is, at best, 57% (Table 2). In contrast, CNNs that employ a larger max-pool size (≥ 10) yields a percent match that is, at worst, 83% (Table 2). Since *in vivo* sequences contain many additional signals compared to the synthetic sequences, we were unable to reliably quantify the percentage of filters that learn *relevant* motifs. Interestingly, we found that each CNN was consistently unable to identify *known* motifs for ARID3A, MEF2A, SP1, and STAT1. However, it is unclear whether this arises because of: experimental or post-processing errors which create label noise in the sequences we assign as having a ChIP-seq peak, the large variance in numbers of sequences for different classes (class imbalance), and/or an inability of the CNNs to learn the correct motif, among the many other possible explanations. Notwithstanding, the same trends in the amount of motif information learned by first layer filters *in vivo* suggests that we have identified a general principle for representation learning by CNNs.

Table 2: Performance of deep learning models on the *in vivo* dataset. The table shows the percentage of matches between the first convolutional layer filters and the JASPAR database (JASPAR) and the percentage of filters that match known motifs for the 12 transcription factors. Motif matches were determined by the Tomtom motif comparison search tool using an *E*-value cutoff of 0.1.

Model	Motif match (JASPAR)	Motif match (Relevant)
CNN-2	30	3
CNN-4	57	30
CNN-10	83	57
CNN-25	83	70
CNN-50	83	67
CNN-100	87	70
CNN ⁹ -4	50	17
CNN ⁹ -25	77	57
CNN-50-2	87	50

Conclusion

By exploring different CNN architectures on a synthetic dataset with a known ground truth, we were able to reveal principles of how architecture design influences representation learning of sequence motifs. Typical deep CNN architectures currently employed in genomics, which employ large filters and small max-pool sizes, tend to learn *distributed* representations of sequence motifs. However, *localist* representations, *i.e.* whole motifs, can be learned by constraining the ability of deeper layers from assembling hierarchical representations of motifs. While we explored the role of architecture, we note that there may be other factors that contribute to the quality of the learned representations, including regularization and optimization algorithms.

Interpreting the representations learned by convolutional filters should be approached with skepticism. Even though a CNN can be designed to preferentially learn whole-motifs in the first convolutional layer, not all filters will learn motifs. Moreover, we showed that using a motif comparison tool does not necessarily provide a reliable way of identifying whether a CNN learns *relevant* motifs. Another aspect of convolutional filters that is often misleading is that the number of filters dedicated to a motif may not be a reliable measure of the importance of the motif. The variation in the number of filters dedicated to a motif on the synthetic sequences, which do not have any class imbalance, suggests that the observed difference is more likely due to random initialization and the difficulty of finding that motif, not the importance of the motif.

The similar performance across the CNNs explored here suggest that motif discovery does not require complicated architectures that learn distributed representations. Nevertheless, we posit that building *distributed* representations may be more beneficial in more complicated tasks, because a wider array of representations can be constructed through combinatorics of partial representations. Moreover, there becomes less dependence on convolutional filter lengths and numbers of filters as long as there exist deeper layers that can build representations hierarchically. In contrast, building *localist* representations means that the CNN is subject to harder constraints set by the numbers of filters and the filter lengths, limiting the amount of representations and their sizes that can be learned. However, when the main features in the dataset is simple, such as whether or not a motif is present, then CNN architectures that learn localist representations achieve an easier to interpret model that still performs competitively.

Methods

Synthetic dataset

The synthetic dataset consists of sequences with known motifs embedded in random DNA sequences to mimic a typical multi-class binary classification task for ChIP-seq datasets. We acquired a pool of 24 PWMs from 12 unique transcription factors (forward and reverse complements) from the JASPAR database (Mathelier

et al., 2016): Arid3a, CEBPB, FOSL1, Gabpa, MEF2A, MAFK, MAX, MEF2A, NFYB, SP1, SRF, STAT1, and YY1. For each sequence, we generated a 200 nt random DNA sequence model with equal probability for each nucleotide. 1 to 5 TF PWMs were randomly chosen with replacement and randomly embedded along the sequence model such that each motif has a buffer of at least 1 nucleotide from other motifs and the ends of the sequence. We generated 25,000 sequence models and simulated a single synthetic sequence from each model. A corresponding label vector of length 12, one for each unique transcription factor, was generated for each sequence with a one representing the presence of a TF's motif or its reverse complement along the sequence model and zero otherwise. The 25,000 synthetic sequences and their associated labels were then randomly split into a training, validation, and test set according to the fractions 0.7, 0.1, and 0.2, respectively.

