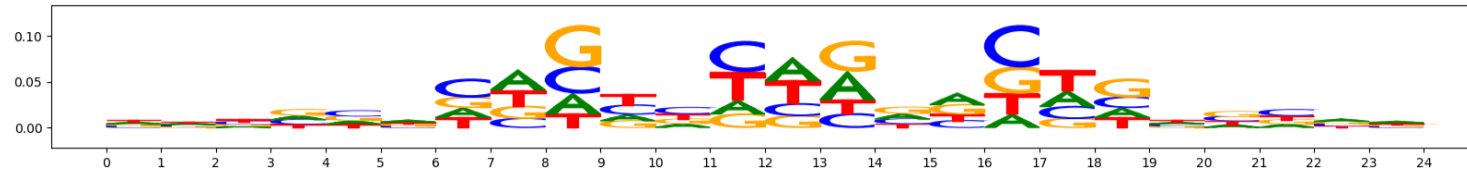


## Bias model training and quality check report

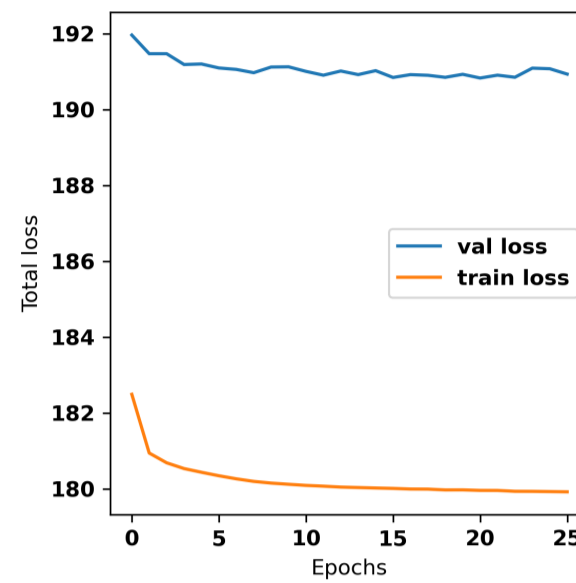
### Preprocessing report

The image below should look closely like a Tn5 or DNase bias enzyme motif.



### Training report

The val loss (validation loss) will decrease and saturate after a few epochs.



### Bias model performance in peaks and non-peaks

**Counts Metrics:** The pearsonr in non-peaks should be greater than 0 (higher the better). The pearsonr in peaks should be greater than -0.3 (otherwise the bias model could potentially be capturing AT bias). MSE (Mean Squared Error) will be high in peaks.

**Profile Metrics** Median JSD (Jensen Shannon Divergence between observed and predicted) lower the better. Median norm JSD is median of the min-max normalized JSD where min JSD is the worst case JSD i.e JSD of observed with uniform profile and max JSD is the best case JSD i.e 0. Median norm JSD is higher the better. Both JSD and median norm JSD are sensitive to read-depth. Higher read-depth results in better metrics.

**What to do if your pearsonr in peaks is less than -0.3?** In the range of -0.3 to -0.5 please be wary of your chrombpnet\_wo\_bias.h5 (that wil potentially be trained with this bias model) TFModisco showing lots of GC rich motifs (> 3 in the top-10). If this is not the case you can continue using the chrombpnet\_wo\_bias.h5. If you end up seeing a lot of GC rich motifs it is likely that bias model has learnt a different GC distribution than your GC-content in peaks. You might benefit from increasing the bias\_threshold\_factor argument input to the *chrombpnet bias pipeline* or *chrombpnet bias train* command used in training the bias model and retrain a new bias model. For more intuition about this argument refer to the [FAQ](#) section in wiki. If the value is less than -0.5 the [chrombpnet training](#) will automatically throw an error.

	nonpeaks.pearsonr	nonpeaks.mse	peaks.pearsonr	peaks.mse
counts metrics	0.28	0.38	-0.51	5.35
	nonpeaks.median jsd	nonpeaks.median norm jsd	peaks.median jsd	peaks.median norm jsd
profile metrics	0.03	0.15	0.49	0.20

