

DeepREG

Boxiang Liu

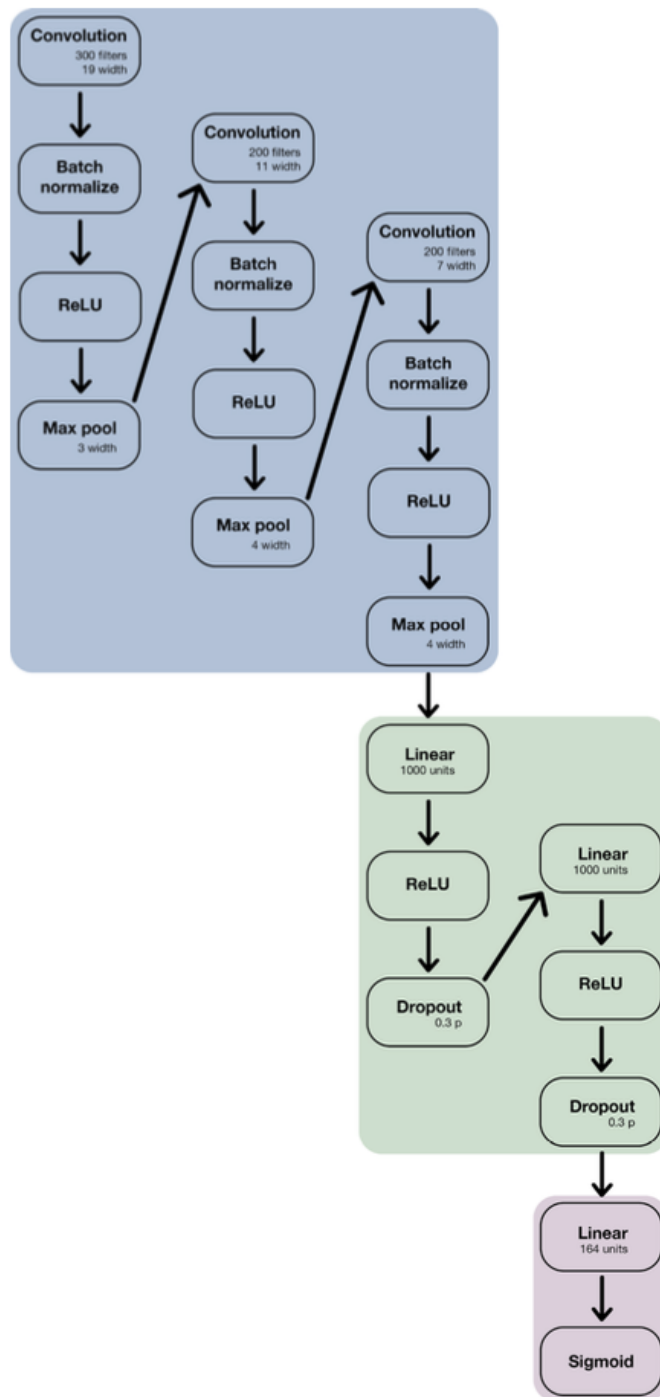
August 4, 2017

Contents

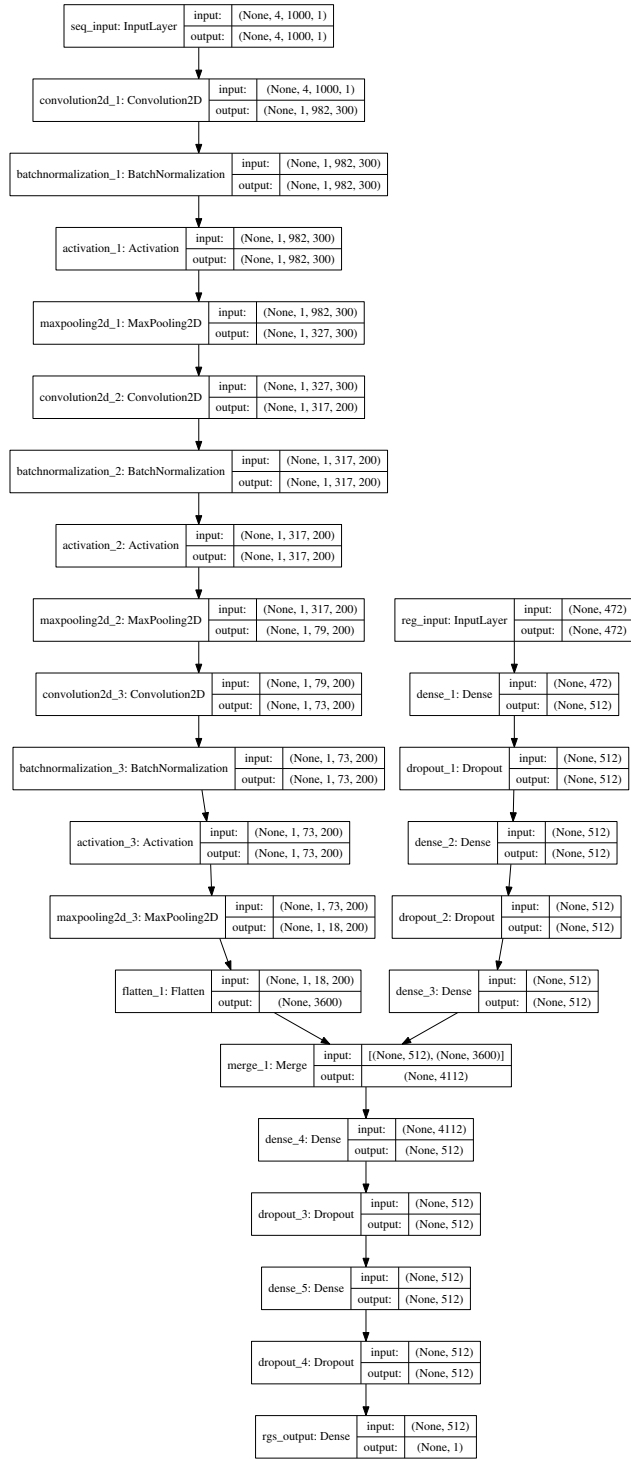
1	Basset	1
2	Reducing CNN layers	5
3	Small filter	6
4	Single layer	8
4.1	Motif discovery	12
4.2	Maxpooling after single layer	16
5	LSTM	17
6	GRU	19
7	Gene-Gene relationship	21
8	deepLIFT	21
9	Improving model accuracy	23

1 Basset

The Basset architecture represent the state-of-the-art for open chromatin predictions. The architecture is as follows:



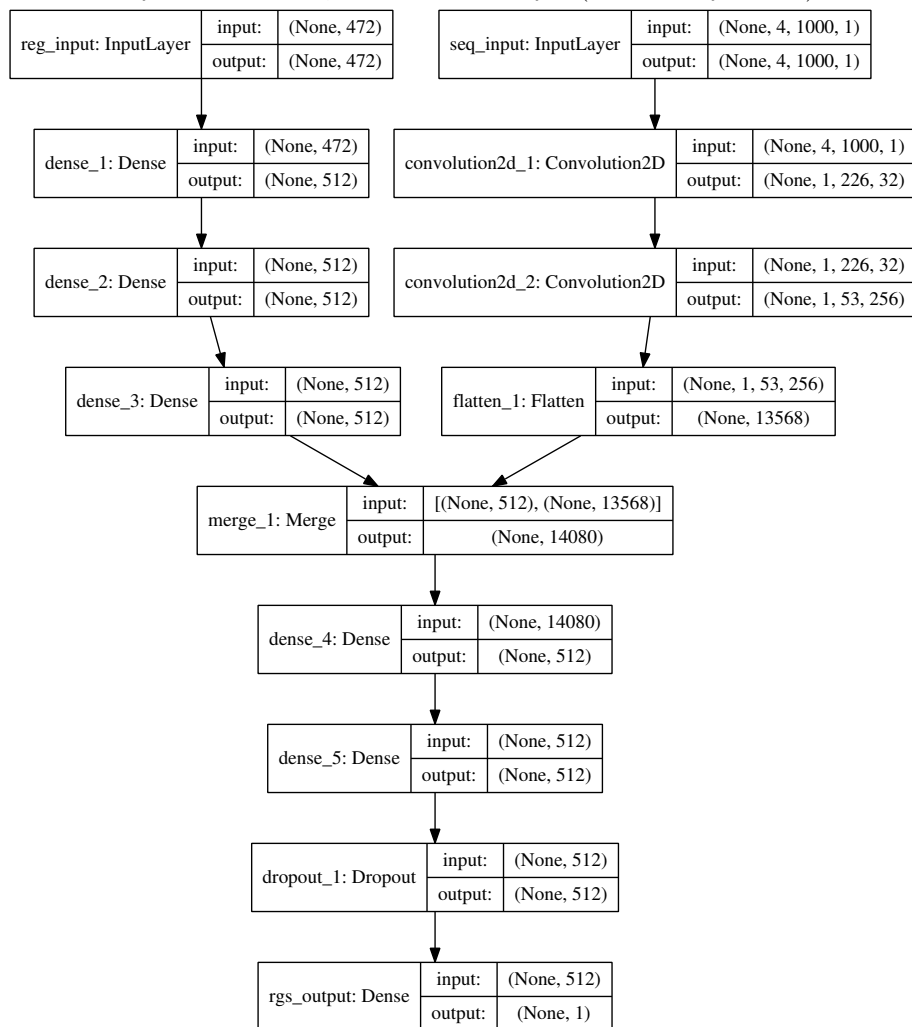
I used it for the CNN part of the network. The detailed network graph is below:



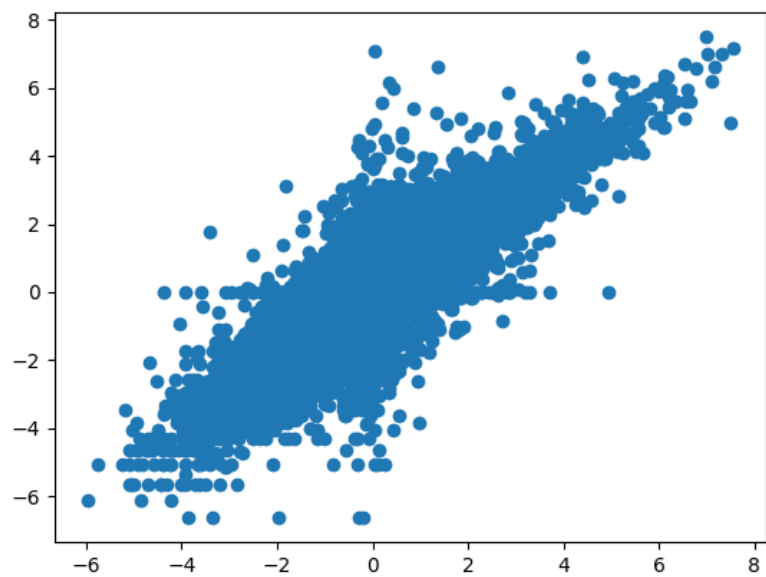
However, the network did not train properly, likely due to too many layers.

2 Reducing CNN layers

Given that 3 layers won't train, I removed one layer (in directory keras1).



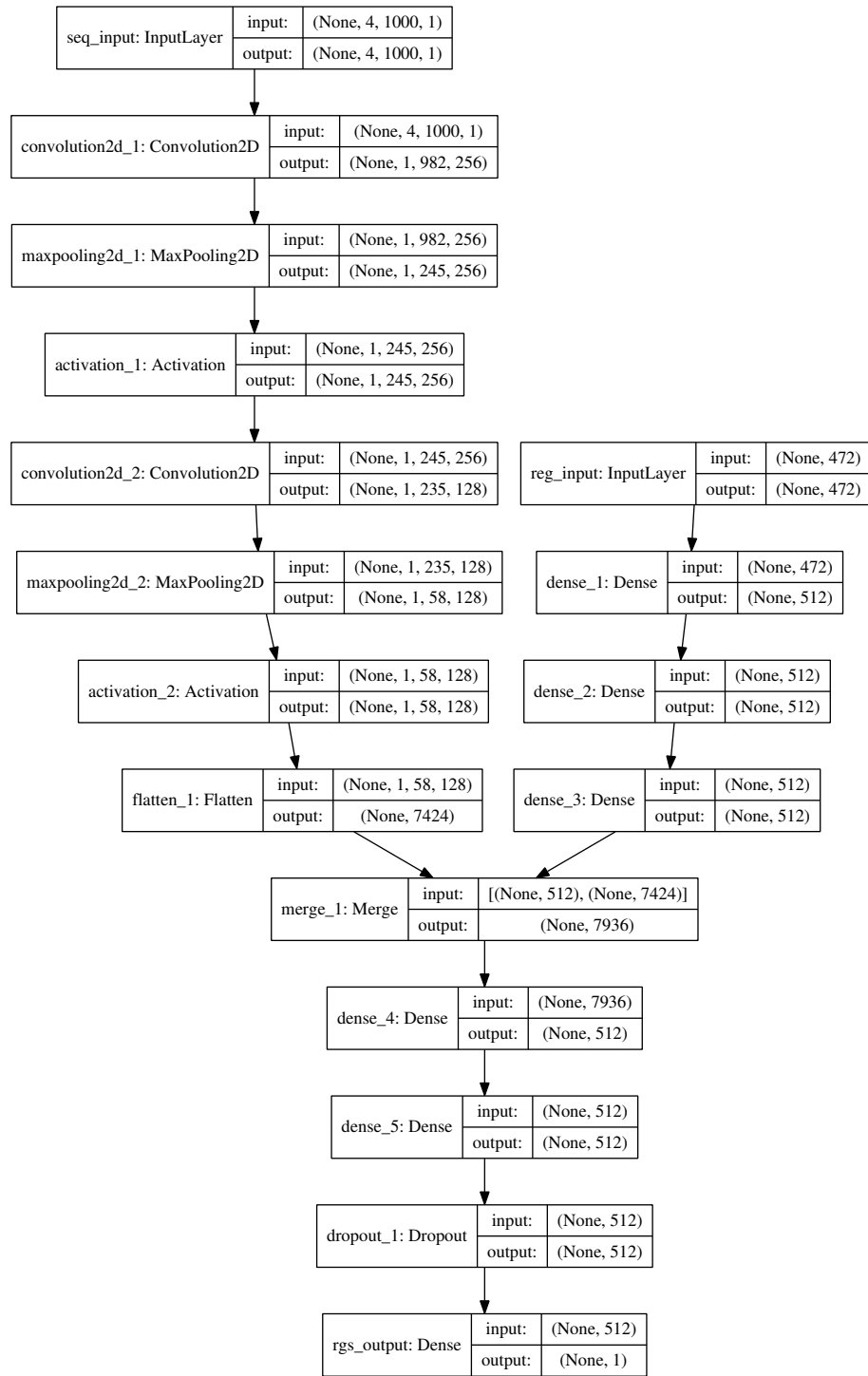
The result is quite promising.



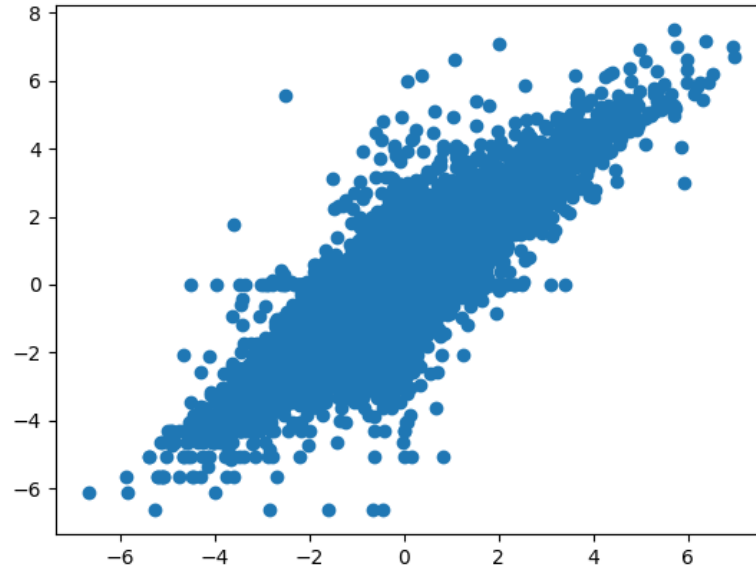
The first layer used filter width of 100, which is quite large.