In vivo dataset

Sequences which contain ENCODE ChIP-seq peaks were downloaded from the DeepSEA dataset via (Zhou and Troyanskaya, 2015). The human reference genome (GRCh37/hg19) was segmented into non-overlapping 200 nt bins. A vector of binary labels for ChIP-seq peaks and DNase-seq peaks was created for each bin, with a 1 if more than half of the 200 nt bin overlaps with a peak region, and 0 otherwise. Adjacent 200 nt bins were then merged to 1,000 nt lengths and their corresponding labels were also merged. Chromosomes 8 and 9 were excluded from training to test chromatin feature prediction performances, and the rest of the autosomes were used for training and validation. We truncated the DeepSea dataset to include only the sequences which contain 12 transcription factor labels: Arid3a, CEBPB, FOSL1, Gabpa, MEF2A, MAFK, MAX, MEF2A, NFYB, SP1, SRF, STAT1, and YY1 (See Supplementary Table S1 for ENCODE filenames and class indices from the original DeepSea dataset). 270,382 (92%) sequences comprise the training set and 23,768 (8%) sequences comprise the test set. Each 1000 nt DNA sequence is one-hot encoded into a 4x1000 binary matrix, where rows correspond to A, C, G and T.

CNN Models

All CNN models take as input a 1-dimensional one-hot-encoded sequence with 4 channels (one for each nucleotide: A, C, G, T), then processes the sequence with two convolutional layers, a fully-connected hidden layer, and a fully-connected output layer with 12 output neurons that have sigmoid activations for binary predictions. Each convolutional layer consists of a 1D cross-correlation operation, which calculates a running sum between convolution filters and the inputs to the layer, followed by batch normalization (Ioffe and Szegedy, 2015), which independently scales the features learned by each convolution filter, and a non-linear activation with a rectified linear unit (ReLU), which replaces negative values with zero.

The first convolutional layer employs 30 filters each with a size of 19 and a stride of 1. The second convolutional layer employs 128 filters each with a size of 5 and a stride of 1. All convolutional layers incorporate zero-padding to achieve the same output length as the inputs. Each convolutional layer is followed by max-pooling with a window size and stride that are equal, unless otherwise stated. The product of the two max-pooling window sizes is equal to 100. Thus, if the first max-pooling layer has a window size of 2, then the second max-pooling window size is 50. This constraint ensures that the number of inputs to the fully-connected hidden layer is the same across all models. The fully-connected hidden layer employs 512 units with ReLU activations.

Dropout (Srivastava *et al.*, 2014), a common regularization technique for neural networks, is applied during training after each convolutional layer, with a dropout probability set to 0.1 for convolutional layers and 0.5 for fully-connected hidden layers. During training, we also employed L₂-regularization with a strength equal to 1e-6. The parameters of each model were initialized according to (He *et al.*, 2015), more commonly known as He initialization.

All models were trained with mini-batch stochastic gradient descent (mini-batch size of 100 sequences) for 100 epochs, updating the parameters after each mini-batch with Adam updates (Kingma and Ba, 2014), using recommended default parameters with a constant learning rate of 0.0003. Training was performed on a NVIDIA GTX Titan X Pascal graphical processing unit with acceleration provided by cuDNN libraries (Chetlur *et al.*, 2014). All reported performance metrics and saliency logos are drawn strictly from the test

set using the model parameters which yielded the lowest binary cross-entropy loss on the validation set, a technique known as early stopping.