### TFModisco motifs learnt from bias model (bias.h5) model

**TFModisco motifs generated from profile contribution scores of the bias model.** cwm\_fwd, cwm\_rev are the forward and reverse complemented consolidated motifs from contribution scores in subset of random peaks. These CWM motifs should be free from any Transcription Factor (TF) motifs and should contain either only bias motifs or random repeats. For each of these motifs, we use TOMTOM to find the top-3 closest matches (match\_0, match\_1, match\_2) from a database consisting of both MEME TF motifs and heterogenous enzyme bias motifs that we have repeatedly seen in our datasets. The qvals (qval0,qval1,qval2) should be high (> 0.0001) if the closest hit is a TF motif (i.e indicating that the closest match is not the correct match) - this is also generally verifiable by eye as the closest match will look nothing like the CWMs. The qvals should be low if the closest hit is enzyme bias motif and generally verifiable that the top match looks like the CWM. The first 3-5 motifs in the list below should look like enzyme bias motif.

#### What to do if you find an obvious TF motif in the list?

Do not use this bias model as it will regress the contribution of the TF motifs (along with bias motifs) from the chrombpnet\_nobias.h5. Reduce the bias\_threshold\_factor argument input to the *chrombpnet bias pipeline* or *chrombpnet bias train* command used in training the bias model and retrain a new bias model. For more intuition about this argument refer to the [FAQ](#) section in wiki.

#### What to do if you are unsure if a given CWM motif is resembling the match\_0 logo for example?

Get marginal footprint on the match\_0 motif logo (using the command *chrombpnet footprints* and make sure that the bias models footprint is closer to that of controls with no motif inserted - for examples look at [FAQ](#) )

pattern	NumSeqs	cwm_fwd	cwm_rev	match0	qval0	match0 logo	match1	qval1	match1 logo	match2	qval2	match2 logo
pos_0	8316			TN5_1	1.012050e-09		TN5_2	5.214110e-09		TN5_7	0.000281	
pos_1	6095			TN5_4	2.464950e-02		TN5_5	2.464950e-02		TN5_8	0.036485	
pos_2	4192			TN5_1	4.948020e-05		TN5_2	4.948020e-05		TN5_3	0.000727	
pos_3	3427			TN5_7	2.160010e-02		TN5_1	2.787090e-02		TN5_3	0.062360	
pos_4	3134			TN5_3	1.290400e-01		K_F4_M40039.3	1.290400e-01		TN5_4	0.129046	
pos_5	2569			TN5_2	6.214200e-11		TN5_4	3.156880e-03		TN5_5	0.003157	
pos_6	2524			TN5_3	4.274200e-09		TN5_4	2.671300e-06		TN5_5	0.000003	
pos_7	912			TN5_3	7.596790e-14		TN5_1	1.311130e-03		TEX21 TEX_3	0.017374	
pos_8	822			TN5_3	1.111940e-04		TN5_1	2.813550e-02		TEX1 TEX_4	0.087364	
pos_9	412			TN5_4	6.944640e-04		TN5_5	6.944640e-04		TN5_1	0.000094	
pos_10	232			ZNF384_M41125.1	6.239300e-02		PRDM6_HUMAN.H11MO.OA	6.239300e-02		STAT1_MOUSE.H11MO.OA	0.090932	
pos_11	209			ZNF384_M41125.1	1.000000e+00		HOXD12_homeodomain_1	1.000000e+00		Hoxc10.mouse_homeodomain_2	1.000000	
pos_12	164			ZNF384_M41125.1	6.609000e-01		FOXQ1_HUMAN.H11MO.OA	6.609000e-01		FOXQ1_MOUSE.H11MO.OA	0.660900	
pos_13	137			POU3F4_POU_2	3.798150e-01		POU3F2_POU_1	4.569570e-01		POU3F1_POU_2	0.574908	
pos_14	119			TN5_3	2.030680e-08		TN5_4	2.524990e-02		TN5_5	0.025250	
pos_15	116			RREB1_M40073.1	1.000000e+00		EGR2_HUMAN.H11MO.OA	1.000000e+00		EGR1_HUMAN.H11MO.OA	1.000000	