3 Small filter

Since most motifs are less than 20 bps, I used a filter width of 19 bps.

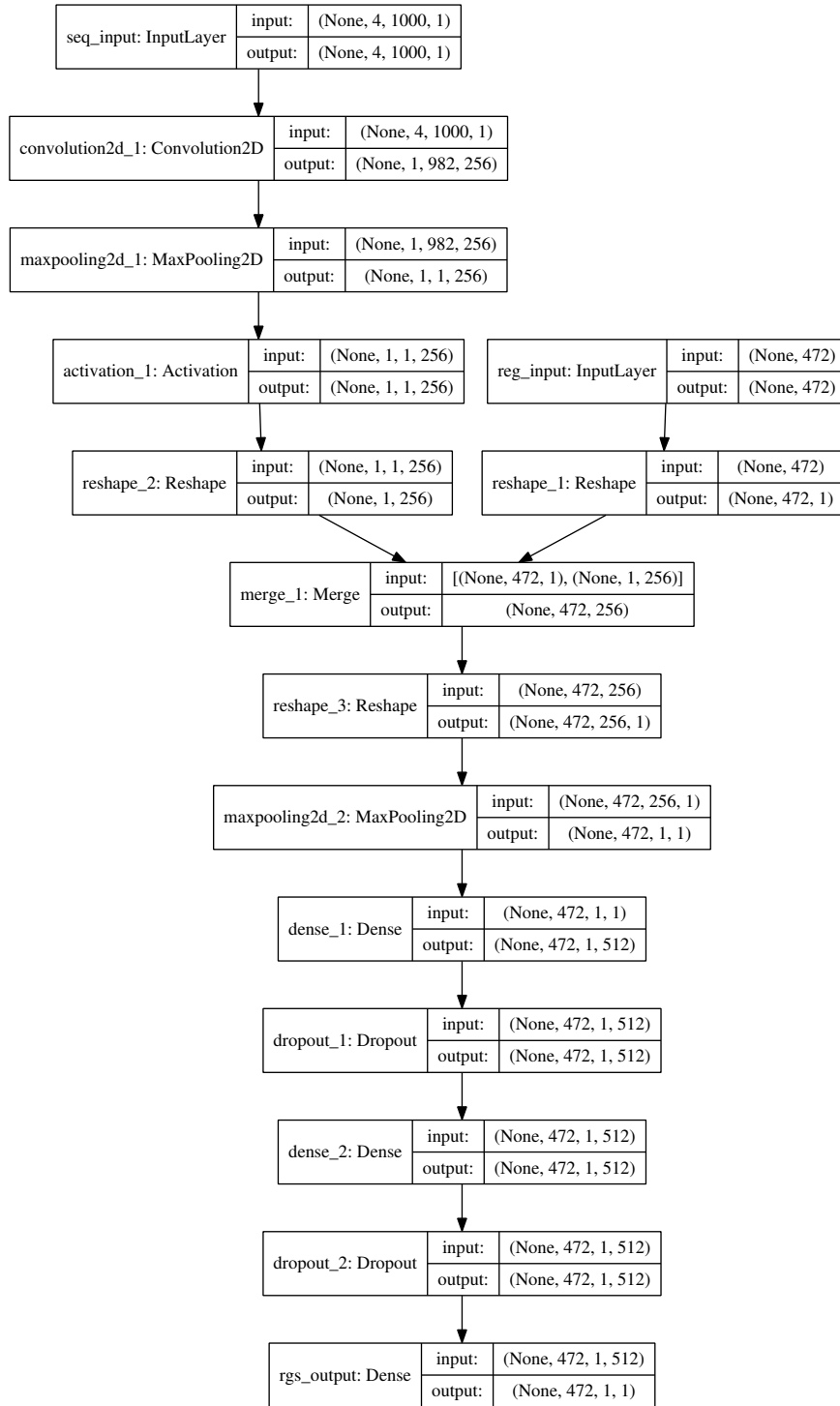


The result is as good as using 100 width filters.

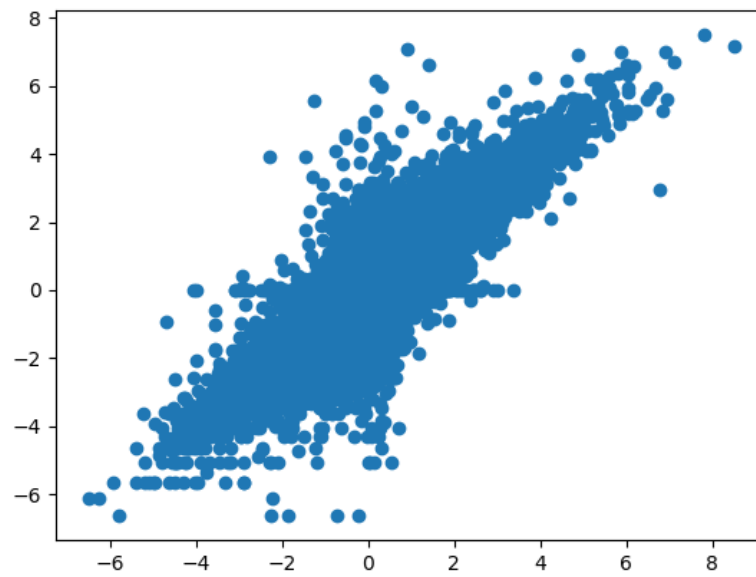


4 Single layer

When using more than one layer, either for the seq or the regulator network, the interpretability is lost. I therefore tried using one conv layer (as motif scanner) for the seq network, and no dense layer for the reg network.



The model worked really well. The training loss dropped to almost zero after 30 epochs, and does not show any sign of plateau (compared to)small filter). However the test loss does not decrease as much indicating overfitting.



```

Epoch 10/30
845208/845208 [=====] - 505s - loss: 0.1892 - val_loss: 0.2426
Epoch 11/30
845208/845208 [=====] - 501s - loss: 0.1679 - val_loss: 0.2351
Epoch 12/30
845208/845208 [=====] - 503s - loss: 0.1470 - val_loss: 0.2500
Epoch 13/30
845208/845208 [=====] - 499s - loss: 0.1279 - val_loss: 0.2880
Epoch 14/30
845208/845208 [=====] - 497s - loss: 0.1111 - val_loss: 0.2439
Epoch 15/30
845208/845208 [=====] - 498s - loss: 0.0966 - val_loss: 0.2307
Epoch 16/30
845208/845208 [=====] - 502s - loss: 0.0848 - val_loss: 0.2442
Epoch 17/30
845208/845208 [=====] - 500s - loss: 0.0749 - val_loss: 0.2221
Epoch 18/30
845208/845208 [=====] - 505s - loss: 0.0667 - val_loss: 0.2182
Epoch 19/30
845208/845208 [=====] - 504s - loss: 0.0600 - val_loss: 0.2152
Epoch 20/30
845208/845208 [=====] - 490s - loss: 0.0536 - val_loss: 0.2418
Epoch 21/30
845208/845208 [=====] - 499s - loss: 0.0491 - val_loss: 0.2178
Epoch 22/30
845208/845208 [=====] - 498s - loss: 0.0450 - val_loss: 0.2156
Epoch 23/30
845208/845208 [=====] - 497s - loss: 0.0417 - val_loss: 0.2163
Epoch 24/30
845208/845208 [=====] - 499s - loss: 0.0388 - val_loss: 0.2151
Epoch 25/30
845208/845208 [=====] - 497s - loss: 0.0361 - val_loss: 0.2129
Epoch 26/30
845208/845208 [=====] - 489s - loss: 0.0340 - val_loss: 0.2110
Epoch 27/30
845208/845208 [=====] - 501s - loss: 0.0321 - val_loss: 0.2430
Epoch 28/30
845208/845208 [=====] - 497s - loss: 0.0304 - val_loss: 0.2117
Epoch 29/30
845208/845208 [=====] - 495s - loss: 0.0287 - val_loss: 0.2128
Epoch 30/30
845208/845208 [=====] - 498s - loss: 0.0274 - val_loss: 0.2093

```

Therefore I created a new model with $l1 (=1e-7)$ and $l2 (=1e-7)$ on all weights regularization. Although this model prevented overfitting on the training set, the test set performance actually worsened.

