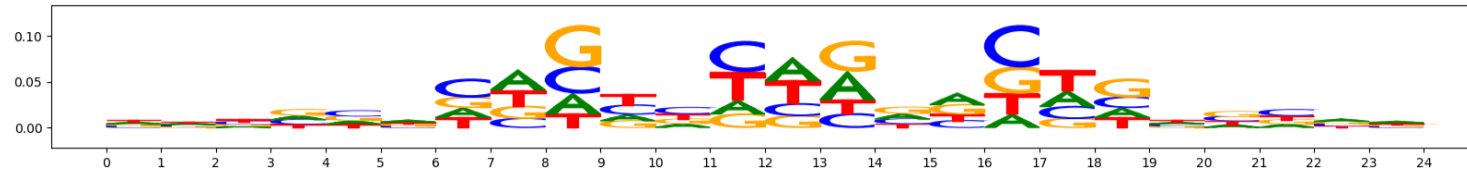


## Bias model training and quality check report

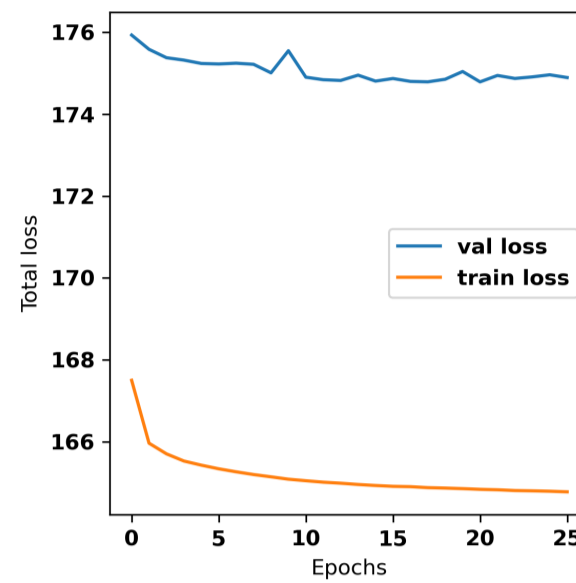
### Preprocessing report

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	0.17	0.51	-0.5	0.09
	nonpeaks.median jsd	nonpeaks.median norm jsd	peaks.median jsd	peaks.median norm jsd
profile metrics	0.04	0.15	0.5	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	9047			TN5 1	1.192810e-09		TN5 2	9.893430e-09		TN5 7	0.000164	
pos 1	6361			TN5 4	2.094450e-02		TN5 5	2.094450e-02		TN5 8	0.027038	
pos 2	4189			TN5 1	1.160210e-07		TN5 3	7.863950e-06		TN5 2	0.002422	
pos 3	3277			TN5 3	2.457370e-01		K_F4 MA0039.3	3.233650e-01		TN5 4	0.323305	
pos 4	3004			TN5 7	2.308130e-02		TN5 1	2.897190e-02		TN5 3	0.080177	
pos 5	2541			TN5 2	2.959900e-20		TN5 4	2.913900e-06		TN5 5	0.000003	
pos 6	1017			TN5 3	2.006130e-06		TN5 4	5.597970e-03		TN5 5	0.005598	
pos 7	949			TN5 3	1.465400e-06		TN5 4	5.462450e-06		TN5 5	0.000005	
pos 8	915			TN5 3	3.280700e-13		TN5 1	2.284950e-03		TEX21 TEX 3	0.026927	
pos 9	767			TN5 3	2.775620e-04		TN5 1	3.042320e-02		TN5 2	0.195808	
pos 10	627			TN5 3	6.143370e-04		TN5 4	2.443650e-02		TN5 5	0.024436	
pos 11	367			TN5 3	3.760040e-03		TN5 4	6.151670e-03		TN5 5	0.006152	
pos 12	146			DNASE 2	9.010930e-01		ZNF384 MA1125.1	9.010930e-01		POU3F3 POU 2	1.000000	
pos 13	133			ZNF384 MA1125.1	7.951510e-02		PRDM6 HUMAN.H11MO.0.C	8.377140e-02		STAT1 MOUSE.H11MO.0.A	0.182921	
pos 14	128			ZNF384 MA1125.1	6.785880e-02		PRDM6 HUMAN.H11MO.0.C	6.785880e-02		ANDR HUMAN.H11MO.0.A	0.213157	

pattern	NumSeqs	cwm_fwd	cwm_rev	match0	qval0	match0 logo	match1	qval1	match1 logo	match2	qval2	match2 logo
pos_15	126			HOMER1_HOMER1_1	1.000000e+00		NFATC1_NFAT_1	1.000000e+00		FOXO2_fothead_1	1.000000	
pos_16	98			TN5_3	6.799950e-02		TN5_1	6.915980e-01		TN5_4	0.948124	
pos_17	90			DNASE_4	1.000000e+00		SOX18_HMG_1	1.000000e+00		Sox10_mouse_HMG_1	1.000000	
pos_18	60			TN5_1	5.881600e-02		TN5_3	5.881600e-02		TN5_7	0.382920	
pos_19	58			TN5_8	7.104100e-04		RARG_HUMAN.H11M0.0.E	1.168770e-01		RARG_MOUSE.H11M0.0.C	0.116877	
pos_20	41			TN5_6	1.670890e-11		TEAD1_HUMAN.H11M0.0.A	6.624160e-01		TEAD2_MA1121.1	0.838196	
pos_21	25			TN5_4	9.094690e-03		TN5_5	9.094690e-03		TN5_7	1.000000	

