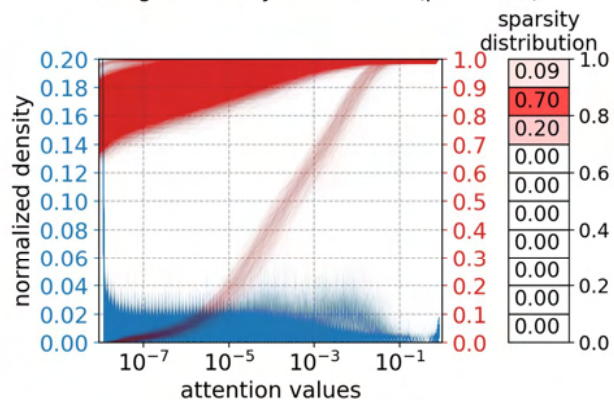
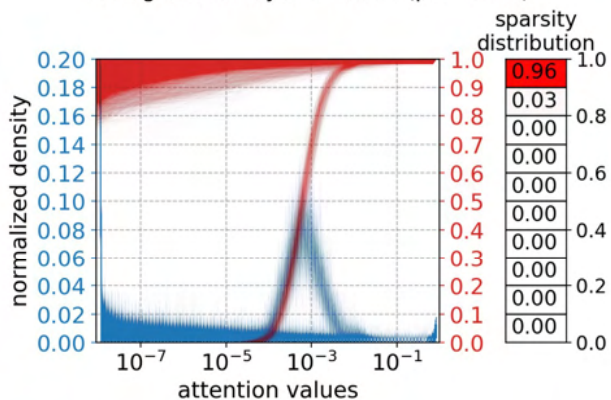


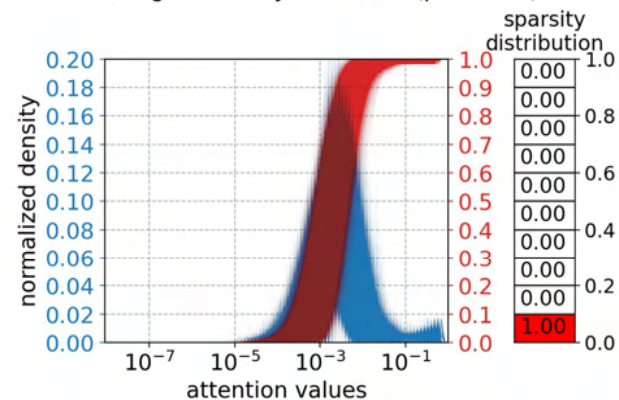
Histogram for layer 0 head 11(per token)



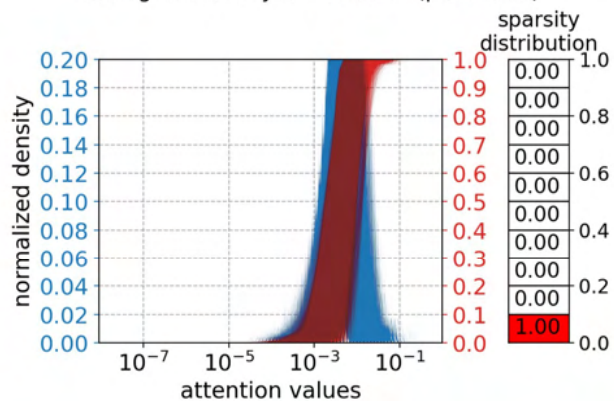
Histogram for layer 1 head 2(per token)



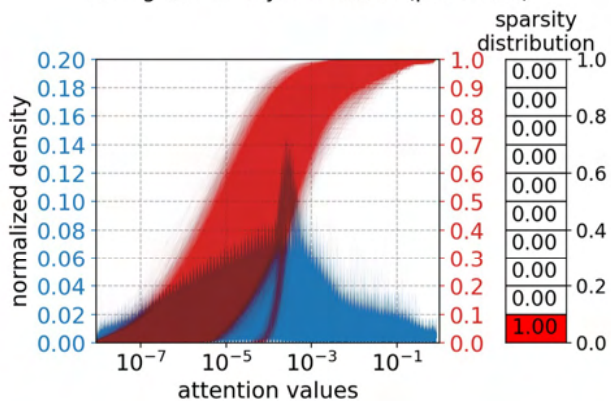
Histogram for layer 1 head 5(per token)



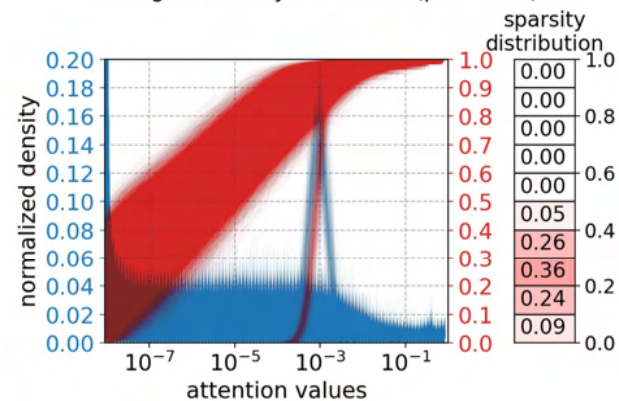
Histogram for layer 0 head 10(per token)



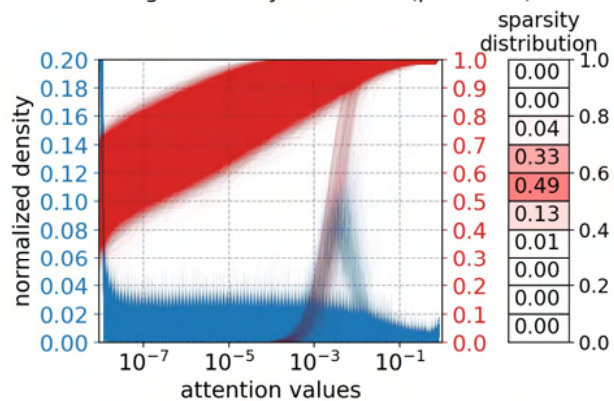
Histogram for layer 1 head 1(per token)



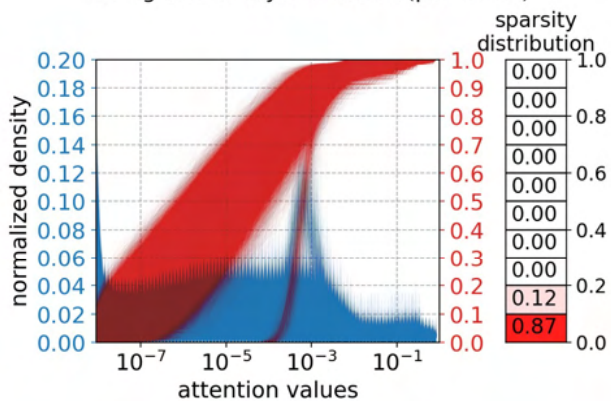
Histogram for layer 1 head 4(per token)



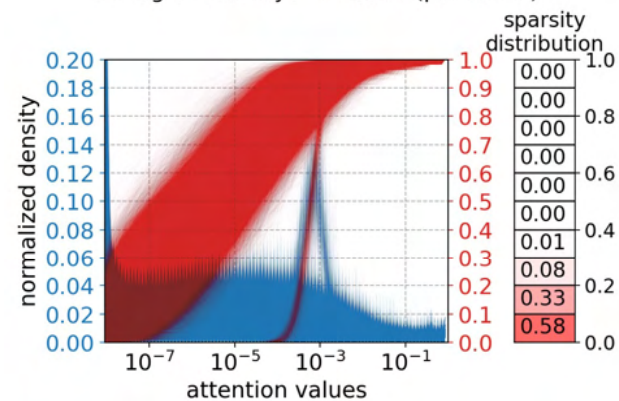
Histogram for layer 0 head 9(per token)



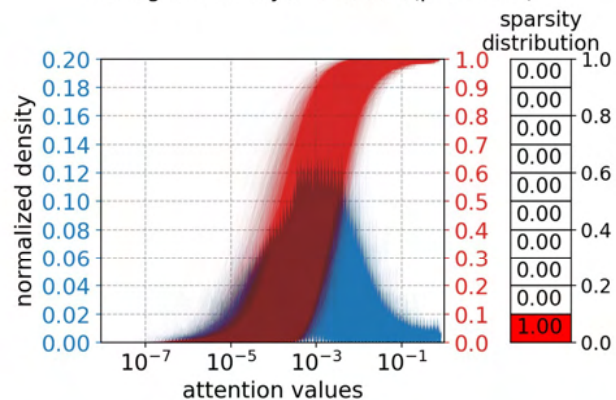
Histogram for layer 1 head 0(per token)



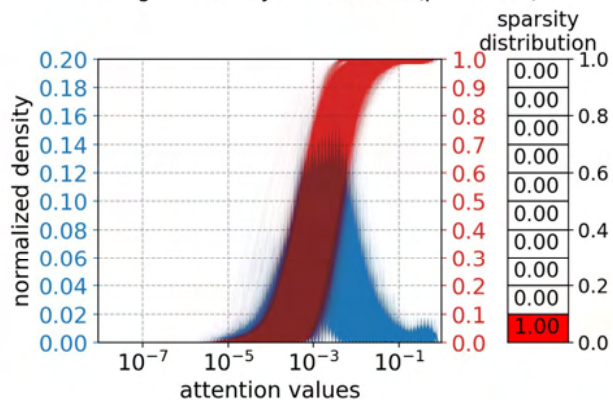
Histogram for layer 1 head 3(per token)



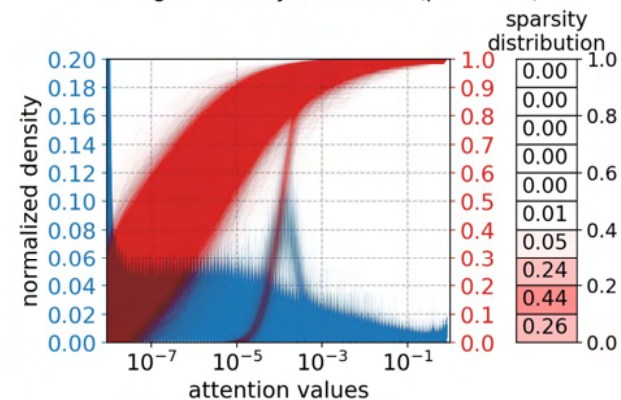
Histogram for layer 1 head 8(per token)



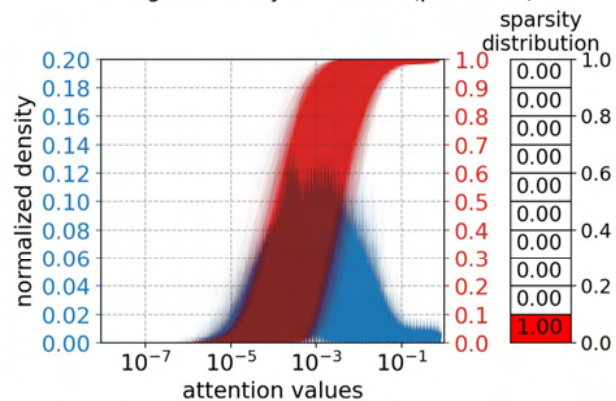
Histogram for layer 1 head 11(per token)



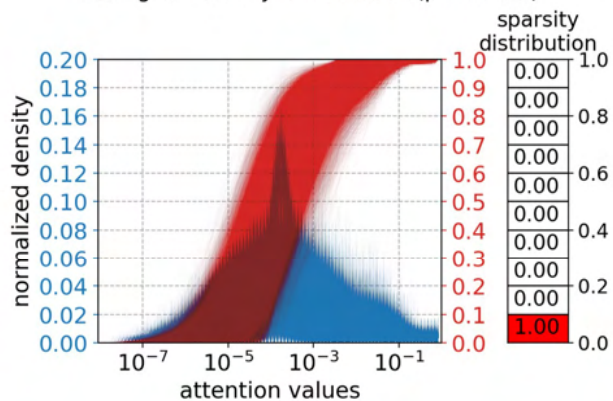
Histogram for layer 2 head 2(per token)



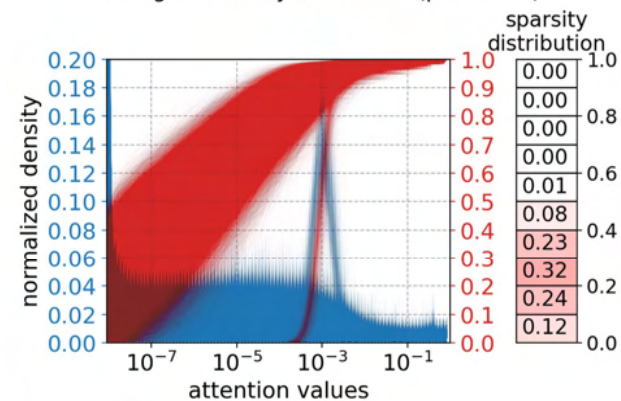
Histogram for layer 1 head 7(per token)