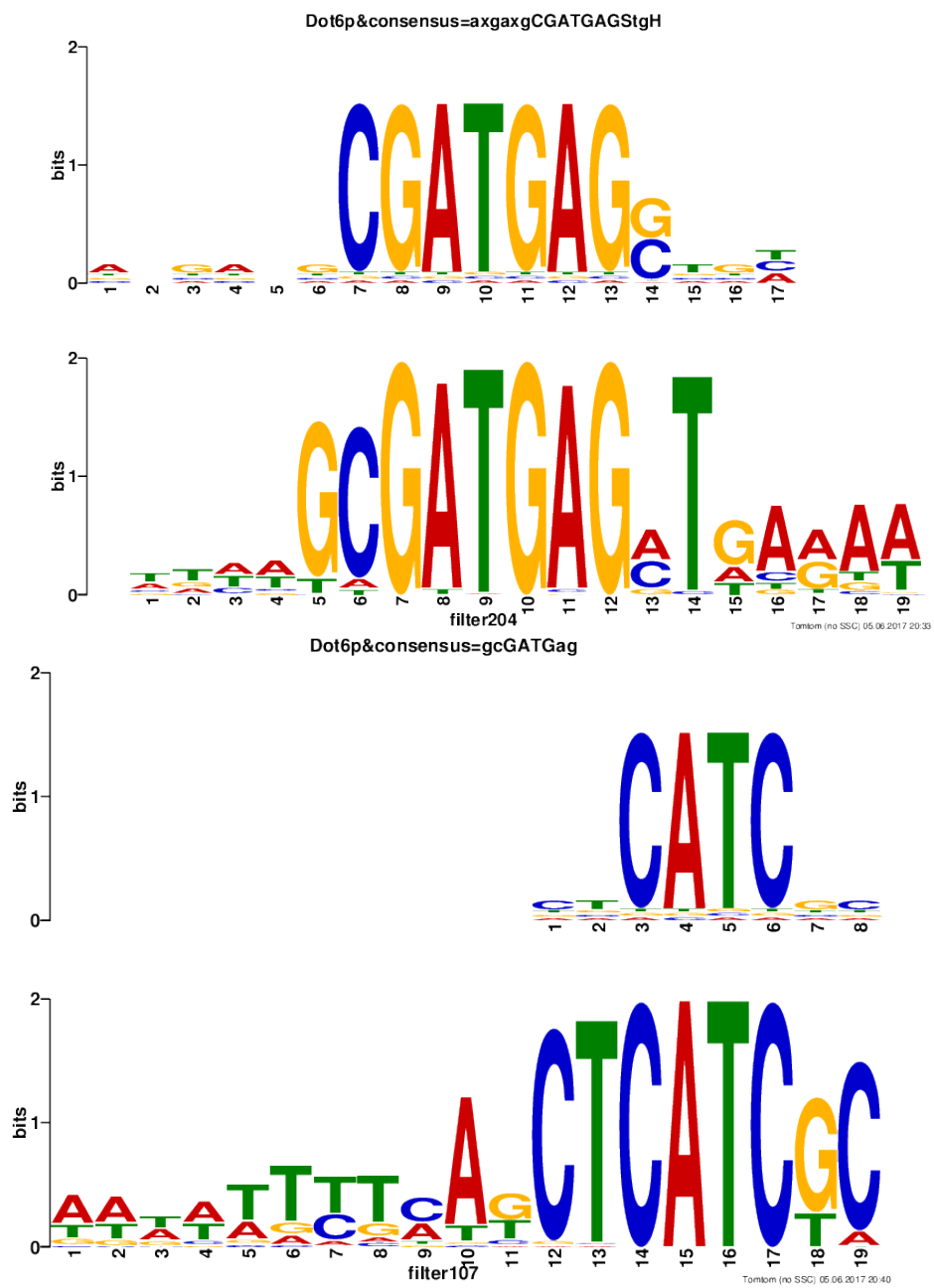
```

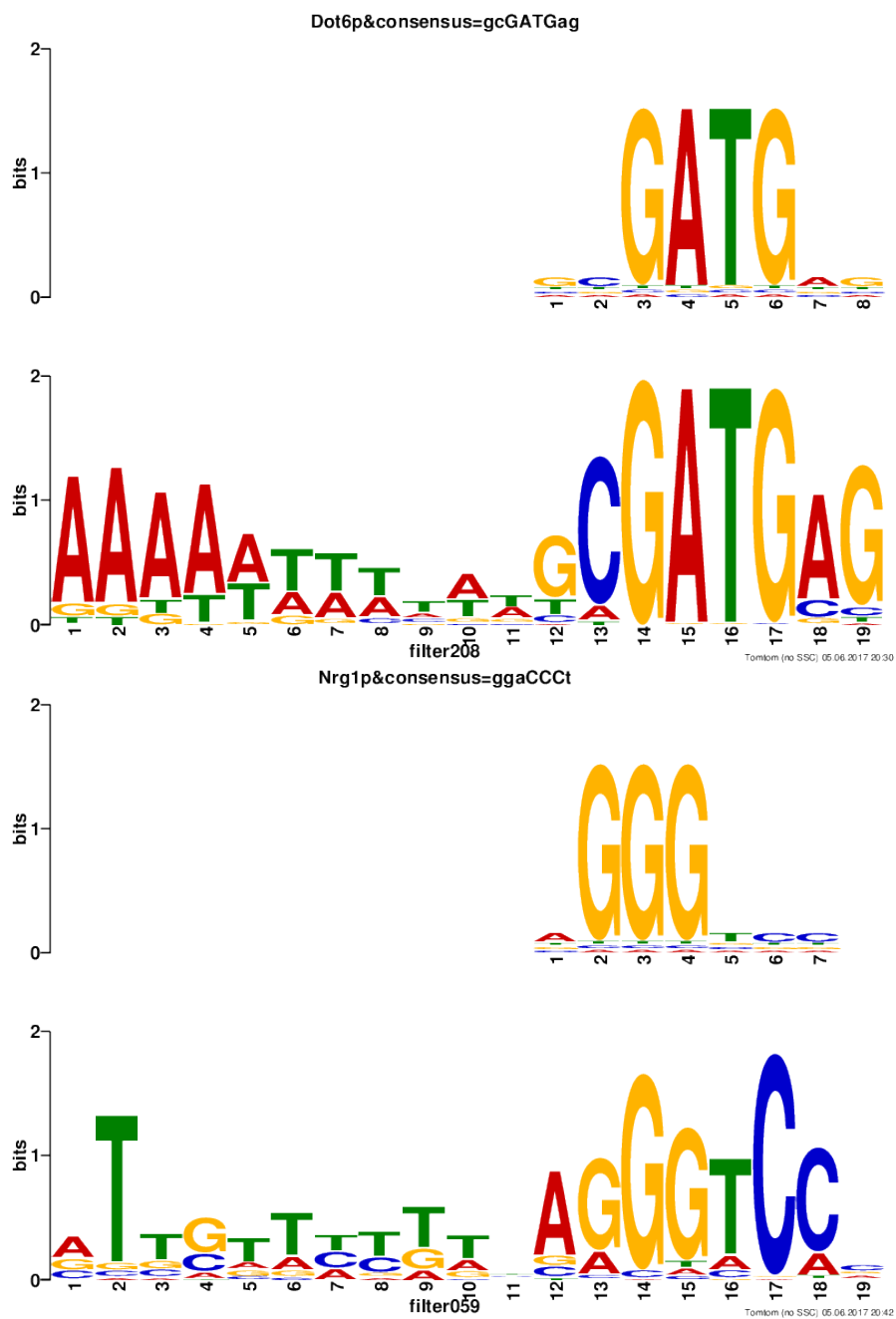
845208/845208 [=====] - 783s - loss: 0.1406 - val_loss: 0.2461
Epoch 28/60
845208/845208 [=====] - 785s - loss: 0.1367 - val_loss: 0.2525
Epoch 29/60
845208/845208 [=====] - 789s - loss: 0.1333 - val_loss: 0.2472
Epoch 30/60
845208/845208 [=====] - 790s - loss: 0.1301 - val_loss: 0.2494
Epoch 31/60
845208/845208 [=====] - 785s - loss: 0.1269 - val_loss: 0.2535
Epoch 32/60
845208/845208 [=====] - 782s - loss: 0.1242 - val_loss: 0.3130
Epoch 33/60
845208/845208 [=====] - 781s - loss: 0.1218 - val_loss: 0.2457
Epoch 34/60
845208/845208 [=====] - 787s - loss: 0.1188 - val_loss: 0.2487
Epoch 35/60
845208/845208 [=====] - 791s - loss: 0.1164 - val_loss: 0.2473
Epoch 36/60
845208/845208 [=====] - 788s - loss: 0.1143 - val_loss: 0.2463
Epoch 37/60
845208/845208 [=====] - 785s - loss: 0.1124 - val_loss: 0.2524
Epoch 38/60
845208/845208 [=====] - 790s - loss: 0.1103 - val_loss: 0.2496
Epoch 39/60
845208/845208 [=====] - 793s - loss: 0.1089 - val_loss: 0.2511
Epoch 40/60
845208/845208 [=====] - 789s - loss: 0.1068 - val_loss: 0.2535
Epoch 41/60
845208/845208 [=====] - 783s - loss: 0.1050 - val_loss: 0.2506
Epoch 42/60
845208/845208 [=====] - 790s - loss: 0.1040 - val_loss: 0.2482
Epoch 43/60
845208/845208 [=====] - 789s - loss: 0.1024 - val_loss: 0.2571
Epoch 44/60
845208/845208 [=====] - 789s - loss: 0.1009 - val_loss: 0.2522
105652/105652 [=====] - 14s

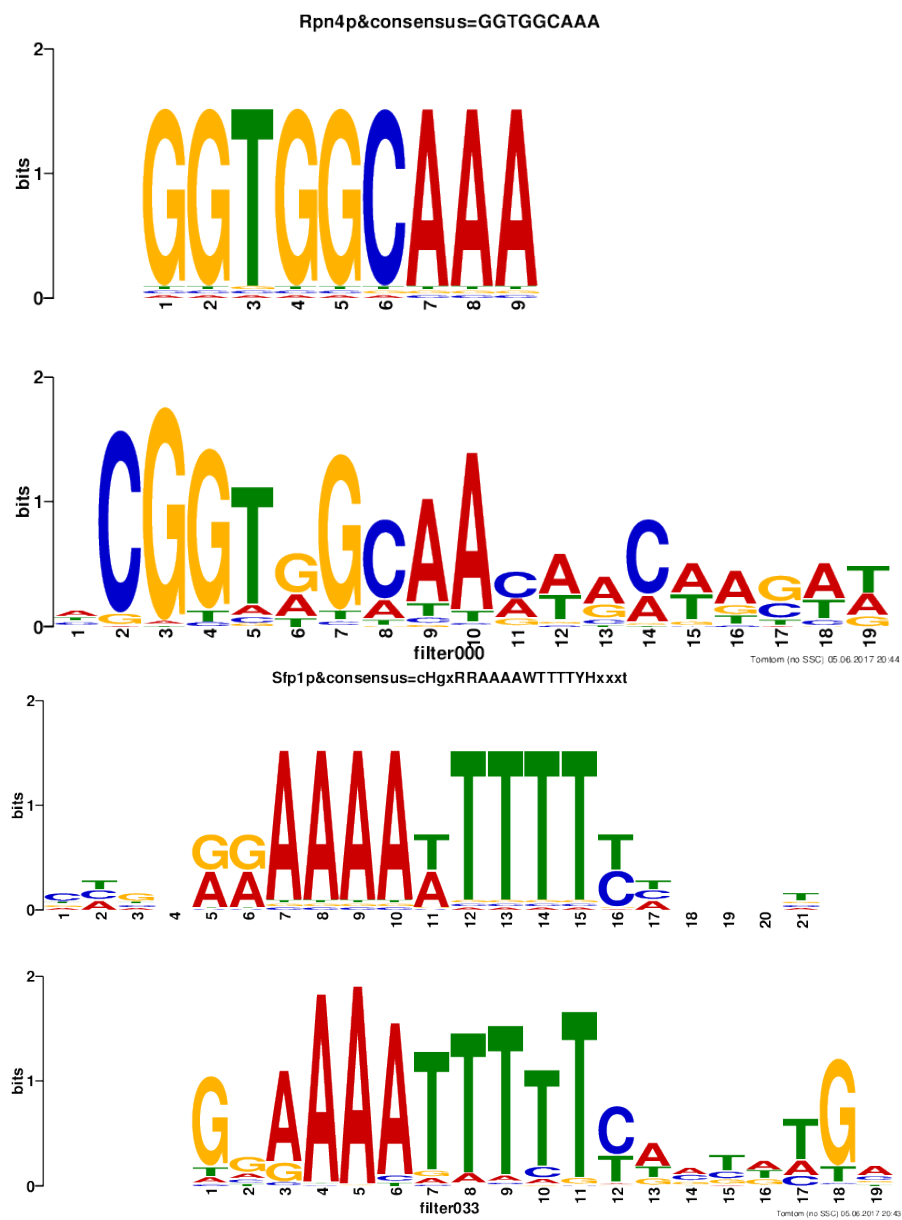
```

