

DeepREG

Boxiang Liu

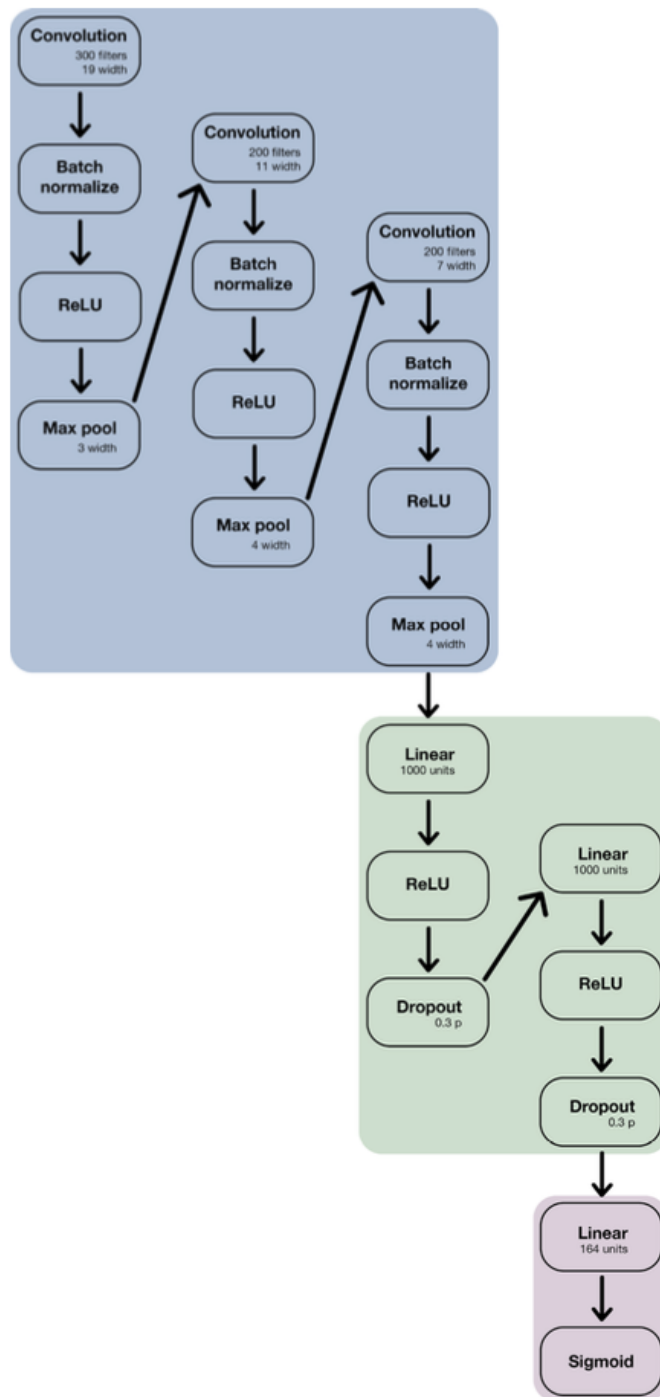
June 5, 2017

Contents

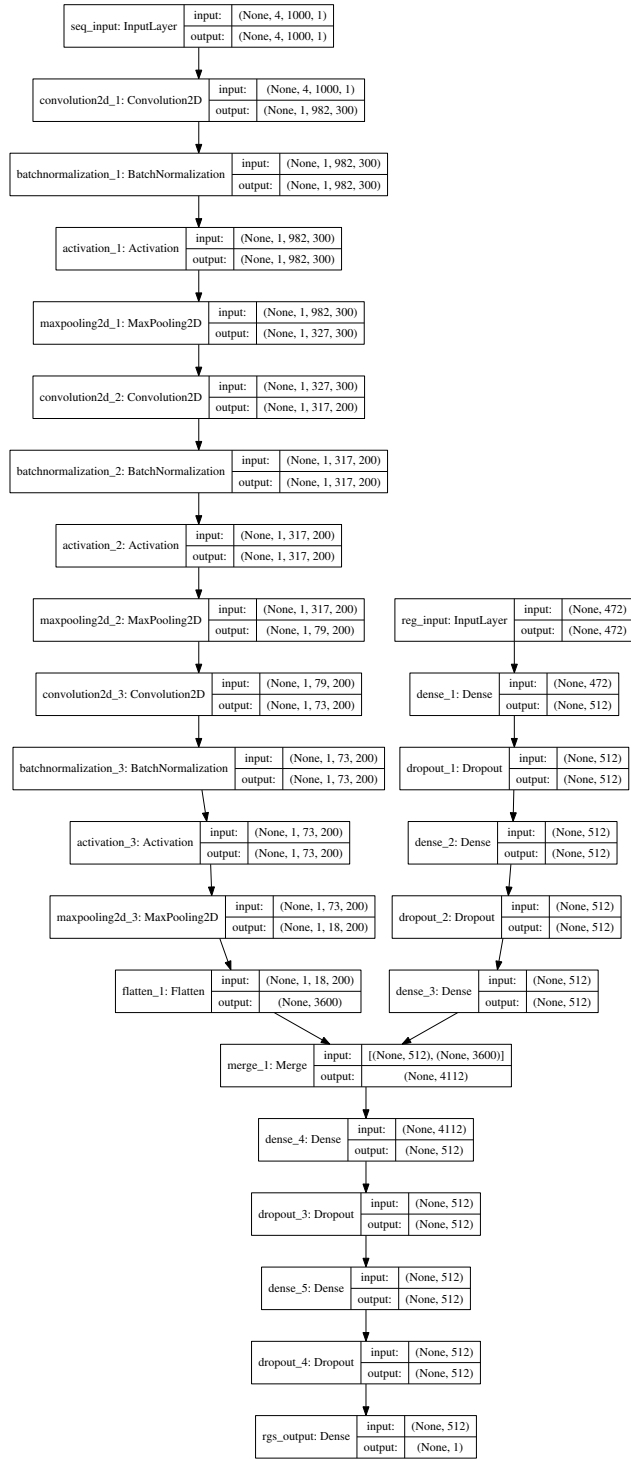
1	Basset	1
2	Reducing CNN layers	5
3	Small filter	6
4	Single layer	8