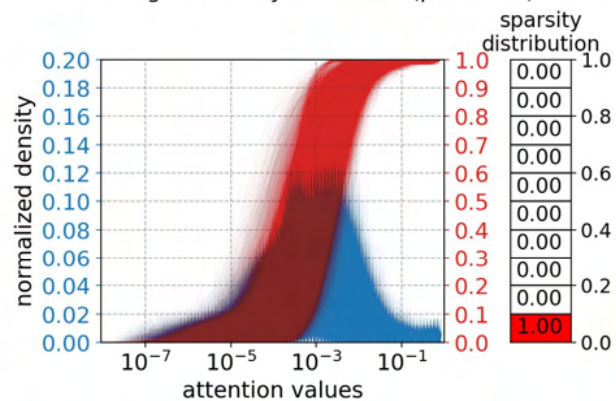
Histogram for layer 1 head 10(per token)



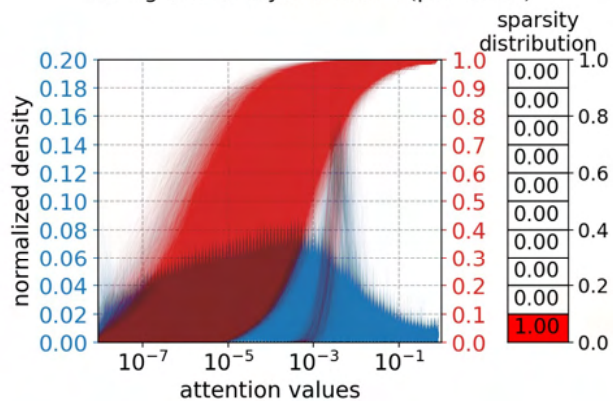
Histogram for layer 2 head 1(per token)



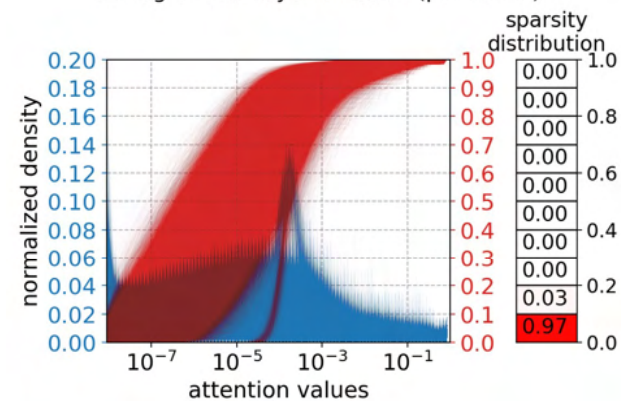
Histogram for layer 1 head 6(per token)



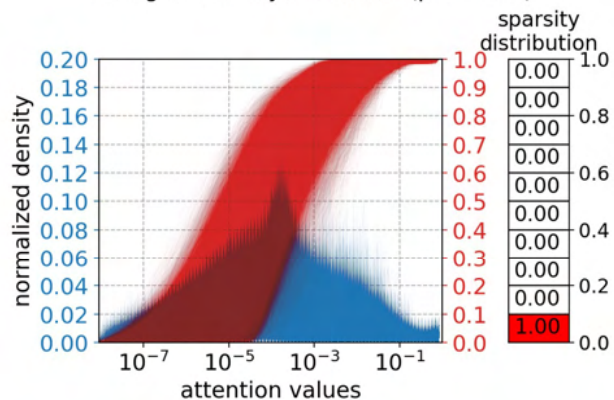
Histogram for layer 1 head 9(per token)



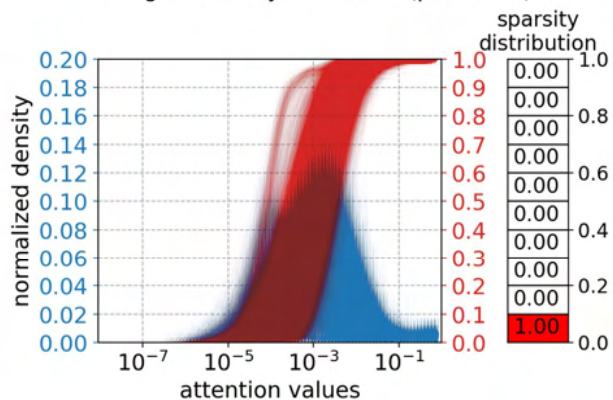
Histogram for layer 2 head 0(per token)



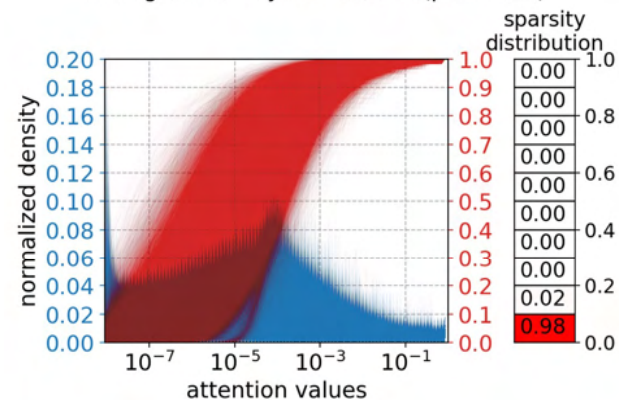
Histogram for layer 2 head 5(per token)



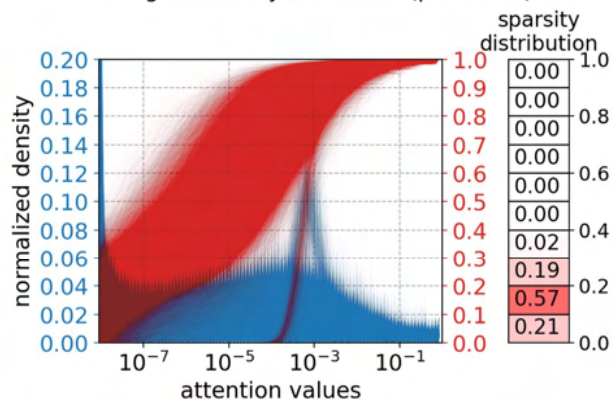
Histogram for layer 2 head 8(per token)



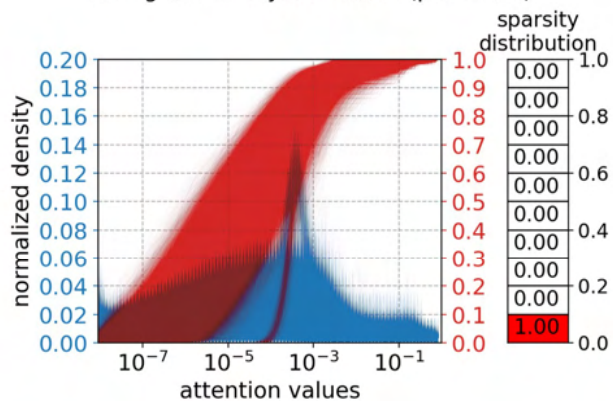
Histogram for layer 2 head 11(per token)



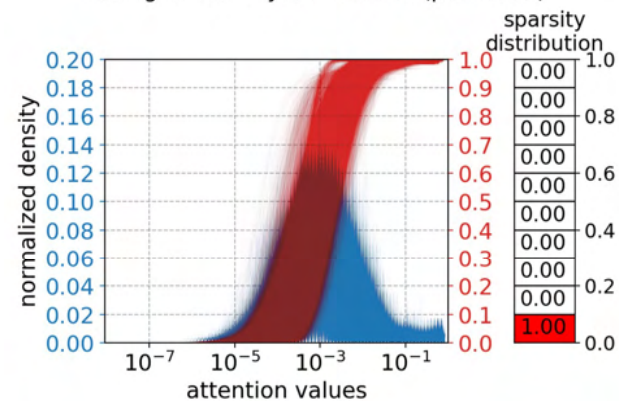
Histogram for layer 2 head 4(per token)



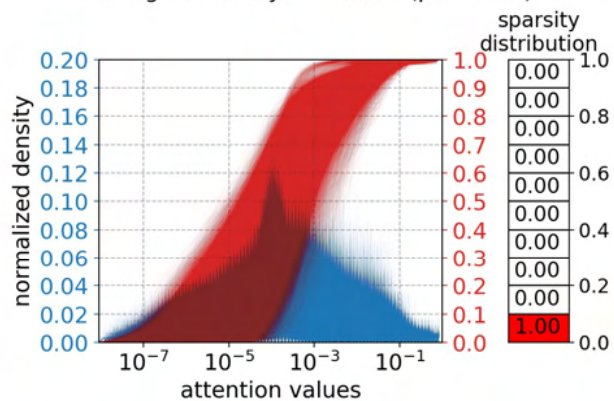
Histogram for layer 2 head 7(per token)



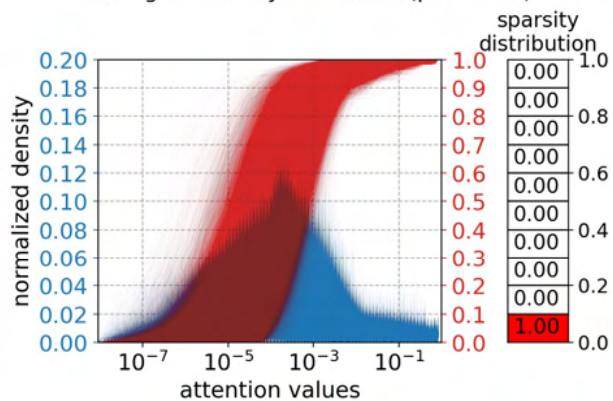
Histogram for layer 2 head 10(per token)



Histogram for layer 2 head 3(per token)



Histogram for layer 2 head 6(per token)



Histogram for layer 2 head 9(per token)

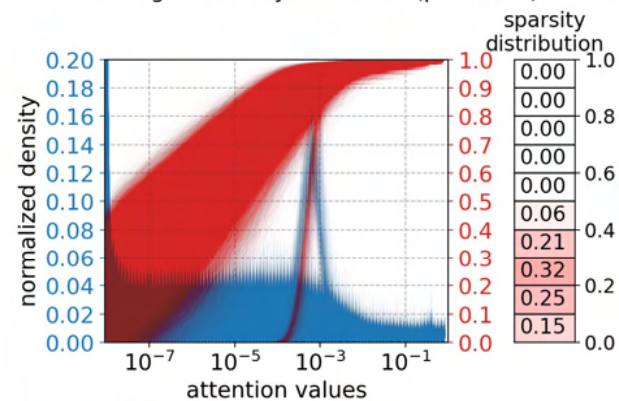


Figure 1 is a plot showing the normalized density of attention values for different sparsity distributions. The x-axis represents attention values on a logarithmic scale from 10^{-7} to 10^{-1} . The y-axis represents the normalized density from 0.00 to 0.20. The plot displays multiple curves corresponding to different sparsity distributions, ranging from 0.0 to 1.0. As the sparsity distribution increases, the distribution of attention values shifts towards higher values, indicating a concentration of attention on a smaller number of tokens.

Figure 1 is a plot showing the normalized density of attention values for different sparsity distributions. The x-axis represents attention values on a logarithmic scale from 10^{-7} to 10^{-1} . The y-axis represents the normalized density from 0.00 to 0.20. A red curve represents the sparsity distribution, and a blue curve represents the attention values. A table on the right shows the sparsity distribution for each attention value.

attention values	sparsity distribution
0.00	1.0
0.00	0.9
0.00	0.8
0.00	0.7
0.00	0.6
0.00	0.5
0.00	0.4
0.00	0.3
0.00	0.2
0.02	0.1
0.98	0.0