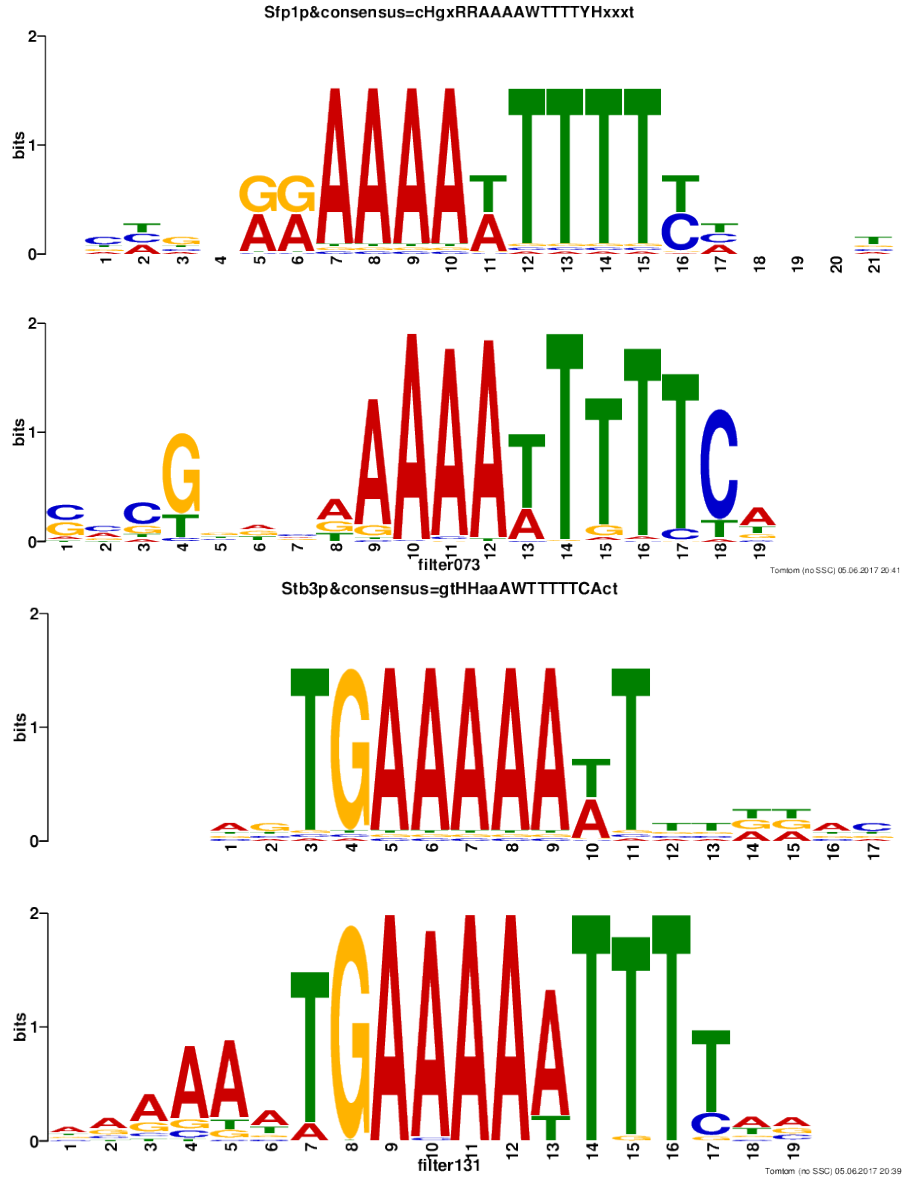
4.1 Motif discovery

Is the network find known motifs? I took top 100 sequences with largest activation for each filter and use TomTom to match them to known motifs.