pattern	NumSeqs	cwm_fwd	cwm_rev	match0	qval0	match0 logo	match1	qval1	match1 logo	match2	qval2	match2 logo
pos_16	97			ZNF384 MA1125.1	1.000000e+00		ONECUT3 CUT 1	1.000000e+00		ONECUT3 MA0757.1	1.000000	
pos_17	51			ZNF384 MA1125.1	1.477330e-01		DNASE 2	1.477330e-01		BARX1 MOUSE.H11M0.0.C	1.000000	
pos_18	31			TNS 3	1.475930e-04		TNS 1	9.001770e-03		TNS 2	0.171835	
pos_19	23			ZNF18 HUMAN.H11M0.0.C	3.088100e-01		ZNF31 HUMAN.H11M0.0.C	3.088100e-01		TEX1 MA0805.1	0.308810	

**TFModisco motifs generated from counts contribution scores of the bias model.** cwm\_fwd, cwm\_rev are the forward and reverse complemented consolidated motifs from contribution scores in subset of random peaks. These motifs should be free from any Transcription Factor (TF) motifs and should contain motifs either weakly related to bias motifs or random repeats. For each of these motifs, we use TOMTOM to find the top-3 closest matches (match\_0, match\_1, match\_2) from a database consisting of both MEME TF motifs and heterogenous enzyme bias motifs that we have repeatedly seen in our datasets. The qvals should be high (> 0.0001) if the closest hit is a TF motif (i.e. indicating that the closest match is not the correct match, this is also generally verifiable by eye and making sure the closest match looks nothing like the CWMs).

#### What to do if you find an obvious TF motif in the list?

Do not use this bias model as it will regress the contribution of the TF motifs (along with bias motifs) from the chrombpnet\_nobias.h5. Reduce the bias\_threshold\_factor argument input to the *chrombpnet bias pipeline* or *chrombpnet bias train* command used in training the bias model and retrain a new bias model. For more intuition about this argument refer to the [FAQ](#) section in wiki.

#### What to do if you are unsure if a given CWM motif is resembling the match\_0 logo for example?

Get marginal footprint on the match\_0 motif logo (using the command *chrombpnet footprints* and make sure that the bias models footprint is closer to that of controls with no motif inserted - for examples look at [FAQ](#) )

pattern	NumSeqs	cwm_fwd	cwm_rev	match0	qval0	match0 logo	match1	qval1	match1 logo	match2	qval2	match2 logo
pos_0	99			SP2 HUMAN.H11M0.0.A	1.388380e-00		SP2 MOUSE.H11M0.0.B	1.388380e-00		SP1 HUMAN.H11M0.0.A	3.987050e-00	
pos_1	78			SP2 HUMAN.H11M0.0.A	2.572490e-00		SP2 MOUSE.H11M0.0.B	2.572490e-00		SP3 HUMAN.H11M0.0.B	2.572490e-00	
pos_2	58			ZSC31 HUMAN.H11M0.0.C	8.805050e-02		RELB MA1117.1	1.000000e+00		SOX10 HMG 3	1.000000e+00	
pos_3	53			WT1 HUMAN.H11M0.0.C	1.527400e-04		WT1 MOUSE.H11M0.0.B	1.527400e-04		SP3 HUMAN.H11M0.0.B	1.527400e-04	
pos_4	47			SP1 HUMAN.H11M0.0.A	1.194880e-02		SP2 HUMAN.H11M0.0.A	1.194880e-02		SP2 MOUSE.H11M0.0.B	1.194880e-02	
pos_5	33			SP2 HUMAN.H11M0.0.A	1.890820e-02		SP2 MOUSE.H11M0.0.B	1.890820e-02		SP1 HUMAN.H11M0.0.A	1.890820e-02	
pos_6	29			SP2 HUMAN.H11M0.0.A	7.030170e-04		SP2 MOUSE.H11M0.0.B	7.030170e-04		SP1 MOUSE.H11M0.0.A	1.200990e-03	