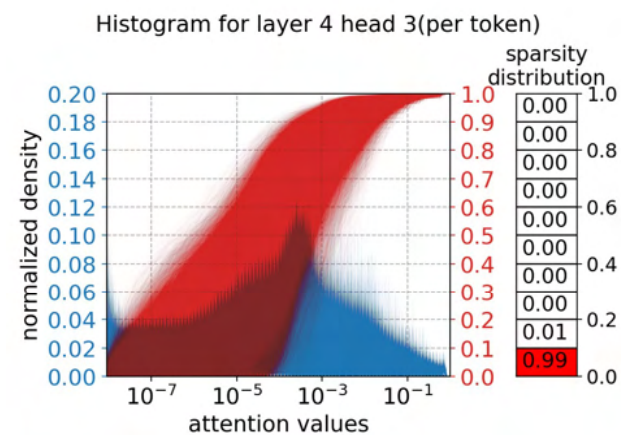
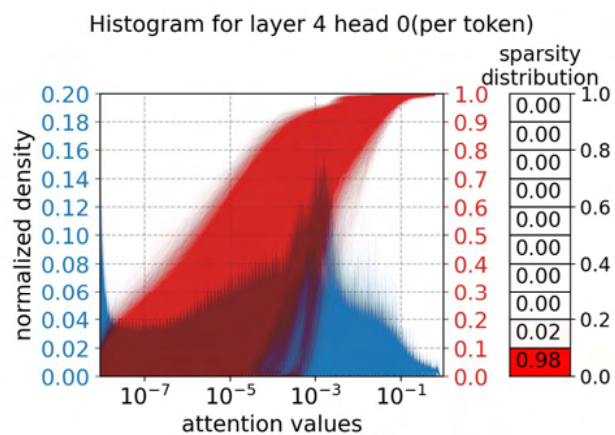
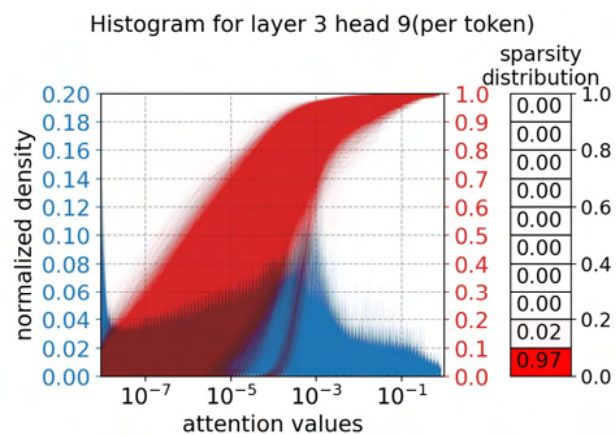
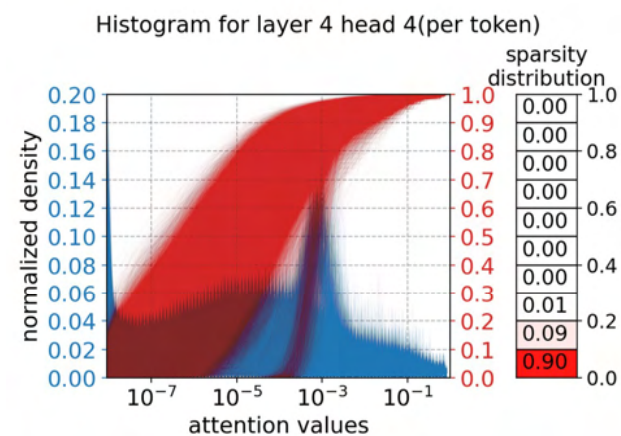
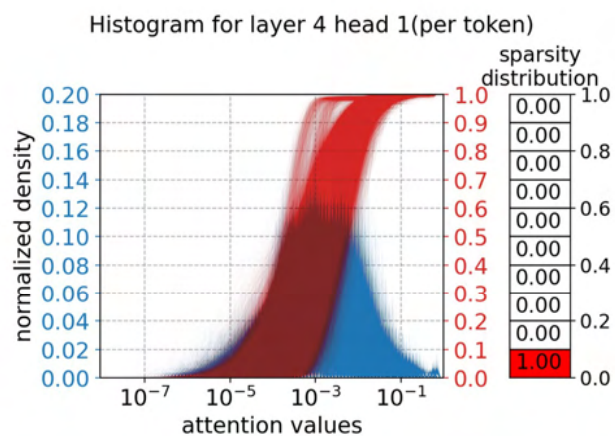
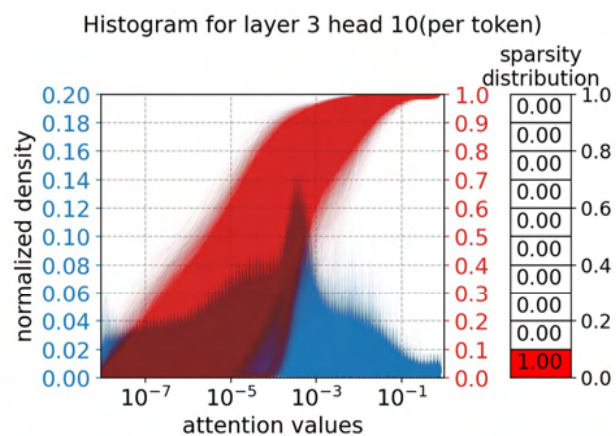
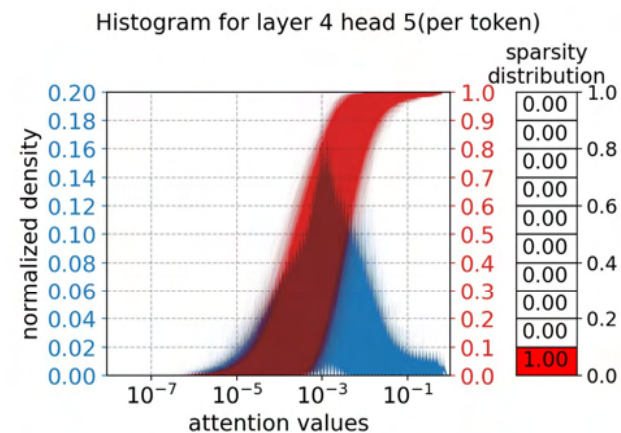
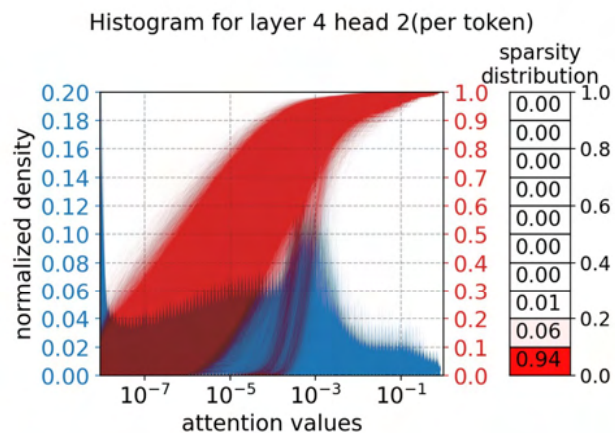
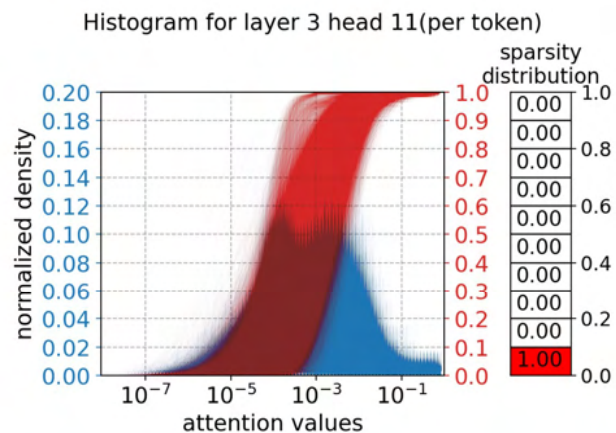


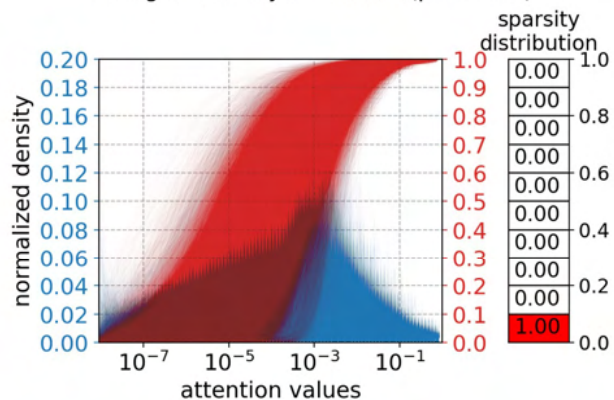
Figure 1: A plot showing the distribution of attention values for different sparsity levels. The x-axis is 'attention values' on a log scale from 10^{-7} to 10^{-1} . The y-axis is 'normalized density' from 0.00 to 0.20. A series of curves are shown for sparsity levels from 0.0 to 1.0. As sparsity increases, the distribution shifts from a broad peak around 10^{-4} to a sharp peak around 10^{-3} . A table on the right shows the sparsity distribution for each curve, with the 1.0 sparsity level highlighted in red.

sparsity	distribution
1.0	0.00
0.9	0.00
0.8	0.00
0.7	0.00
0.6	0.00
0.5	0.00
0.4	0.00
0.3	0.00
0.2	0.00
0.1	0.00
0.0	1.00

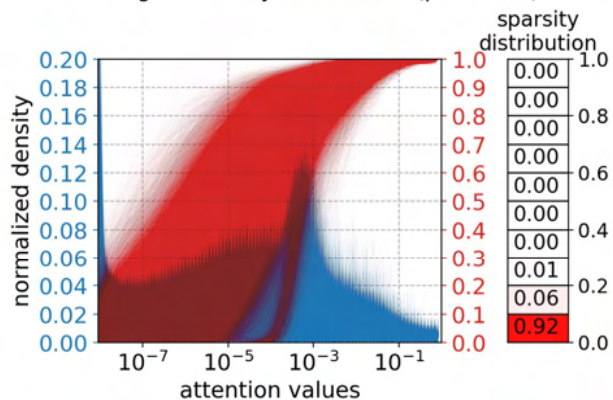
Figure 1: A plot showing the normalized density of attention values for different sparsity levels. The x-axis is 'attention values' on a log scale from 10^{-7} to 10^{-1} . The left y-axis is 'normalized density' from 0.00 to 0.20. The right y-axis is 'sparsity distribution' from 0.0 to 1.0. A red curve represents the density for sparsity 0.92, and a blue curve represents the density for sparsity 0.07. The red curve is shifted to the right (higher attention values) compared to the blue curve.



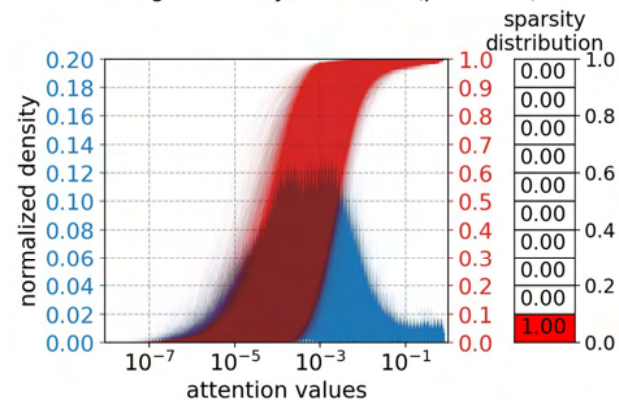
Histogram for layer 4 head 8(per token)



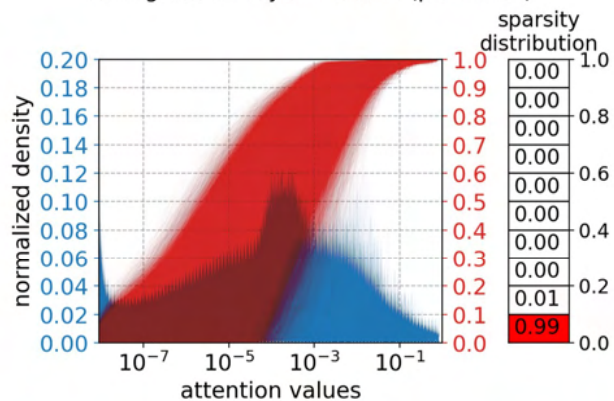
Histogram for layer 4 head 11(per token)



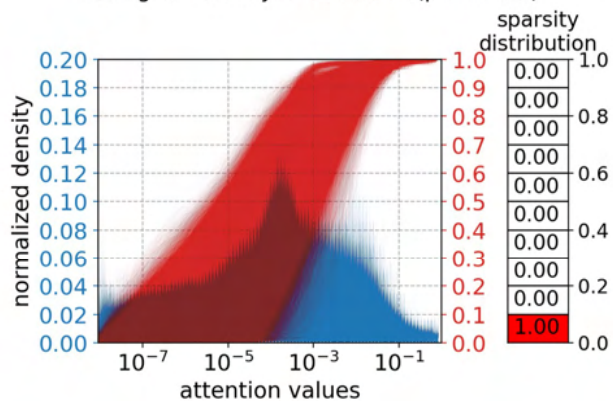
Histogram for layer 5 head 2(per token)



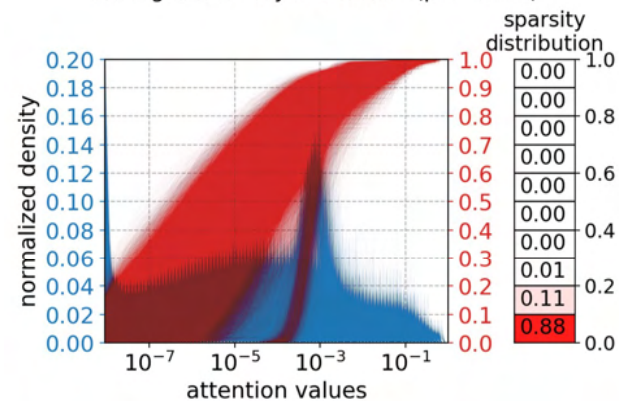
Histogram for layer 4 head 7(per token)



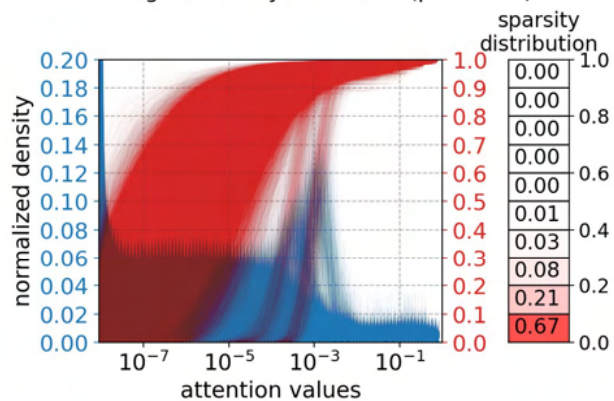
Histogram for layer 4 head 10(per token)