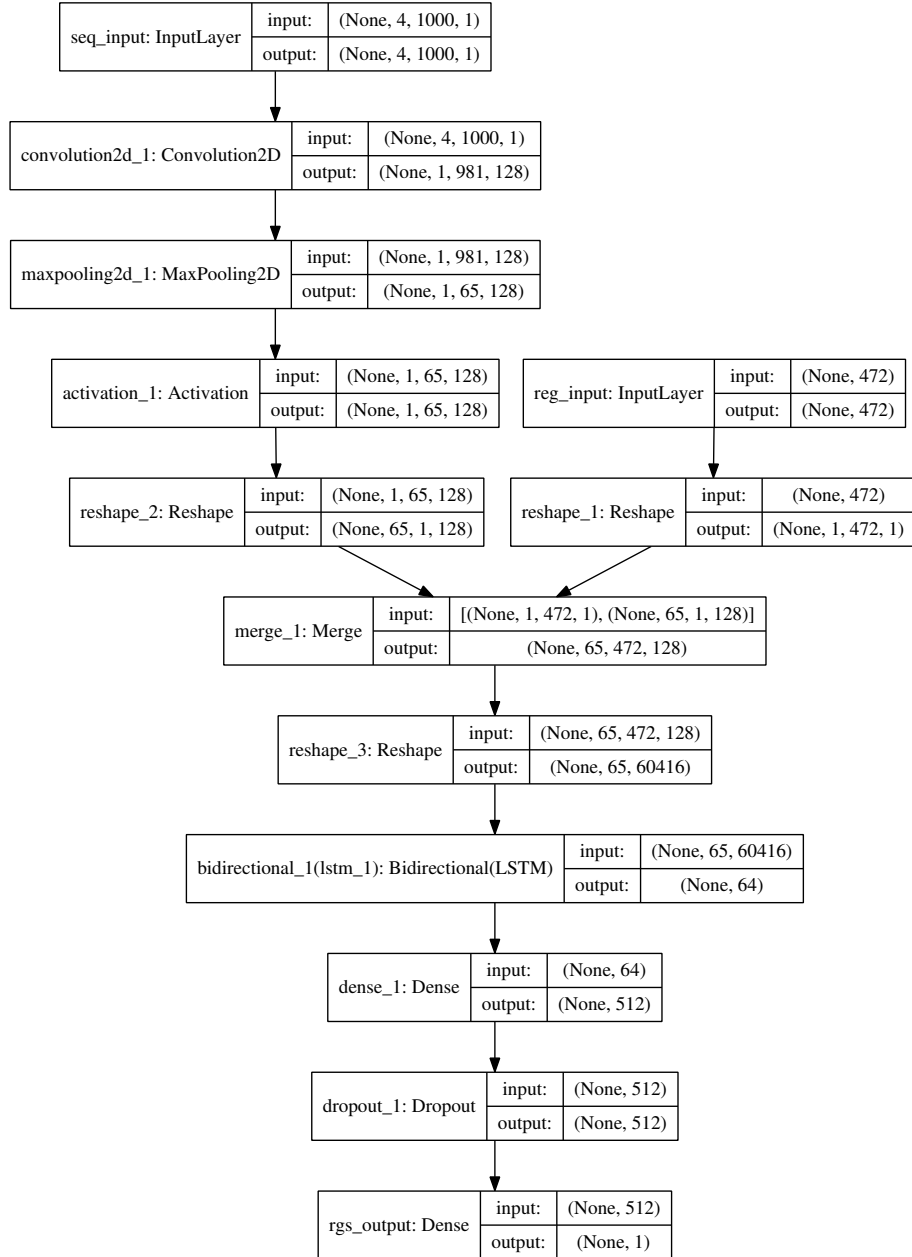


4.2 Maxpooling after single layer

The best validation loss is 0.4233 which is worse than just tensor product network.

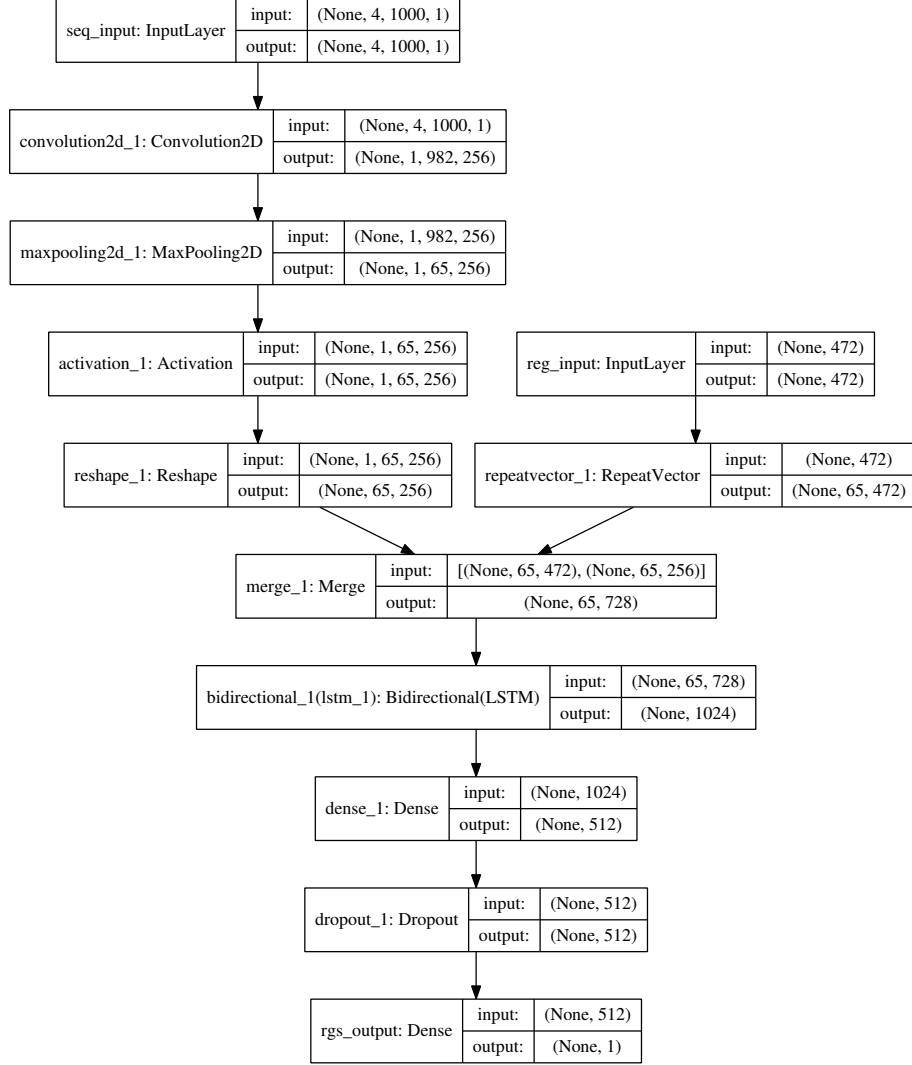
5 LSTM

Using outer product network ignores the spatial information along the genome. An LSTM with temporal dimension set to the genomic coordinate should be able to incorporate additional information in theory. The network is as follows:



This initial model uses an outer product as the input dimension and the

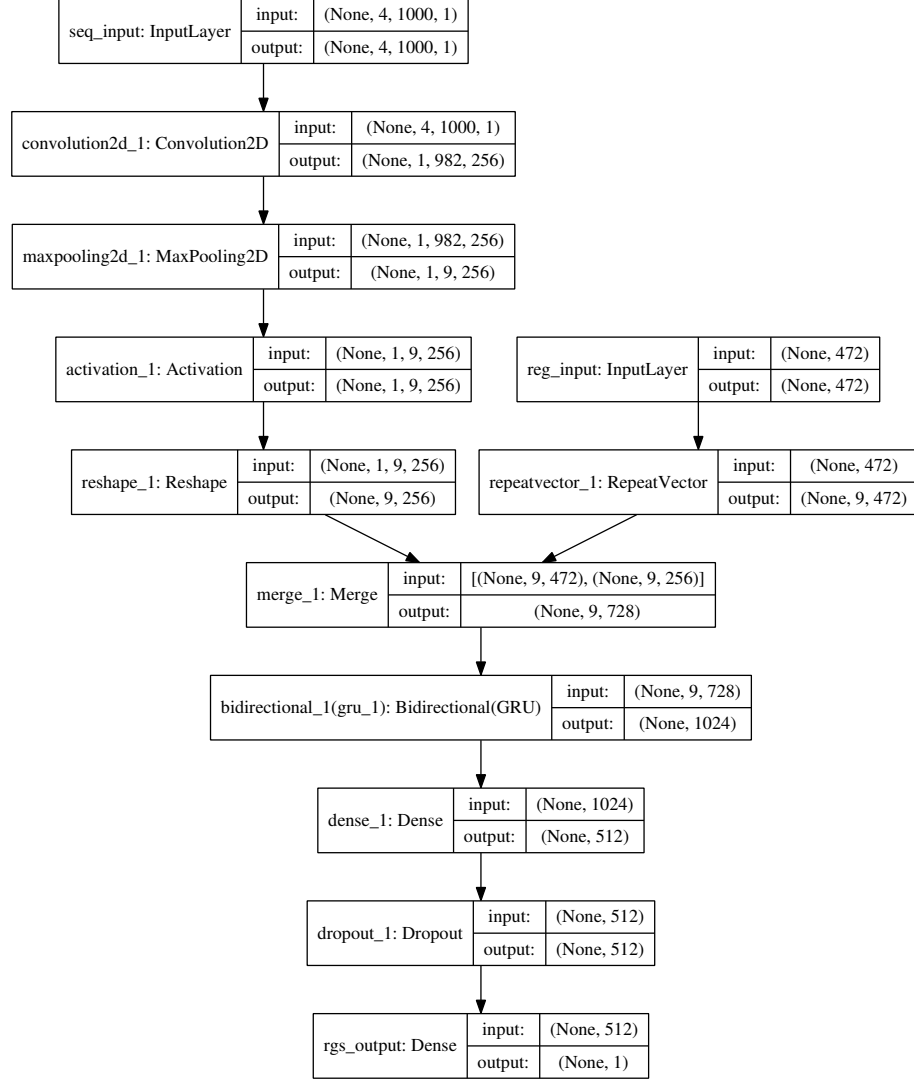
genomic coordinate as the time dimension. The size of input dimension is $256 \times 472 = 120832$, too large to fit into memory of GeForce GTX 970. A simpler variant replaces the outer product dimension with concatenation, effectively reducing the memory requirement by two orders of magnitude.

