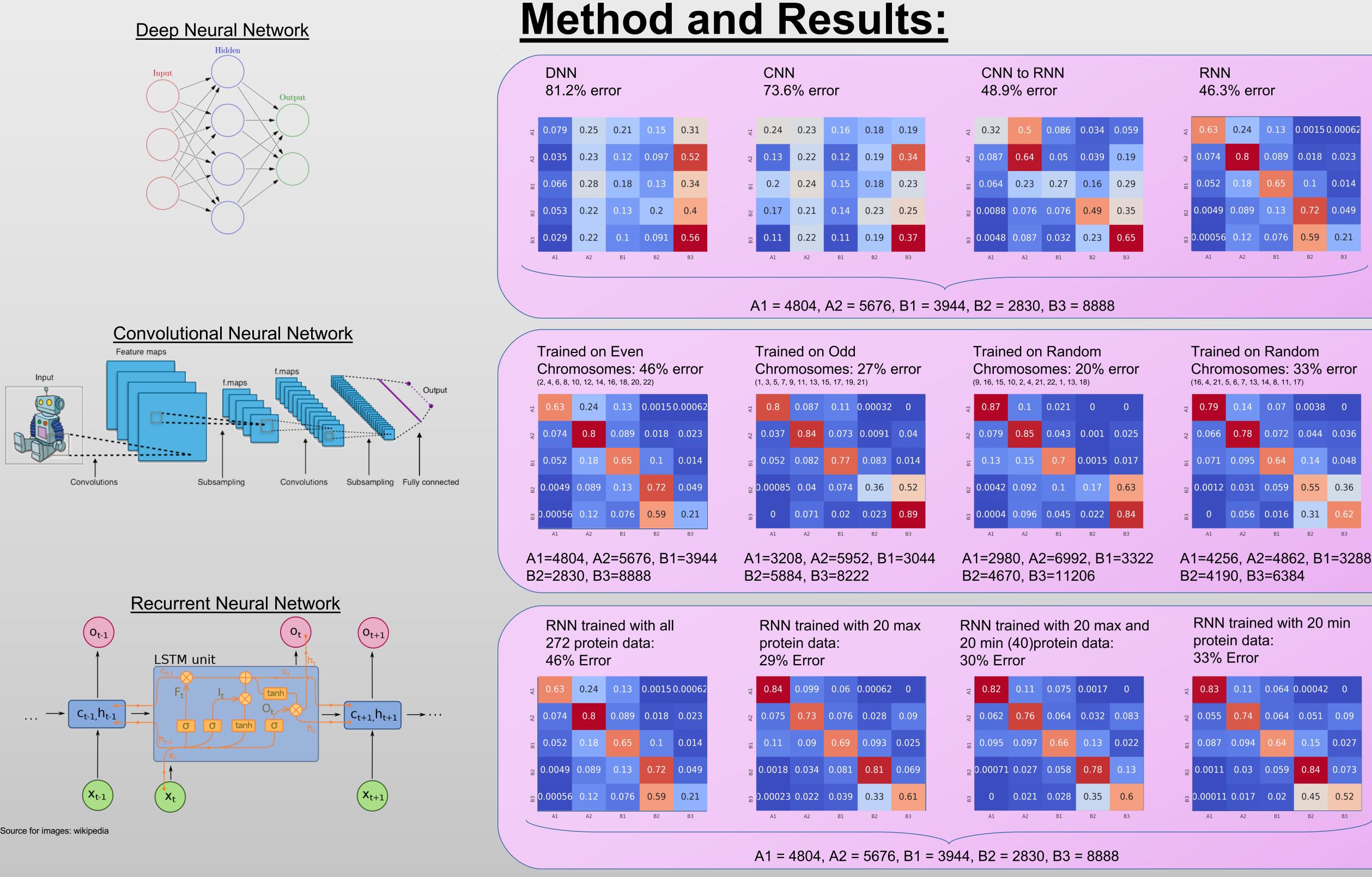
UNIVERSITY of HOUSTON

Abstract:

Modern molecular biology produces large amounts of data, which can be difficult to derive any useful information from. We are investigating correlations that exist between genetic annotations of human DNA and chromosome structural features. Chromatin Immuno-Precipitation Sequencing(ChIP-Seq) data tracks, made available through the ENCODE project, characterize the biochemical nature of chromosomal loci. Chromatin can be categorized into types that we call type A and type B which we further classify into chromatin sub-types (A1, A2, B1, B2, and B3). It has been previously shown that these chromatin structural types are strongly related to the overall genome architecture of cells.