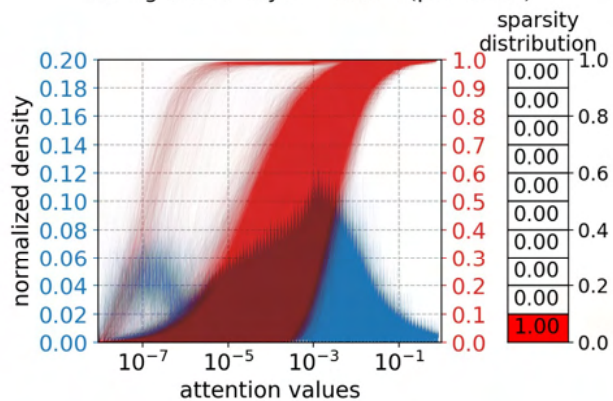
Histogram for layer 5 head 1(per token)



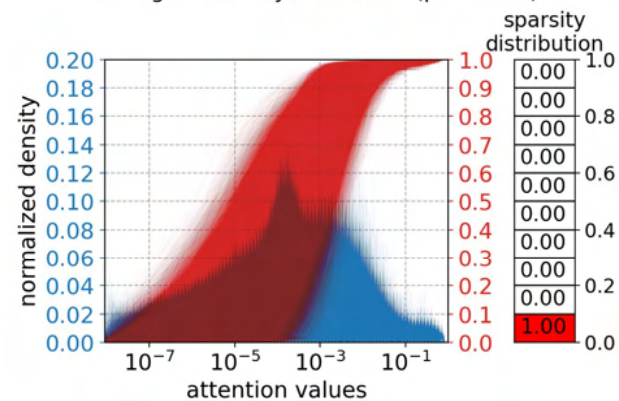
Histogram for layer 4 head 6(per token)

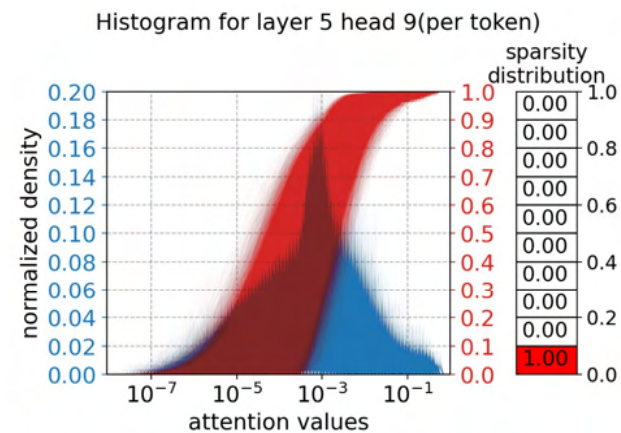
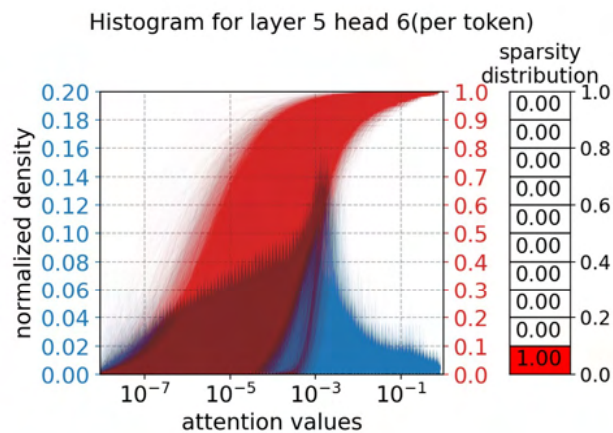
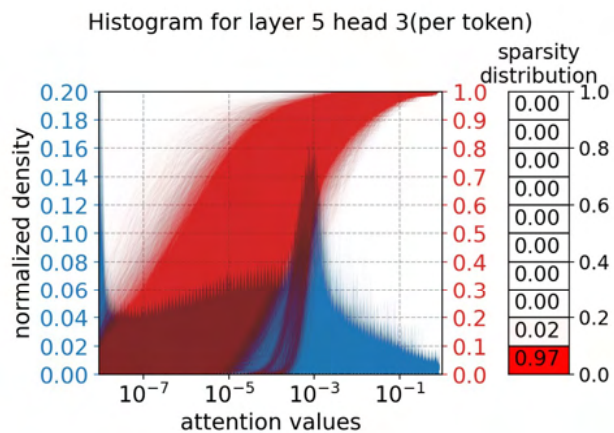
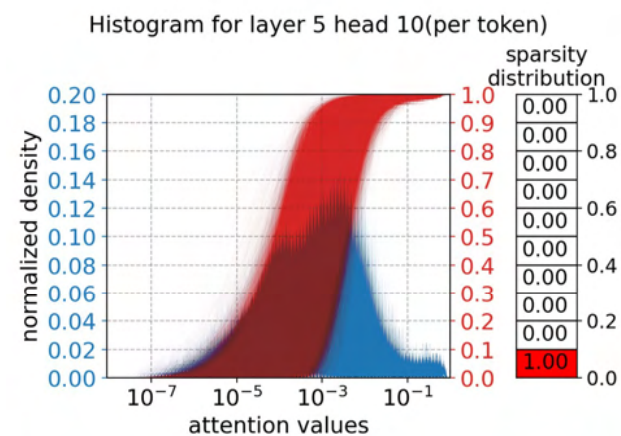
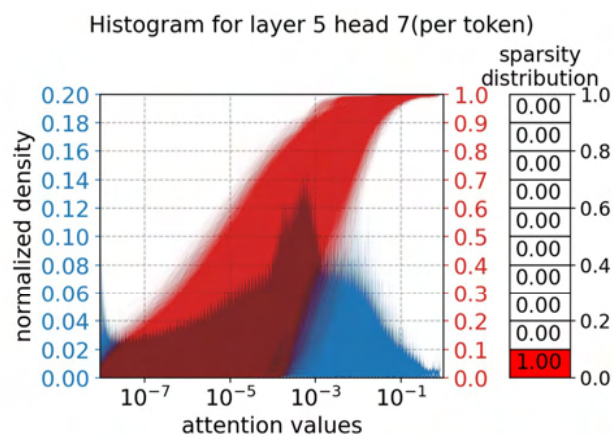
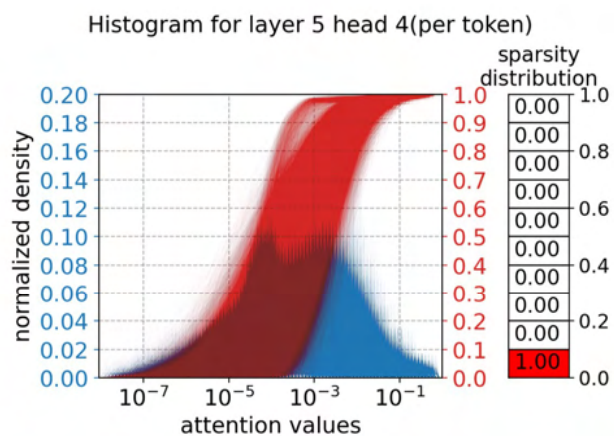
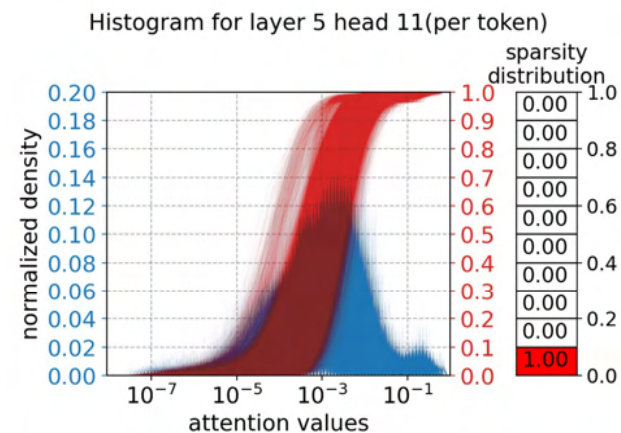
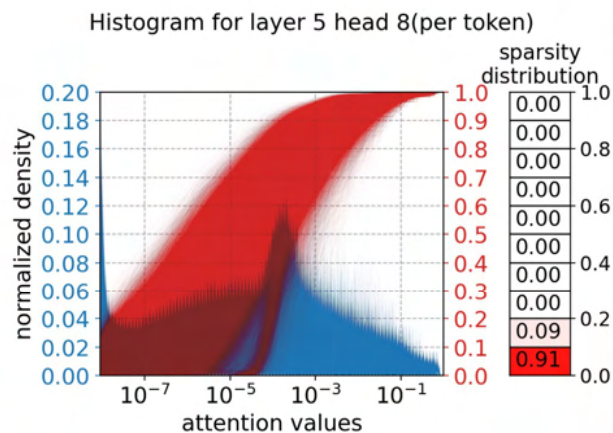
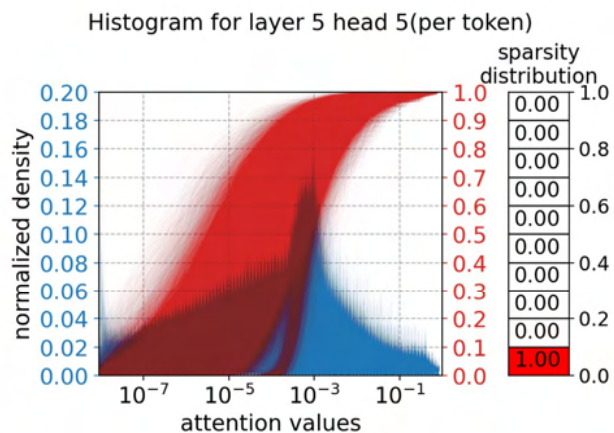


Histogram for layer 4 head 9(per token)

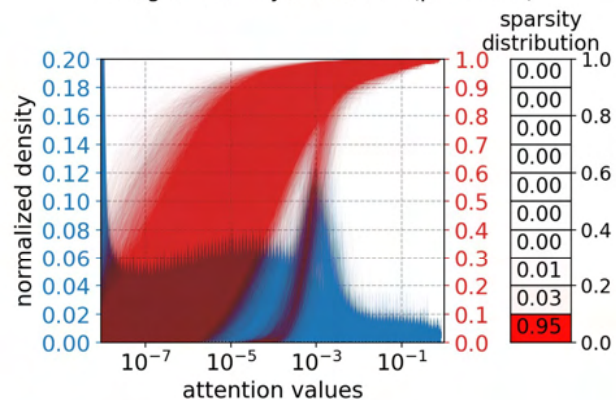


Histogram for layer 5 head 0(per token)

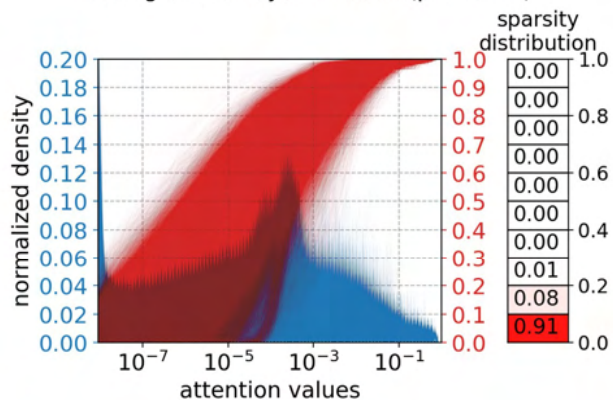




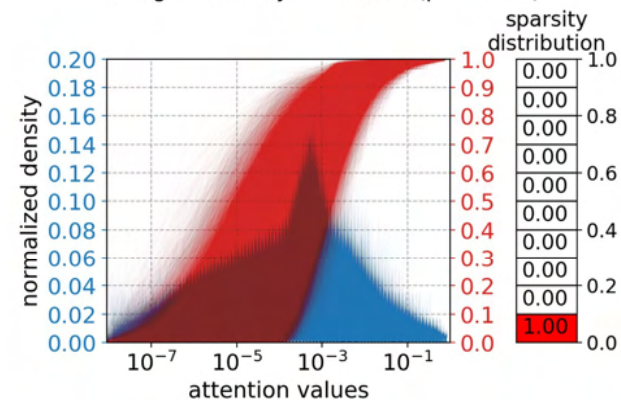
Histogram for layer 6 head 2(per token)



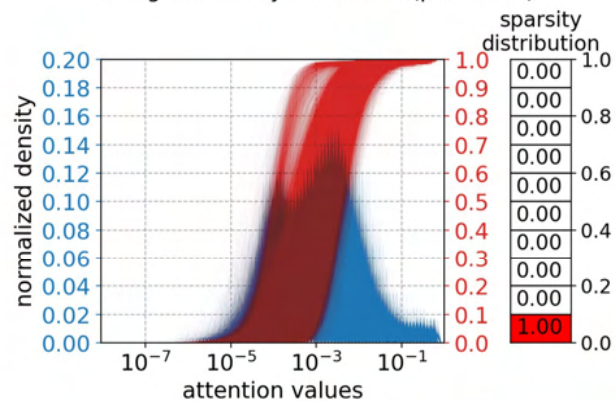
Histogram for layer 6 head 5(per token)