4.1	Motif discovery	11

1 Basset

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



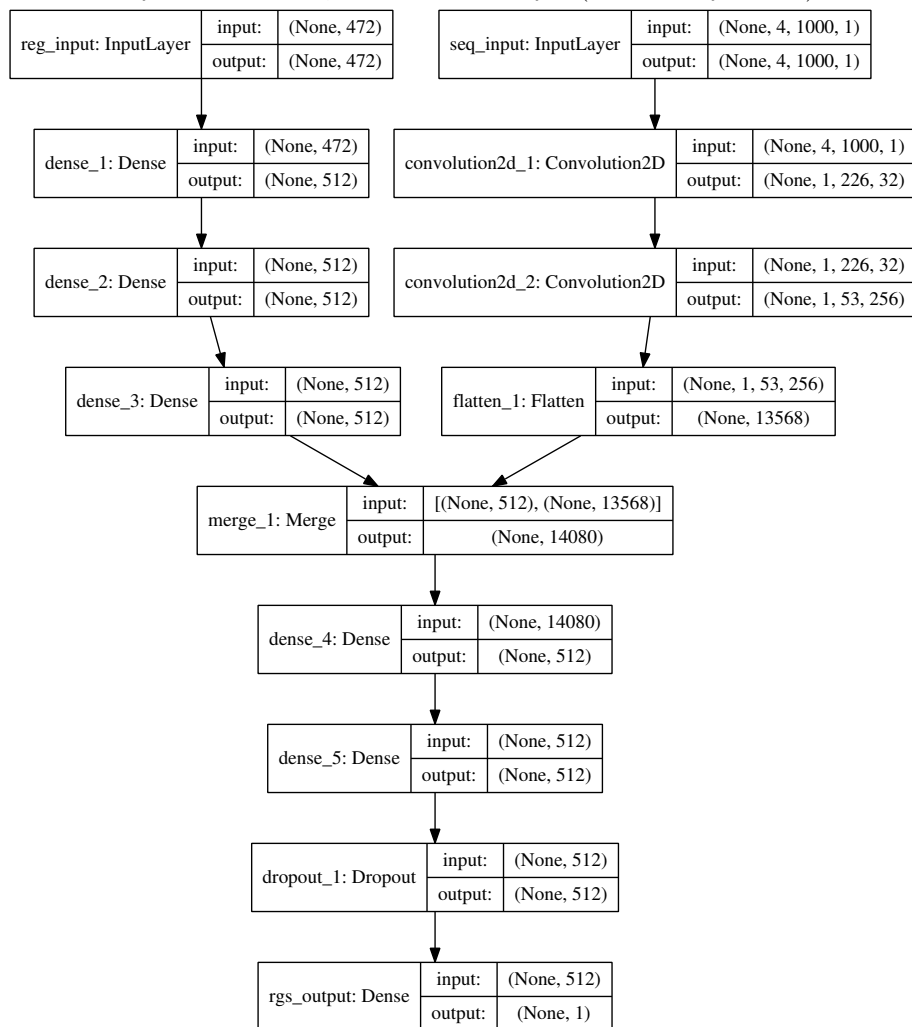
I used it for the CNN part of the network. The detailed network graph is below:



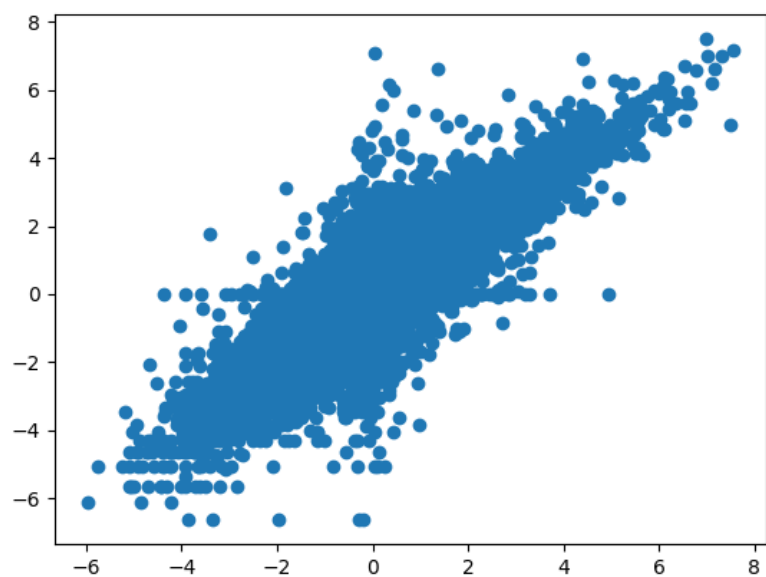
However, the network did not train properly, likely due to too many layers.

2 Reducing CNN layers

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



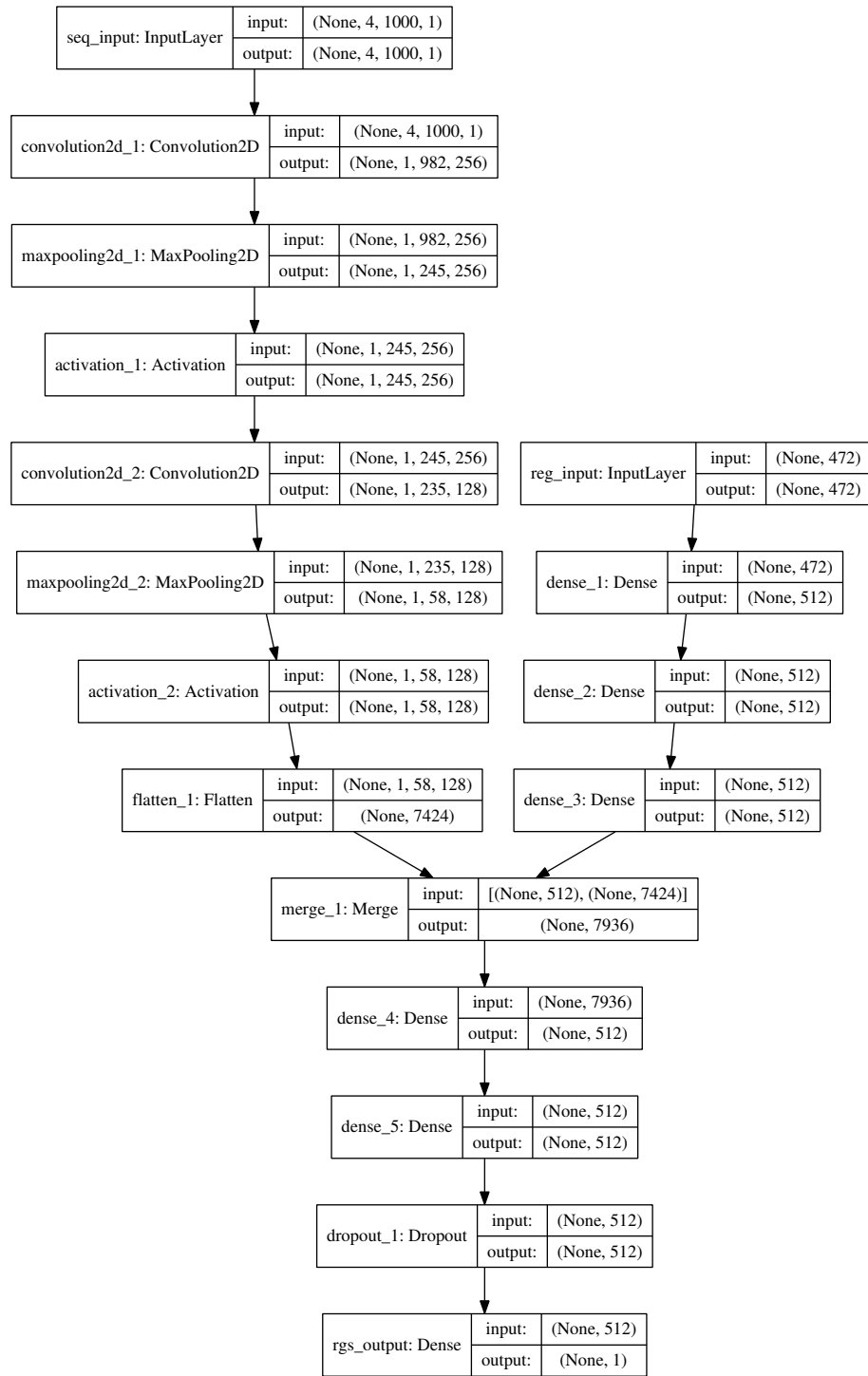
The result is quite promising.