**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	81			SP2_HUMAN.H11M0.0.A	7.049960e-07		SP2_MOUSE.H11M0.0.E	7.049960e-07		SP1_HUMAN.H11M0.0.A	4.239940e-06	
pos_1	70			ZFX_MOUSE.H11M0.0.E	1.000000e+00		TN5_2	1.000000e+00		None	NaN	
pos_2	69			PRDM6_HUMAN.H11M0.0.C	2.561030e-02		ZNF384_MA1125.1	7.721220e-02		ANDR_HUMAN.H11M0.0.A	1.101040e-01	
pos_3	67			MAZ_HUMAN.H11M0.0.A	1.127900e-05		MAZ_MOUSE.H11M0.0.A	1.127900e-05		SP5_MOUSE.H11M0.0.C	1.127900e-05	
pos_4	62			ZFX_MOUSE.H11M0.0.E	4.061940e-02		SP2_HUMAN.H11M0.0.A	4.061940e-02		SP2_MOUSE.H11M0.0.E	4.061940e-02	
pos_5	62			WT1_HUMAN.H11M0.0.C	1.476980e-02		WT1_MOUSE.H11M0.0.E	1.476980e-02		ZFX_MOUSE.H11M0.0.E	1.065530e-01	
pos_6	60			USF2_HUMAN.H11M0.0.A	1.000000e+00		USF2_MOUSE.H11M0.0.A	1.000000e+00		ZN331_HUMAN.H11M0.0.C	1.000000e+00	
pos_7	57			AP2B_HUMAN.H11M0.0.E	9.282070e-02		AP2B_MOUSE.H11M0.0.E	9.282070e-02		MX1_HUMAN.H11M0.0.A	1.579230e-01	
pos_8	53			VDR_MA0693.2	4.127990e-01		ZFX_MOUSE.H11M0.0.E	4.127990e-01		Zfx_MA0146.2	1.000000e+00	
pos_9	49			Zfx_MA0146.2	3.344920e-03		ZFX_MOUSE.H11M0.0.E	3.344920e-03		SP2_HUMAN.H11M0.0.A	1.292590e-02	
pos_10	47			ZN467_HUMAN.H11M0.0.C	3.245590e-04		VEZF1_HUMAN.H11M0.0.C	1.997320e-03		SP1_MOUSE.H11M0.0.A	1.997320e-03	
pos_11	47			KLF15_HUMAN.H11M0.0.A	4.757260e-01		KLF15_MOUSE.H11M0.0.A	4.757260e-01		TEF_MA0108.2	4.757260e-01	
pos_12	45			ASC_1_MA1100.1	5.426520e-01		SUH_MOUSE.H11M0.0.A	5.426520e-01		FOXO3_fothead_3	5.426520e-01	
pos_13	45			ZFX_MOUSE.H11M0.0.E	1.000000e+00		E2F4_MA0470.1	1.000000e+00		None	NaN	
pos_14	42			SP1_HUMAN.H11M0.0.A	4.237310e-03		EGR1_MOUSE.H11M0.0.A	6.153530e-03		KLF1_HUMAN.H11M0.0.A	7.581450e-03	
pos_15	42			TN5_1	3.121450e-01		NR1_2_HUMAN.H11M0.0.C	4.107380e-01		NR1_2_MOUSE.H11M0.0.C	4.107380e-01	
pos_16	41			SP1_HUMAN.H11M0.0.A	1.814860e-01		SP2_HUMAN.H11M0.0.A	1.814860e-01		SP2_MOUSE.H11M0.0.E	1.814860e-01	
pos_17	35			TN5_2	6.725560e-04		TN5_8	9.555670e-01		TN5_1	1.000000e+00	
pos_18	30			BREB1_MA0073.1	1.000000e+00		ZNF524_C2H2_2	1.000000e+00		PAX5_MA0014.3	1.000000e+00	
pos_19	22			MAX+MYC_MA0059.1	1.000000e+00		MAX_MA0058.3	1.000000e+00		MAX_DH_H_2	1.000000e+00	
pos_20	22			VEZF1_HUMAN.H11M0.0.C	7.575740e-04		ZN467_HUMAN.H11M0.0.C	7.575740e-04		MAZ_HUMAN.H11M0.0.A	9.410710e-04	