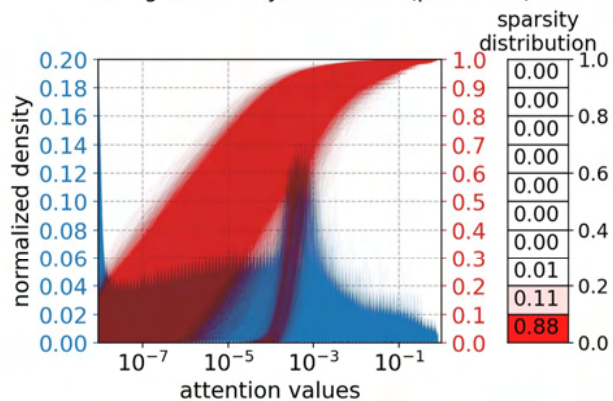
Histogram for layer 6 head 8(per token)



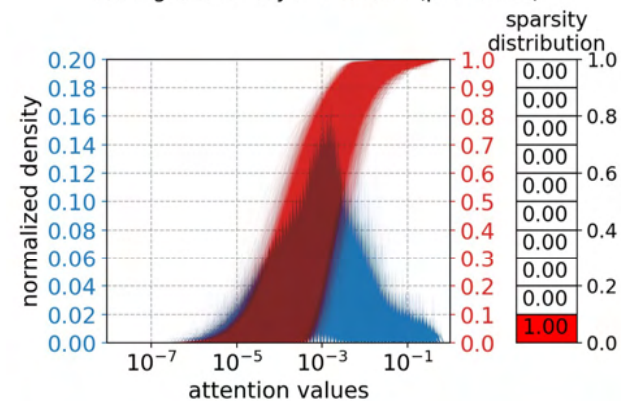
Histogram for layer 6 head 1(per token)



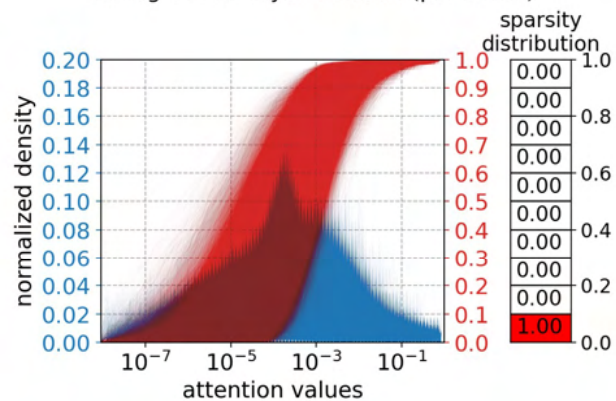
Histogram for layer 6 head 4(per token)



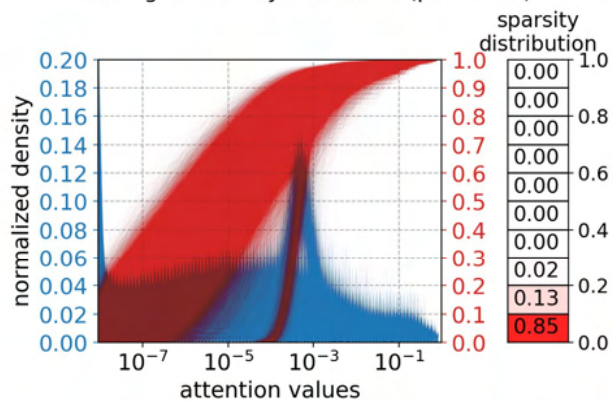
Histogram for layer 6 head 7(per token)



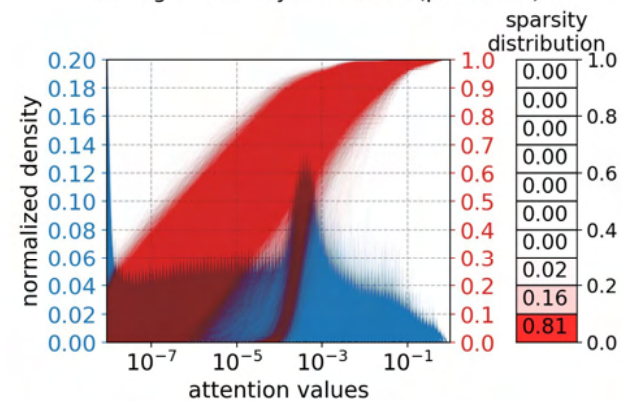
Histogram for layer 6 head 0(per token)

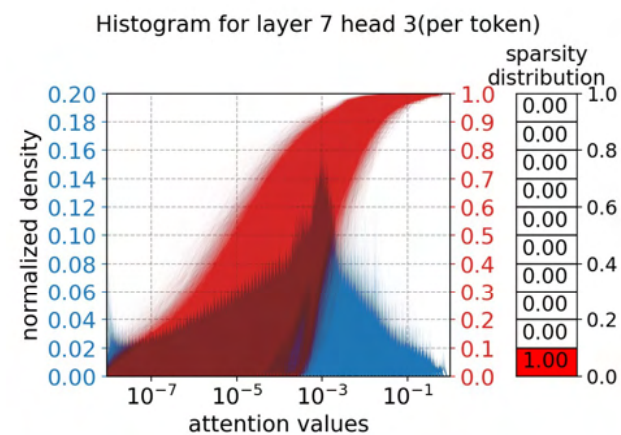
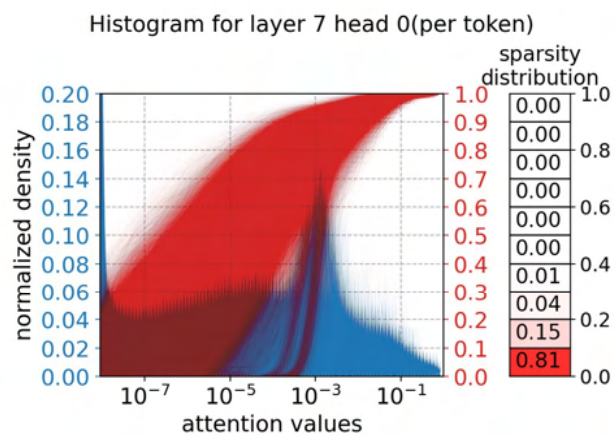
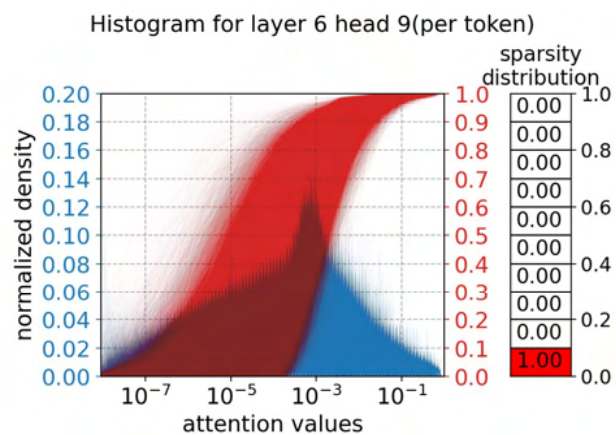
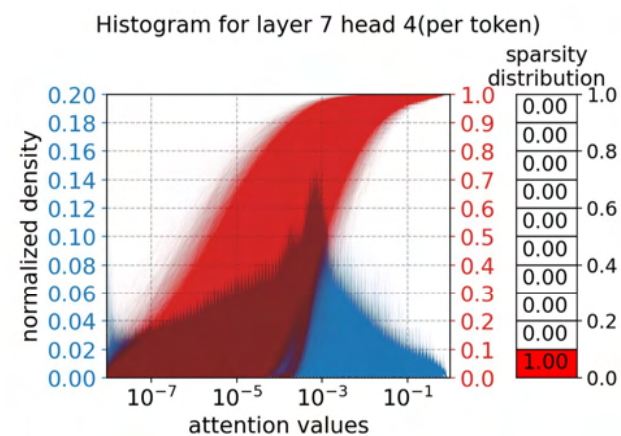
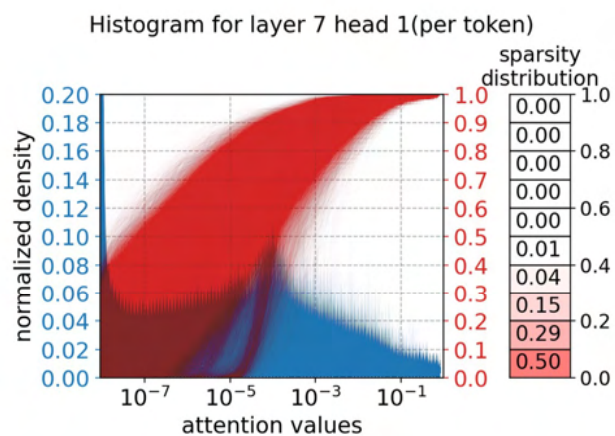
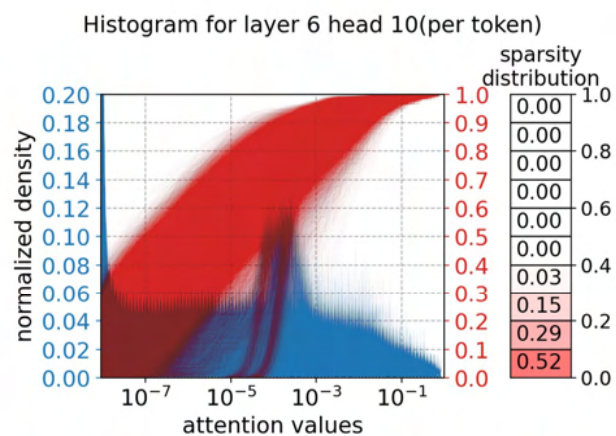
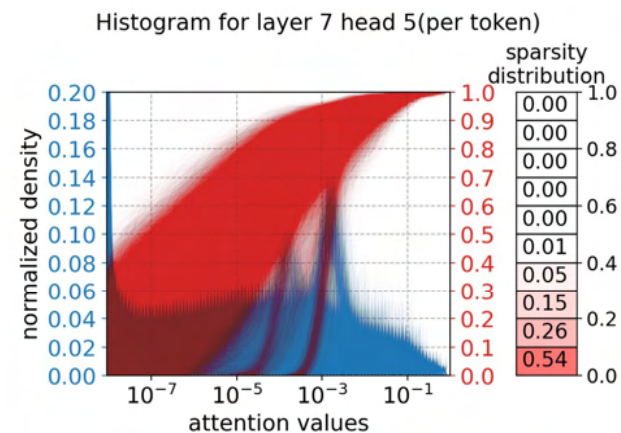
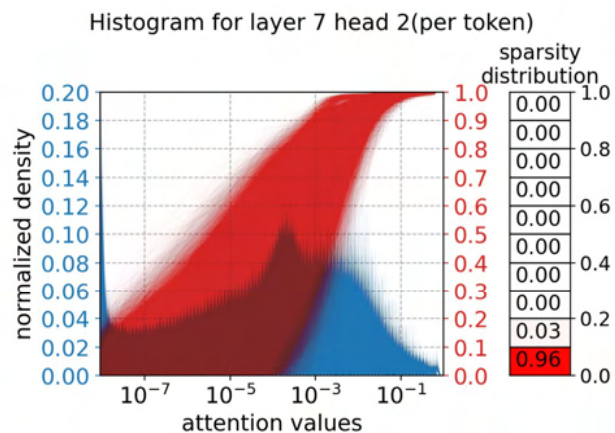
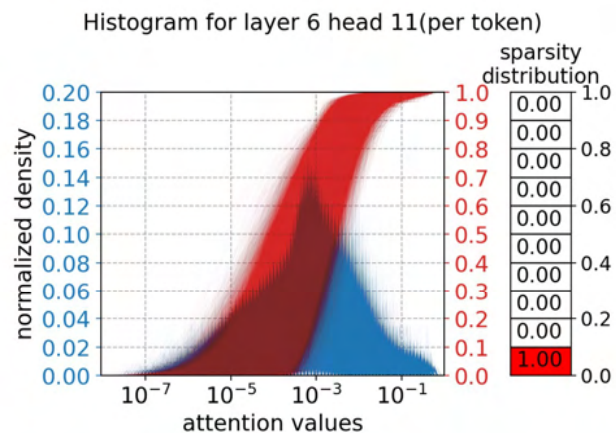


Histogram for layer 6 head 3(per token)

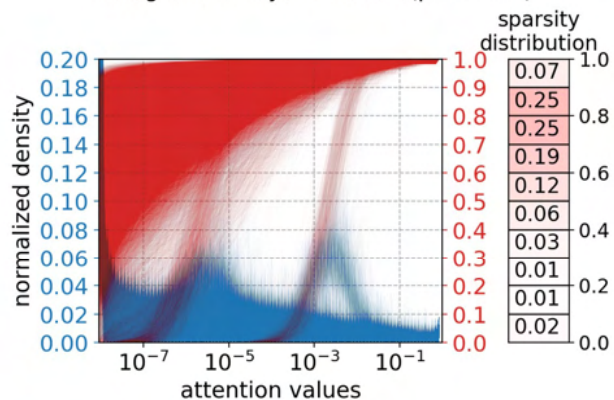


Histogram for layer 6 head 6(per token)