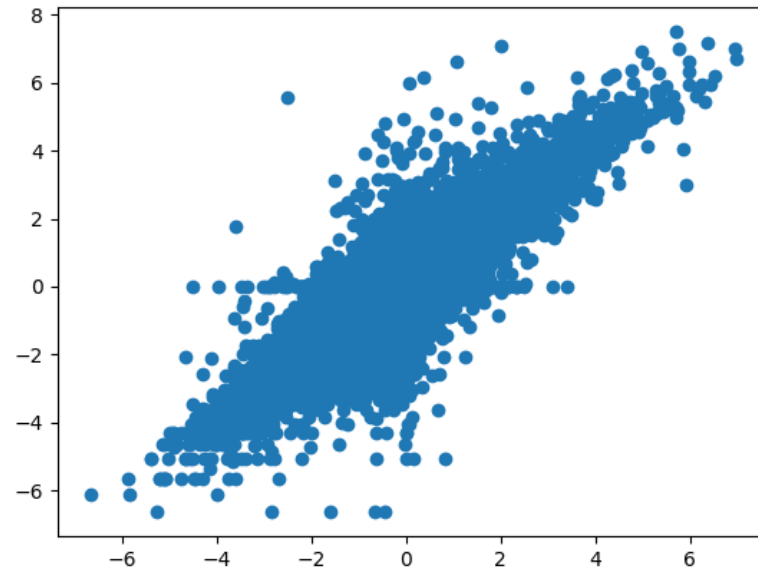
The first layer used filter width of 100, which is quite large.