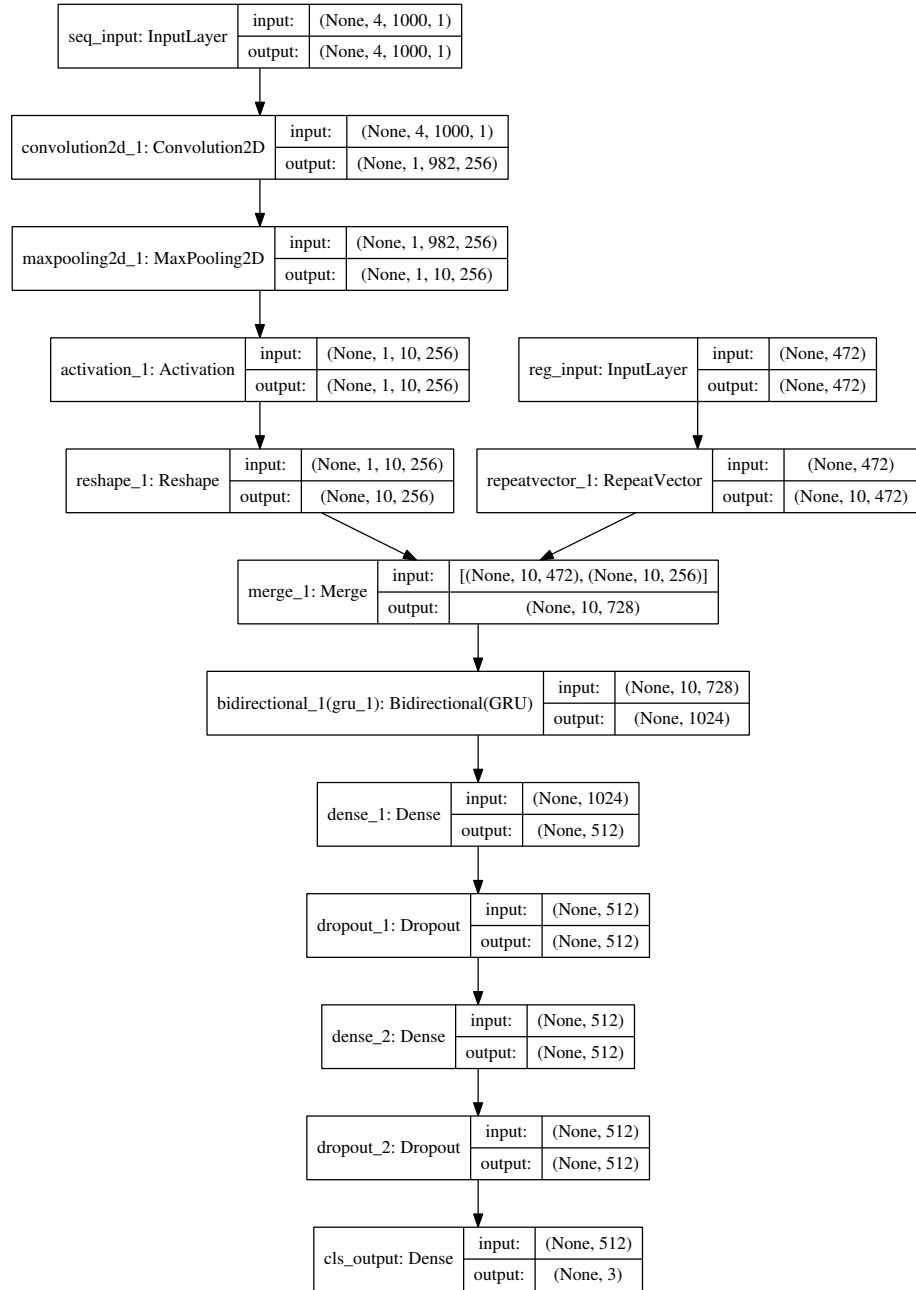


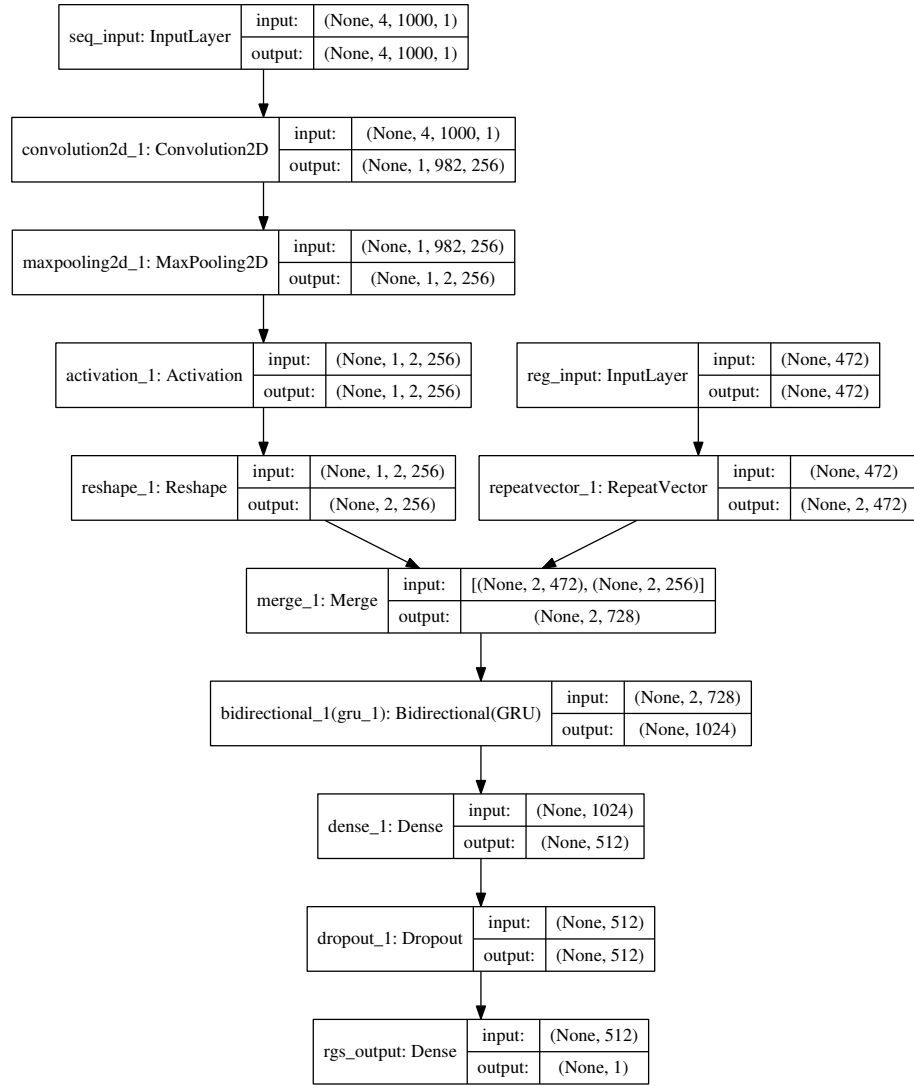
The concatenation model was trained with both SGD and ADAM. ADAM worked better than SGD. The best validation loss for LSTM is 0.3799, worse than 0.21 for tensor product network.

6 GRU

Since the GRU uses two gates, one gate less than LSTM, it is believed to be more computational efficient. I replaced the LSTM with GRU, and used maxpool = 15, 100, 491. In theory, when we increase the subsample ratio to 982, the GRU should become equivalent to the concatenation network.







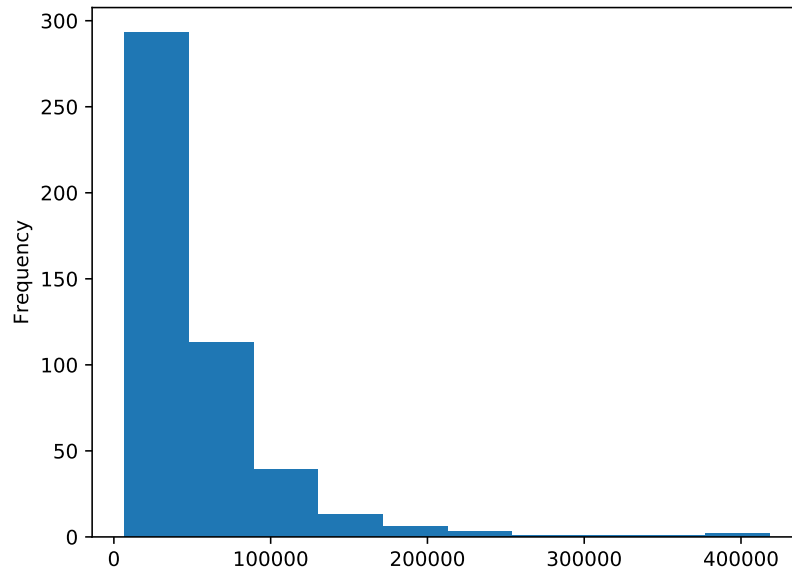
Performance-wise, the best validation loss for GRU is 0.3556, worse than 0.21 for concatenation network.

7 Gene-Gene relationship

8 deepLIFT

I rerun the **concatenation** network with *valid* padding (see *modeling/concatnation*) and calculated deeplift scores with respect to the regulator layer. Each of the 472 regulators receives 1056511 scores, one for each [gene,conditions] pair (there are

173 conditions \times 6107 genes). I used the sums of absolute values as the overall importance score for each regulator.



The distribution is highly skewed to the right, indicating that several important (potentially master regulators) dominates the model weights.

The top 10 regulators are

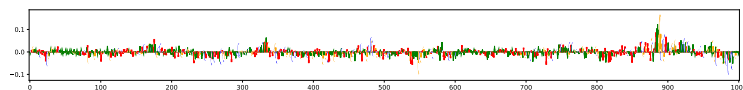
1. GAC1 / YOR178C
2. RAS1 / YOR101W
3. USV1 / YPL230W
4. MSN2 / YMR037C
5. PDR3 / YBL005W
6. YVH1 / YIR026C
7. PRR2 / YDL214C
8. SLT2 / YHR030C
9. BAS1 / YKR099W
10. SIP2 / YGL208W

Within these, SIP2 / YGL208W, SLT2/YHR030C, USV1 / YPL230W, GAC1 / YOR178C, MSN2 / YMR037C are also noted in Kundaje (2006) as top parents.

Since MSN2 is a known master regulator with over a hundred known targets, we tested whether MSN2 have higher deepLIFT scores for known targets. We downloaded a total of 381 genes from (<http://www.yeastgenome.org/locus/S000004640/overview>). Indeed, the known target has higher deepLIFT scores other genes (50.2 vs 49.9, $p=0.0016$).

I also tried running deepLIFT on the regression model. When running w.r.t the regulator layer with background nucleotide frequency as the reference, the following are the top 10 genes with highest cumulative deepLIFT scores: YDR277C 6.0 MTH1 / YDR277C YGL096W 2.0 TOS8 / YGL096W YGL099W 7.0 LSG1 / YGL099W YGR123C 3.0 PPT1 / YGR123C YHR136C 10.0 SPL2 / YHR136C YIR026C 4.0 YVH1 / YIR026C YKL109W 1.0 HAP4 / YKL109W YLR452C 9.0 SST2 / YLR452C YOR101W 8.0 RAS1 / YOR101W YPL230W 5.0 USV1 / YPL230W Of these, MTH1, PPT1, YVH1, HAP4, USV1 are also in Kundaje 2006. For MSN2, I compared its score between known targets and non-targets, the p-value is 0.00046!

I tried calculating deepLIFT score for the seq layer using genomics default and gradient*input. This did not work very well. Below is an example genomics default:

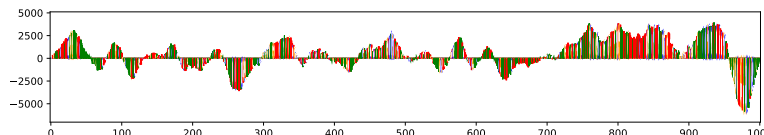


grad time input:

Later, I tried summing over the gradient*input over all 173 experiments.

□

This seems to have worked. We can observe patches of + and - scores. However, when I summed over all experiments and genes.



It seems the negative and positive patches are systematic!

9 Improving model accuracy

In the deep learning book, Ian give the following guidelines

First, performance metric is chosen along with the desired values. In the Street View project, he chose coverage $\hat{c}=95\%$.

Second, a baseline model is established. For his project, a convolutional neural network with ReLU were used.

Third, the baseline model is iteratively refined. One test whether each change makes an improvement.

Fourth, compare the train and test error. If the train and test error are similar, the model is underfitting or something is wrong with the training data. Consider using a more expressive model. If that does not work, look at the worst errors because something might be wrong with the training data. If the training error is lower than the test error, that indicates overfitting and regularization can be added.