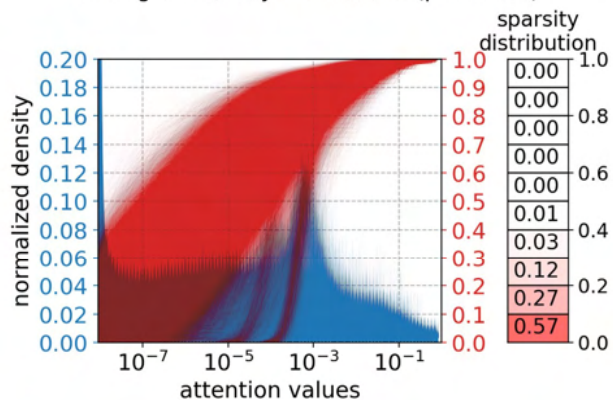




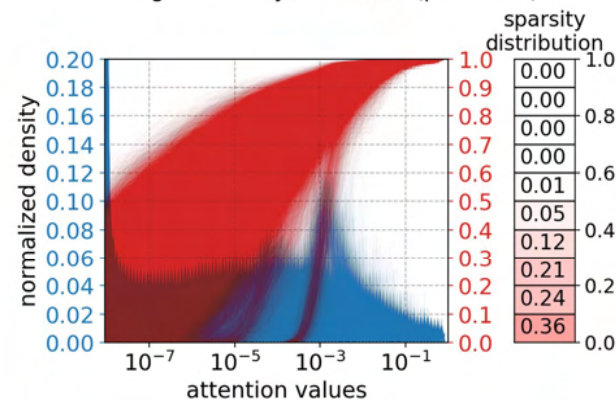
Histogram for layer 7 head 8(per token)



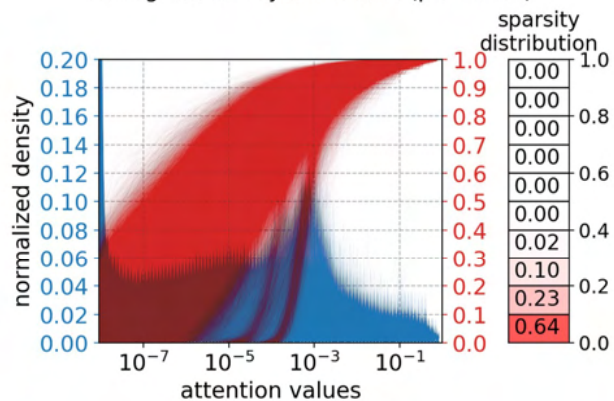
Histogram for layer 7 head 11(per token)



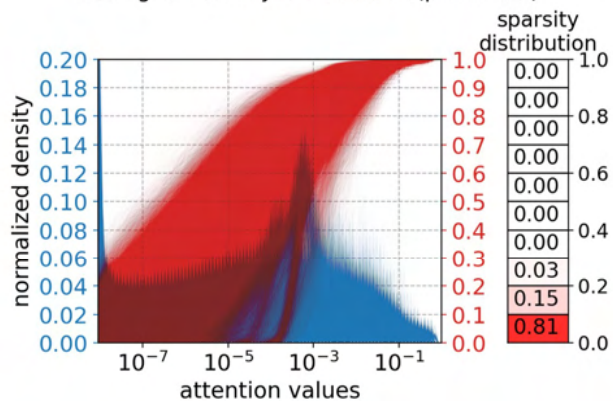
Histogram for layer 8 head 2(per token)



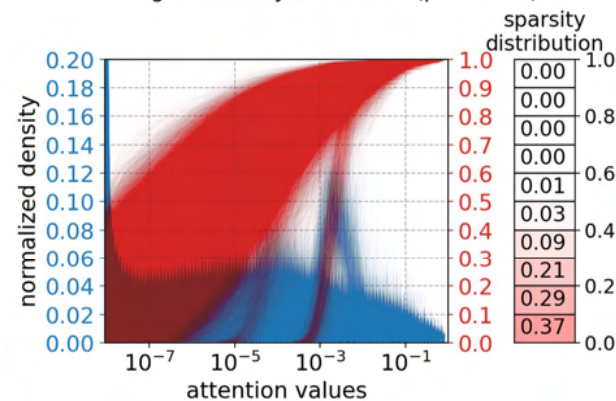
Histogram for layer 7 head 7(per token)



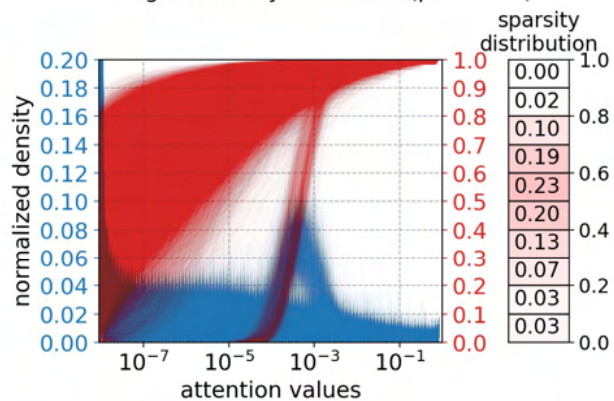
Histogram for layer 7 head 10(per token)



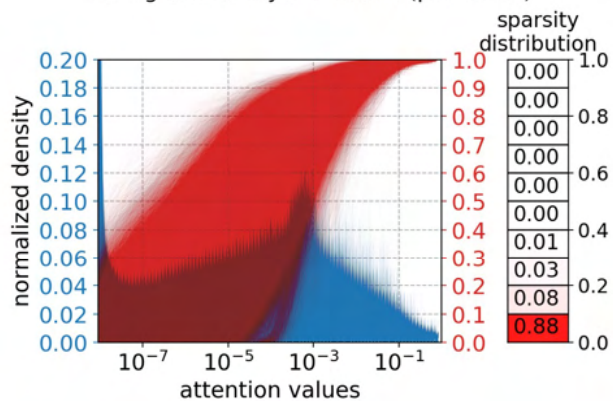
Histogram for layer 8 head 1(per token)



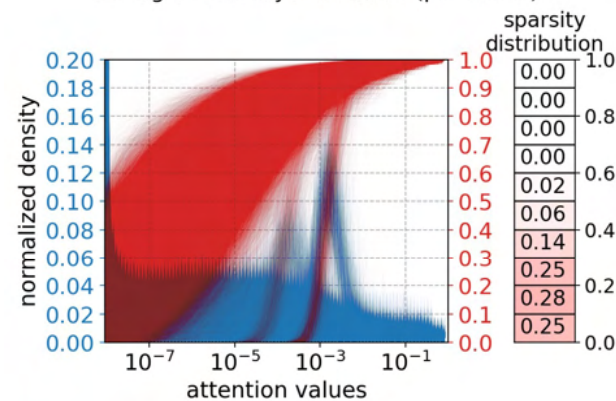
Histogram for layer 7 head 6(per token)



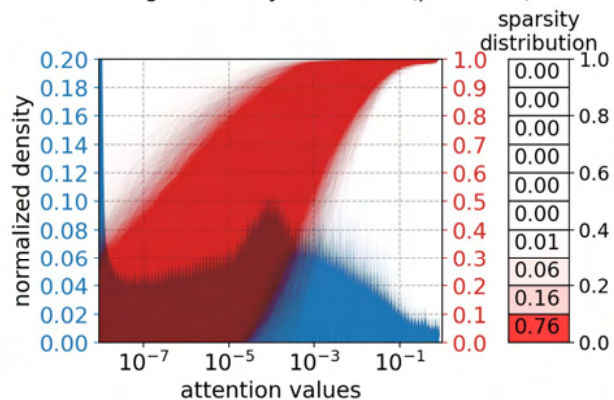
Histogram for layer 7 head 9(per token)



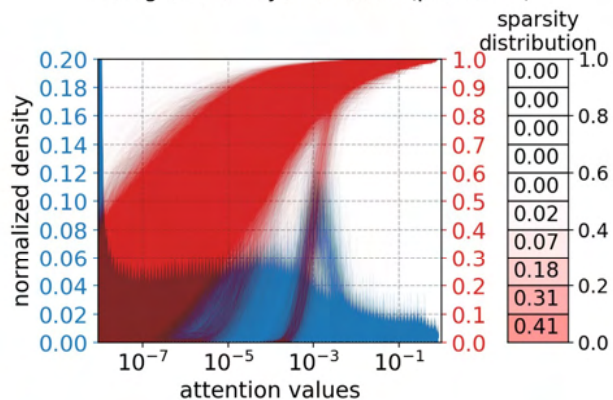
Histogram for layer 8 head 0(per token)



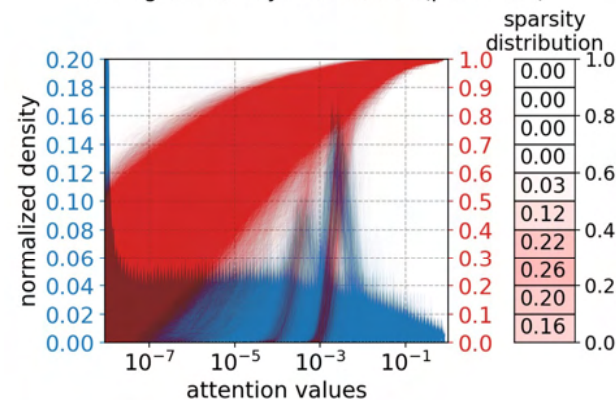
Histogram for layer 8 head 5(per token)



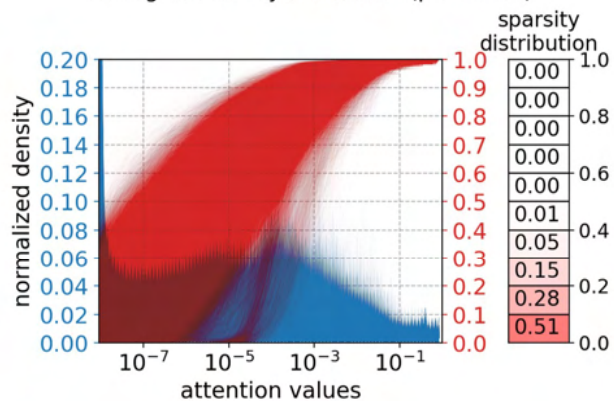
Histogram for layer 8 head 8(per token)



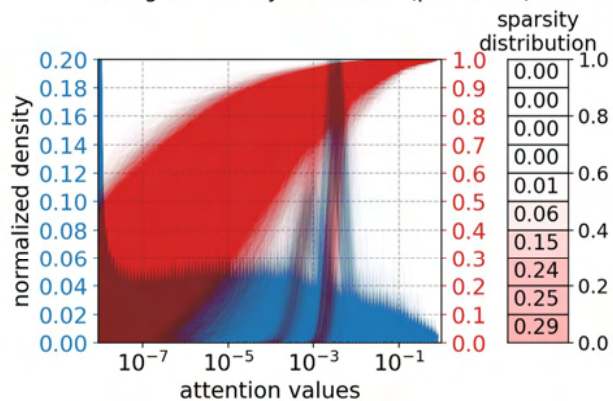
Histogram for layer 8 head 11(per token)



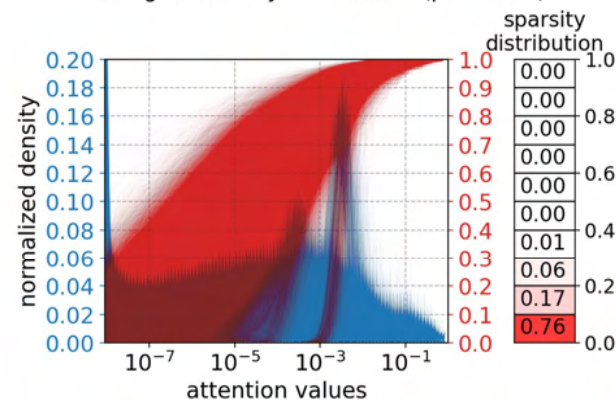
Histogram for layer 8 head 4(per token)



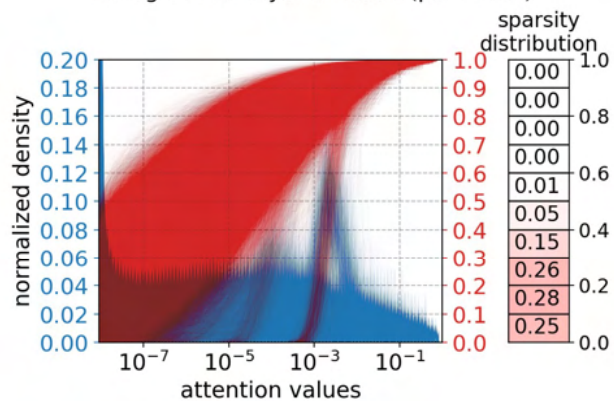
Histogram for layer 8 head 7(per token)