Machine learning algorithms have proven to be especially adept at "learning" from correlations in very large data sets. We constructed a number of machine learning models and tested how accurately each performed when identifying chromatin sub-types. Our best approach so far is a recurrent neural network which produced a total error of less than 28% when classifying chromatin sub-types.



Bibliography:

Di Pierro, M. et al. De novo prediction of human chromosome structures: Epigenetic marking patterns encode genome architecture., 1714980114, (2017).

Acknowledgments:



NSF PHY-1427654 NSF ACI-1531814

Discussion:

The best model we have developed so far is a recurrent neural network made up of two bidirectional LSTM layers, the first having 272 memory units and the second having 350. The output layer is a Dense, time distributed layer with 5 nodes, correlating to the 5 chromatin sub-types, and a softmax activation function. Binary cross entropy is the loss function and Adam is used as the optimization function in this model. By reducing the protein data set, increasing the number of training epochs and decreasing the batch size, we were able to obtain a model that produces only 28% total error and correctly distinguishes between type A and type B with high accuracy.

Using machine learning to analyze large amounts of data in this manner has been shown to be very insightful in the fields of epigenetics and genomics. Applying this approach to similar problems will help us to learn more about chromosome structure and how it is related to gene expression.

Evaluating Machine Learning Approaches for Structural Genomics

Jonathan C. Pickett^{*ab}, Arya Haji Taheri^{ab}, Ryan R. Cheng^a, Michele Di Pierro^a Vinícius G. Contessoto^a, and José Onuchic^a

^aCenter for Theoretical Biological Physics – CTBP – BRC, Rice University, Houston, TX ^bUniversity of Houston - UH - Houston, TX **j.pickett.2004@gmail.com*



Model Comparison:

There are many types of machine learning model structures. Some of the most common are deep neural networks(DNN), convolutional neural networks(CNN) and recurrent neural networks(RNN). Each of these has a number of variants and limitless combinations that can be used to address unique problems.

Here, DNN, CNN, and RNN, as well as a model that combines aspects of all three, are compared through the use of confusion matrices. All of these models were trained with ChIP-Seq data from the even numbered chromosomes and tested with ChIP-Seq data from the odd numbered chromosomes. The RNN model outperforms any other model that was tested.

Model Consistency:

The ability to make accurate predictions on data that a model has not been previously exposed to is the ultimate goal of any machine learning model. We used different sets of ChIP-Seq data to train the RNN model, showing that it produces accurate results on testing data regardless of whether it is trained on the same data or different data. A curious observation from these results is the preferential prediction of B2 over B3 or vice-versa. The B3 sub-type is by far the most numerous, so when the model preferentially predicts B2 it is bound to have more errors. The models that preferentially predict B3 have far better accuracy.

Protein Data Optimization:

In our RNN there are weights matrices which contain values which strongly express the importance of specific protein experiments. By evaluating which of these weights are given the highest values we can determine which protein experimental values are more important to training our model. This evaluation not only helps us to optimize our model's predictive capabilities, it also gives us a clue about the important role that some of these proteins may play in determining chromosome structure in vivo. Reducing the size of the data set also seems to improve the quality of our results, which might imply that training on larger data sets is not yet optimal. Further investigation is needed to determine the actual meaning of these findings.

A1	0.82	0.11	0.066	0.0017	0		
A2	0.068	0.7	0.084	0.03	0.12		Maj
B1	0.1	0.078	0.67	0.11	0.04		<u> </u>
B2	0.0021	0.022	0.065	0.81	0.1	A	0.
B3	0	0.019	0.033	0.25	0.7	ß	0.
	A1	A2	B1	B2	B3		