3 Small filter

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

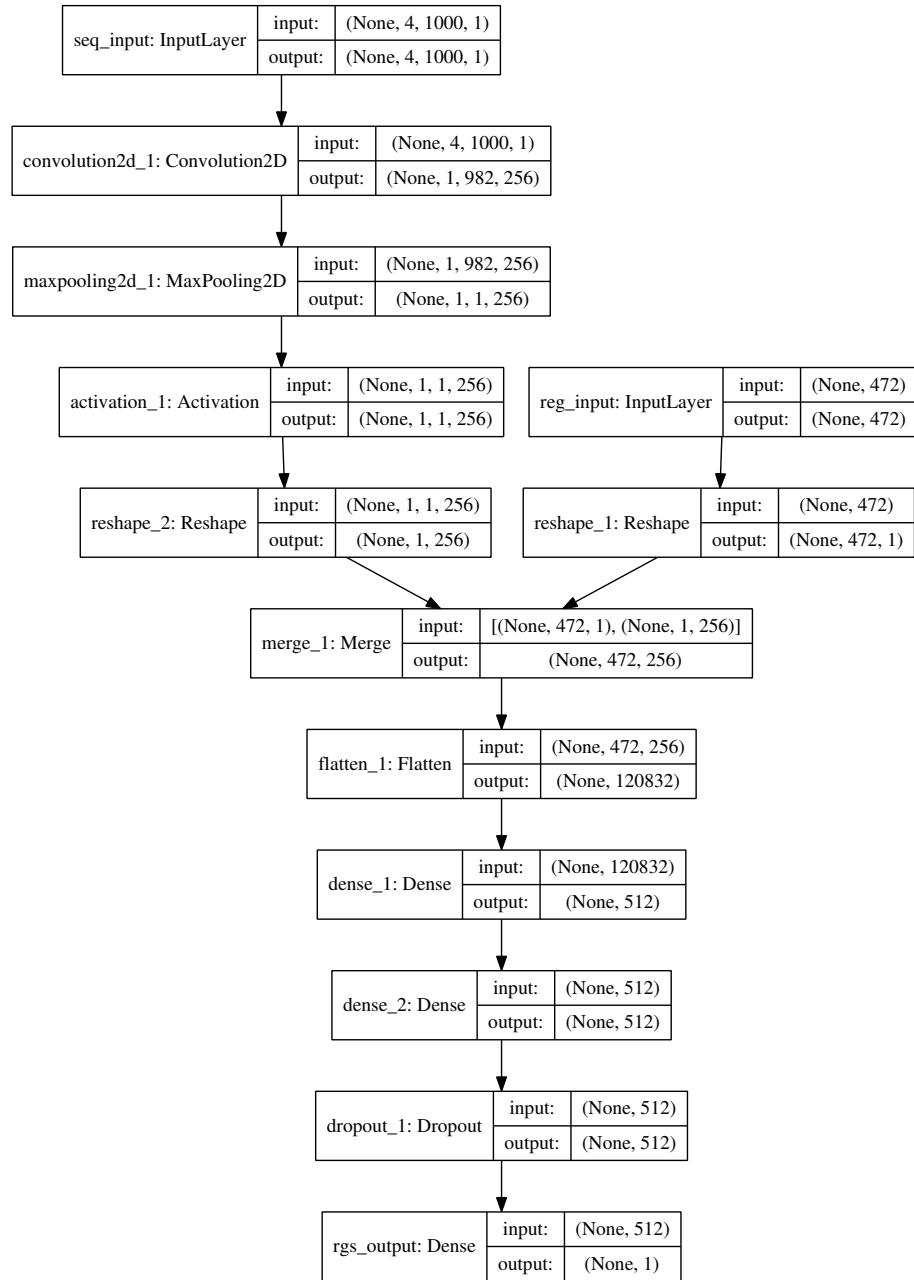


The result is as good as using 100 width filters.

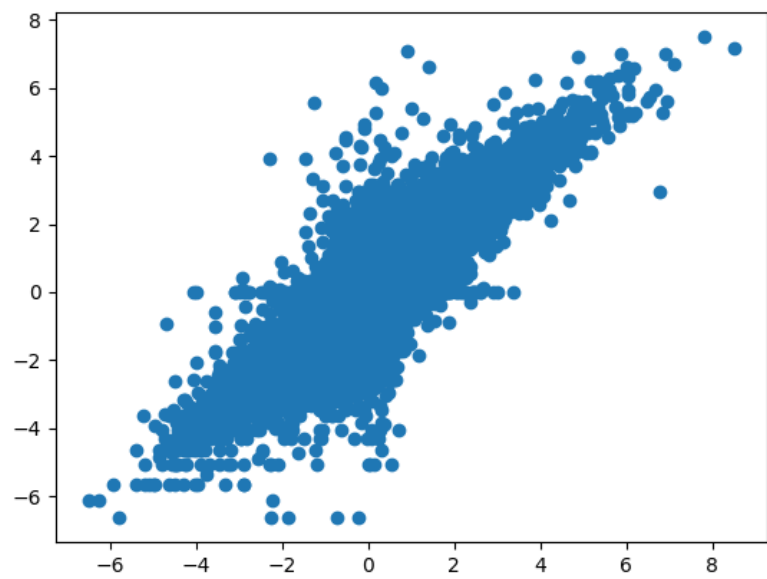


4 Single layer

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.