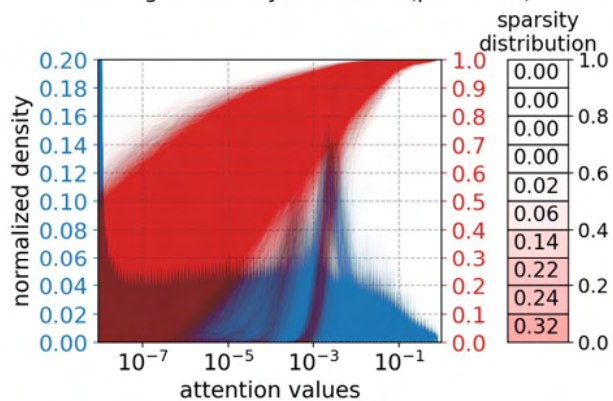
Histogram for layer 8 head 10(per token)



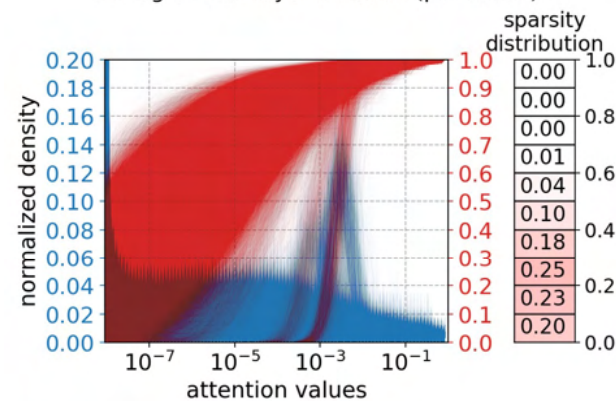
Histogram for layer 8 head 3(per token)



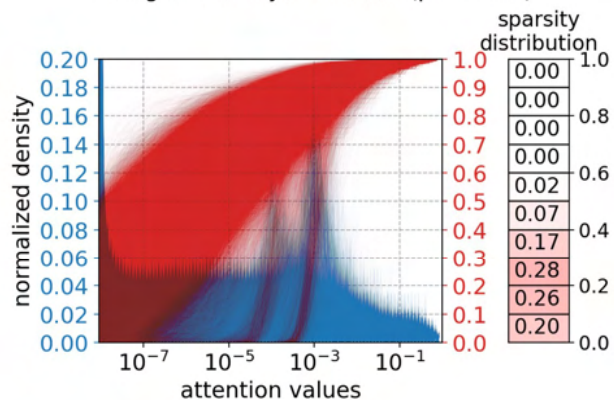
Histogram for layer 8 head 6(per token)



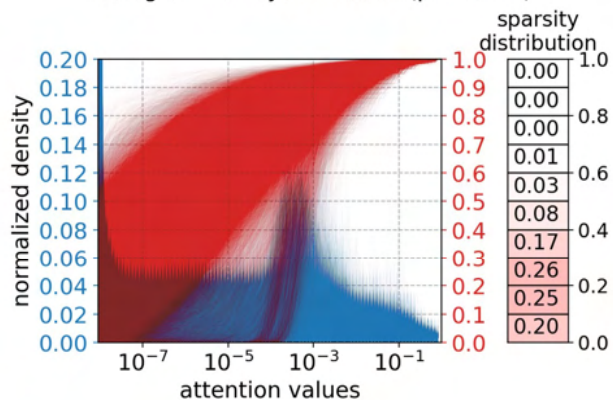
Histogram for layer 8 head 9(per token)



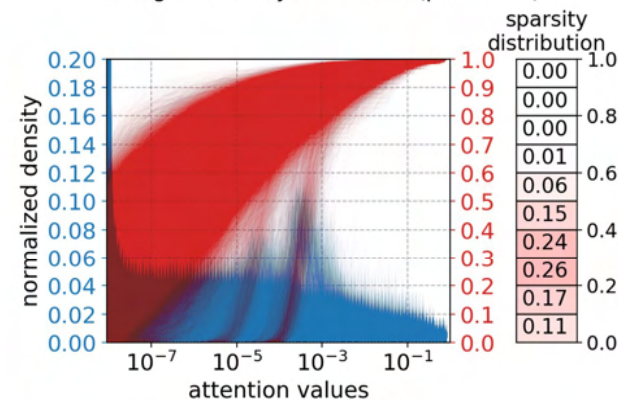
Histogram for layer 9 head 2(per token)



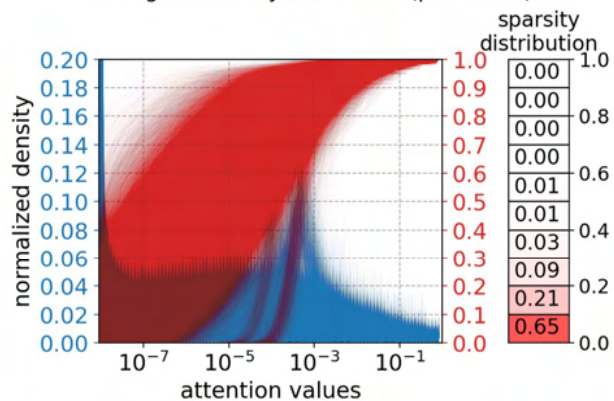
Histogram for layer 9 head 5(per token)



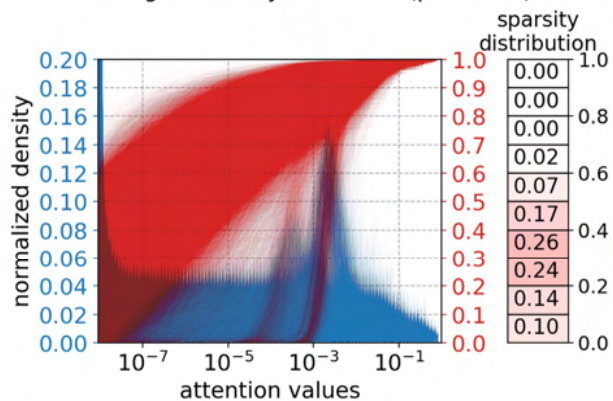
Histogram for layer 9 head 8(per token)



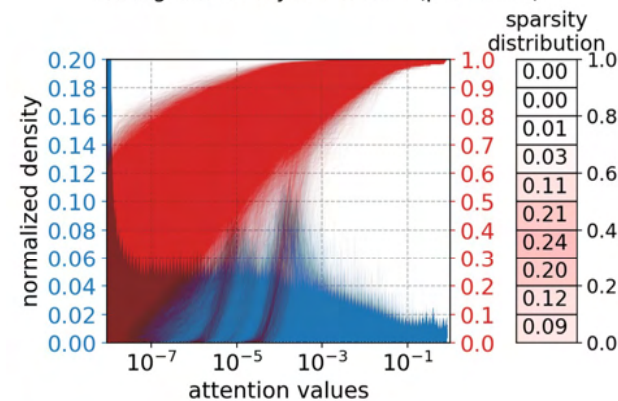
Histogram for layer 9 head 1(per token)



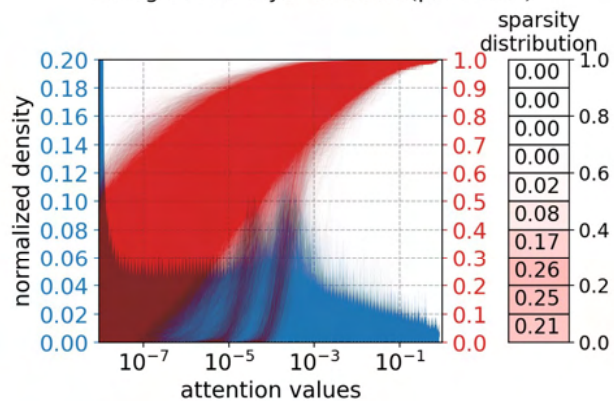
Histogram for layer 9 head 4(per token)



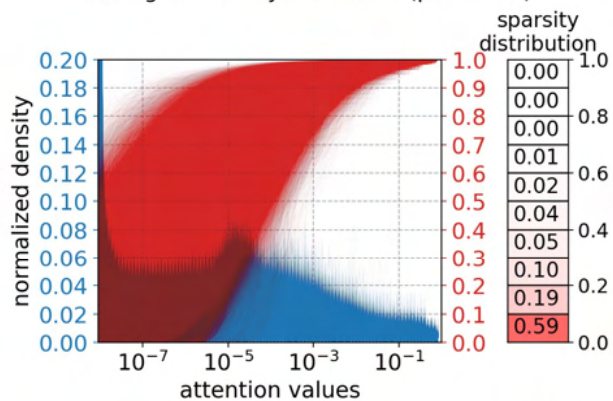
Histogram for layer 9 head 7(per token)



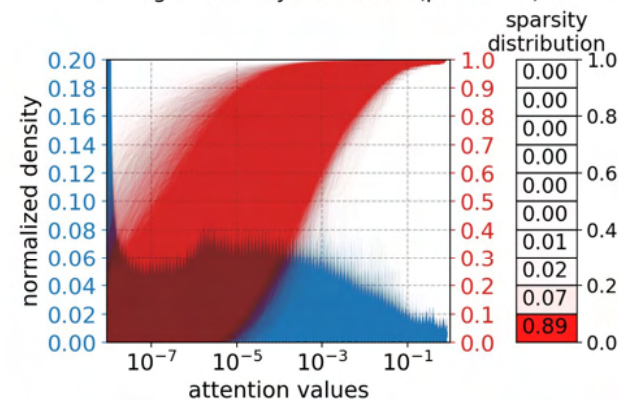
Histogram for layer 9 head 0(per token)