Visualizing saliency analysis and 1st layer filters

Saliency analysis is performed by calculating the gradients of a neuron-of-interest with respect to the input one-hot representation. We use a variant of saliency analysis, called guided-backpropagation (Springenberg *et al*, 2014), which rectifies negative gradients through each ReLU activation. To generate a saliency logo, we calculated the saliency map using guided-backpropagation from the logits of a given class to the inputs. We then normalized the saliency map by dividing the maximum absolute value across the saliency map. Next, we applied an exponential filter according to: $\hat{S} = \exp\left[\lambda \frac{S}{\max|S|}\right]$, where \hat{S} is the normalized saliency map, S is the saliency map generated by guided-backprop, λ is a scaling factor that we set to 3 for all of saliency logos in this paper. We then separately normalized each position across by dividing the sum of the filtered saliency map across nucleotides, thereby providing a probability for each nucleotide at each position. To generate a sequence logo, each amino acid a at each nucleotide position i is scaled according to: $\hat{S}_{a,i} \times H_i$, where $H_i = 2 + \sum_a \hat{S}_{a,i} \log_2 \hat{S}_{a,i}$. First layer convolution filters were normalized and visualized following the same procedure, with the exception that the filter was used instead of the guided-backprop saliency map.

Availability

Python scripts to process the datasets and TensorFlow code to build, train, and evaluate the CNNs can be found via https://github.com/p-koo/learning_sequence_motifs.

Acknowledgements

The authors thank Tim Dunn and Soohyun Cho for helpful feedback on the manuscript.

References

- Alipanahi B, Delong A, Weirauch MT, Frey BJ (2015) Predicting the sequence specificities of DNA-and RNA-binding proteins by deep learning. *Nature Biotechnology* **33**: 831–838
- Angermueller C, Pärnamaa T, Parts L, Stegle O (2016) Deep learning for computational biology. *Molecular Systems Biology* **12**: 878
- Chetlur S, Woolley C, Vandermersch P, Cohen J, Tran J, Catanzaro B, Shelhamer E (2014) cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv* **1410.0759**
- Foat BC, Morozov AV, Bussemaker HJ (2006) Statistical mechanical modeling of genome-wide transcription factor occupancy data by MatrixREDUCE. *Bioinformatics* **22**: e141–e149
- Ghandi M, Mohammad-Noori M, Ghareghani N, Lee D, Garraway L, Beer MA (2016) gkmSVM: an R package for gapped-kmer SVM. *Bioinformatics* **32**: 2205–2207
- Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble WS (2007) Quantifying similarity between motifs. *Genome Biology* **8**
- He K, Zhang X, Ren S, Sun J (2015) Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision*
- Hinton G, McClelland J, Rumelhart D (1986) Distributed representations. *Parallel distributed processing Explorations in the microstructure of cognition* **1**: 77–109
- Hiranuma N, Lundberg S, Lee S (2017) DeepATAC: A deep-learning method to predict regulatory factor binding activity from ATAC-seq signals. *bioRxiv* **172767**

- Ioffe S, Szegedy C (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv* **1502.03167**
- Kelley DR, Snoek J, Rinn JL (2016) Bassett: Learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Research* **26**: 990–999
- Kingma D, Ba J (2014) Adam: A method for stochastic optimization. *arXiv preprint arXiv* **1412.6980**
- Mathelier A, Fornes O, Arenillas DJ, Chen Cy, Denay G, Lee J, Shi W, Shyr C, Tan G, Worsley-Hunt R, et al (2016) JASPAR 2016: a major expansion and update of the open-access database of transcription factor binding profiles. *Nucleic Acids Research* **44**: D110–D115
- Quang D, Xie X (2016) DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. *Nucleic Acids Research* **44**: 107
- Simonyan K, Vedaldi A, Zisserman A (2013) Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv* **1312.6034**
- Springenberg JT, Dosovitskiy A, Brox T, Riedmiller M (2014) Striving for simplicity: The all convolutional net. *arXiv preprint arXiv* **1412.6806**
- Srivastava N, Hinton GE, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* **15**: 1929–1958
- Zeng H, Edwards MD, Liu G, Gifford DK (2016) Convolutional neural network architectures for predicting DNA–protein binding. *Bioinformatics* **32**: i121–i127
- Zhou J, Troyanskaya OG (2015) Predicting effects of noncoding variants with deep learning-based sequence model. *Nature Methods* **12**: 931–934