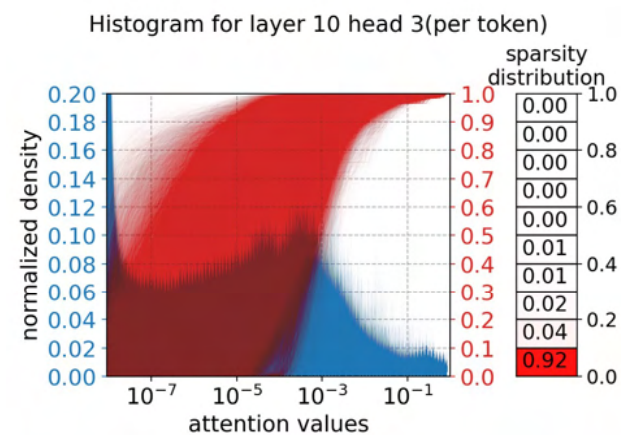
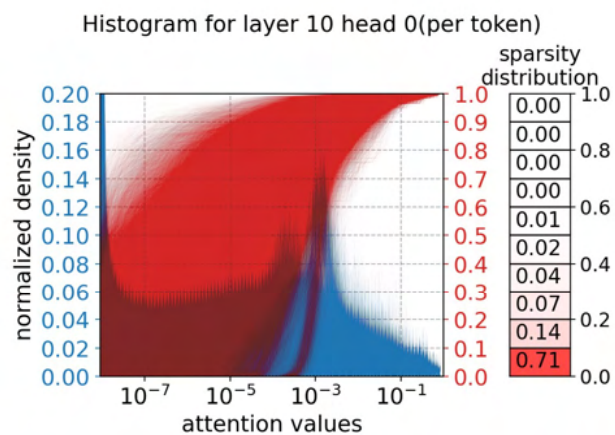
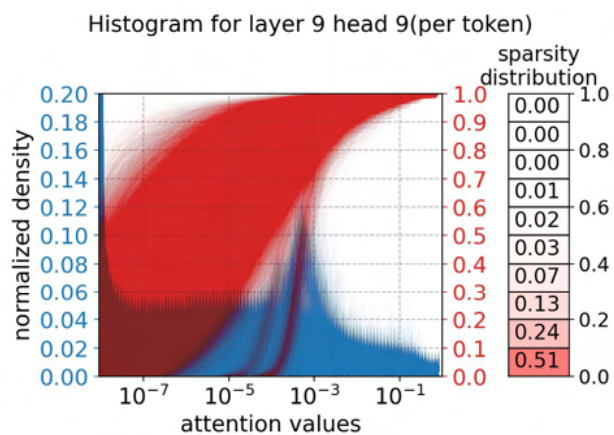
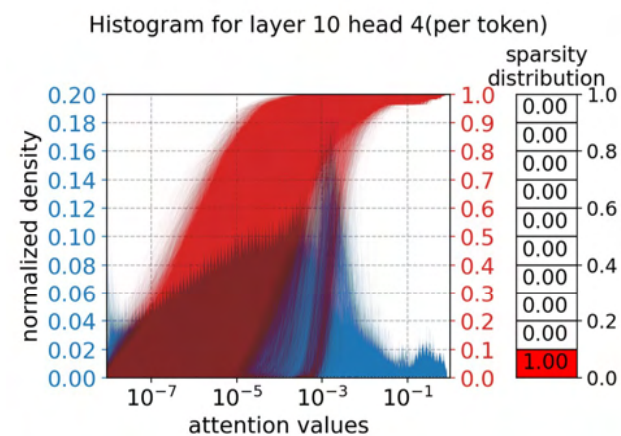
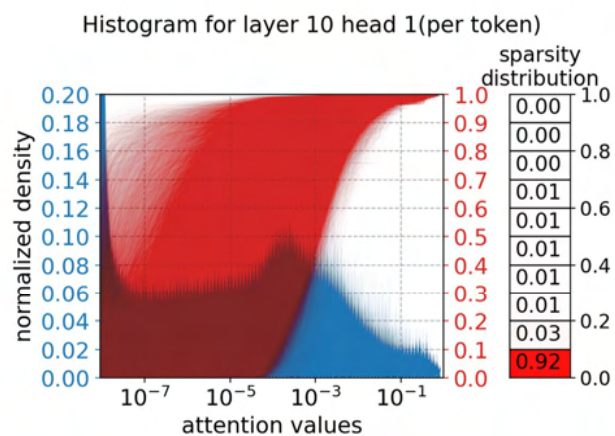
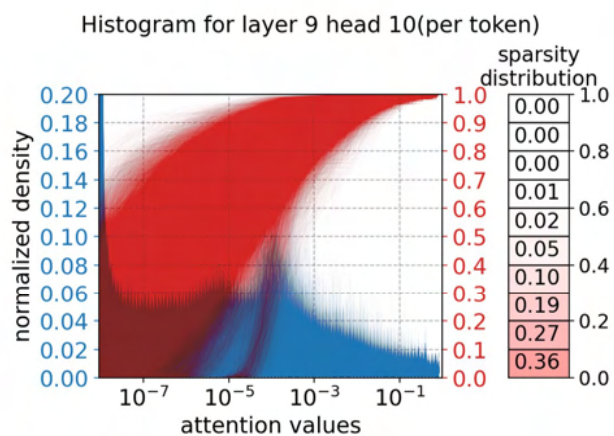
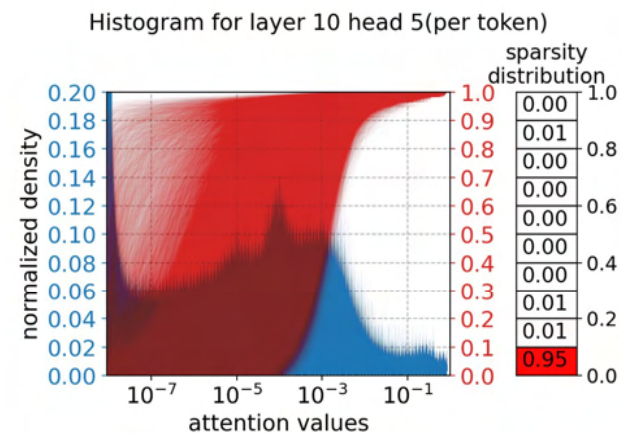
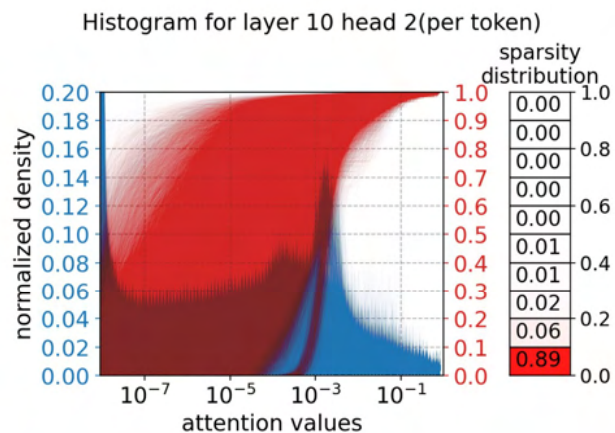
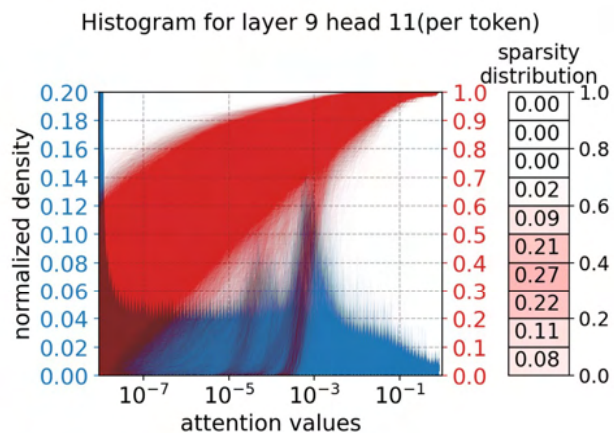


Histogram for layer 9 head 3(per token)



Histogram for layer 9 head 6(per token)





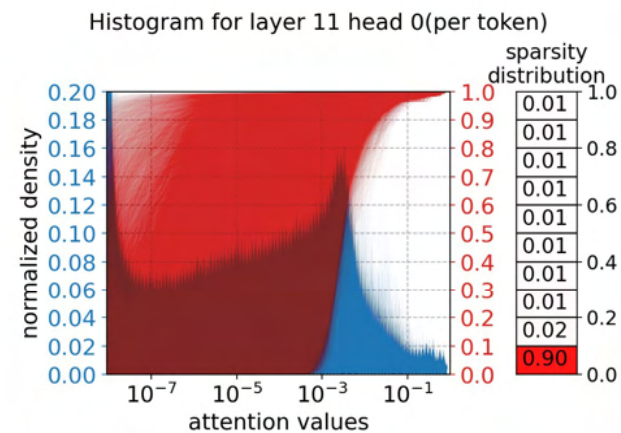
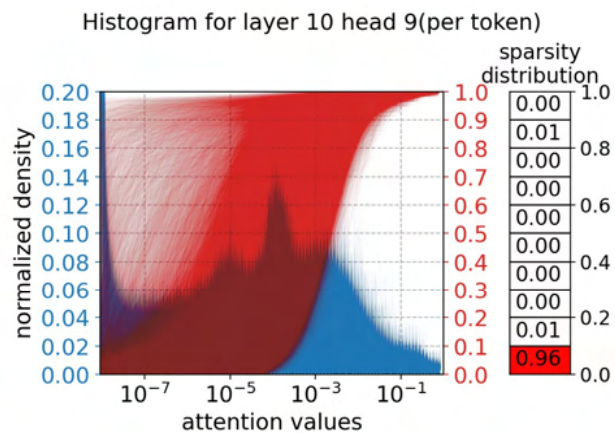
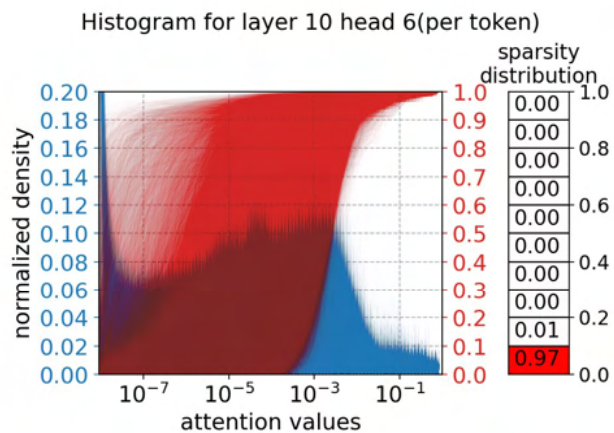
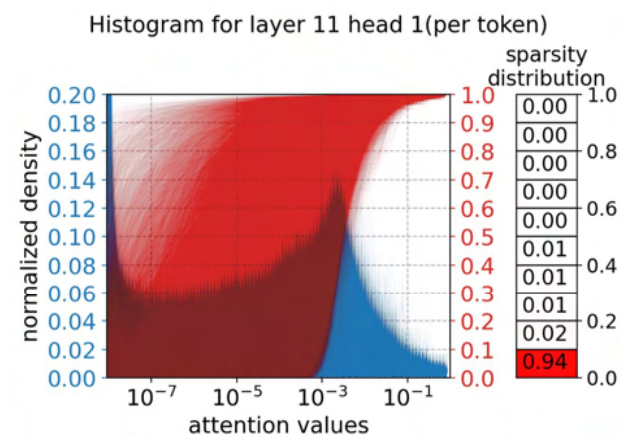
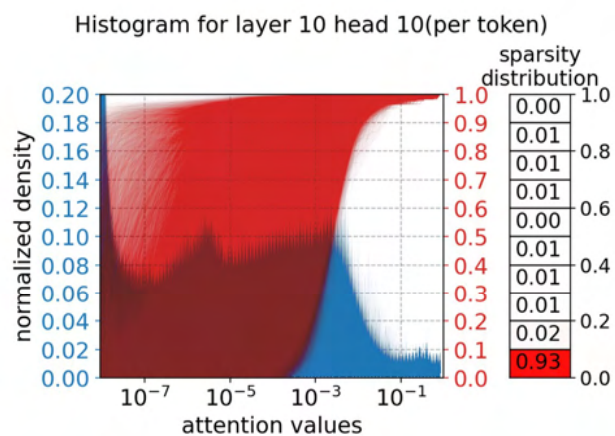
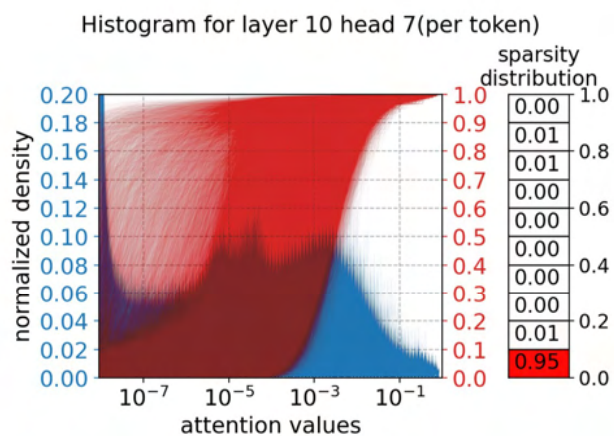
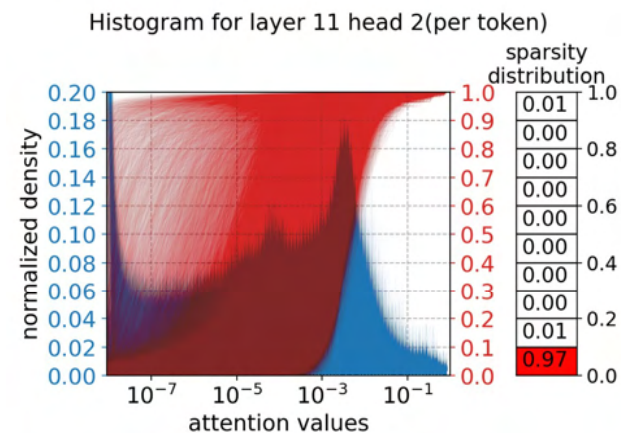
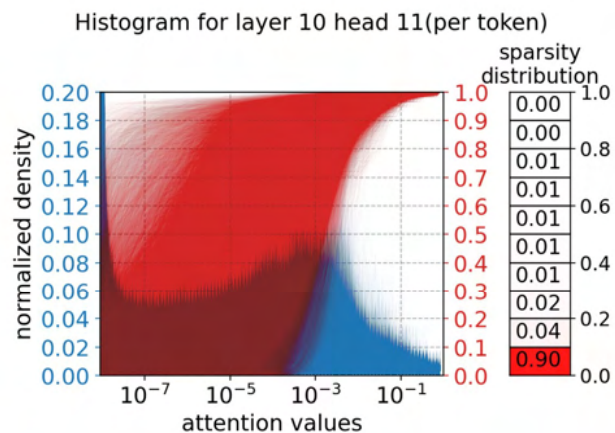
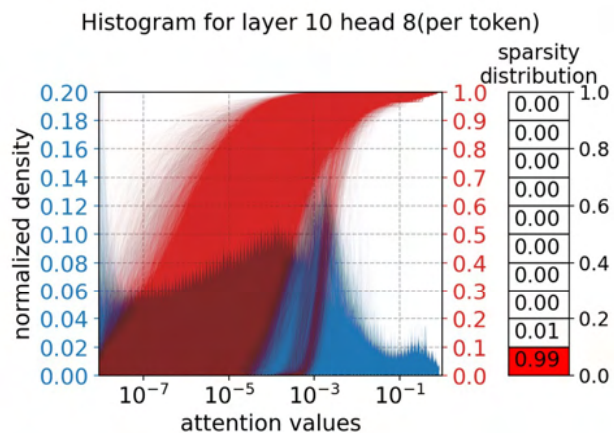


Figure 1: A 2D density plot showing the distribution of attention values (x-axis, logarithmic scale from 10^{-7} to 10^{-1}) and normalized density (y-axis, linear scale from 0.00 to 0.20). The plot is colored by sparsity, with a color bar on the right ranging from 0.0 (blue) to 1.0 (red). The distribution shows a sharp transition from high sparsity (red) to low sparsity (blue) around an attention value of 10^{-3} . A table on the right lists the sparsity distribution for different attention values, with the highest value (0.87) highlighted in red.

attention values	sparsity distribution
0.00	1.0
0.01	0.9
0.01	0.8
0.01	0.7
0.01	0.6
0.01	0.5
0.01	0.4
0.02	0.3
0.03	0.2
0.04	0.1
0.87	0.0

Figure 1 is a 2D density plot showing the distribution of attention values (x-axis, logarithmic scale from 10^{-7} to 10^{-1}) and normalized density (y-axis, linear scale from 0.00 to 0.20). The plot is color-coded by sparsity distribution, with a color bar on the right ranging from 0.0 (dark red) to 1.0 (dark blue). The plot shows a dense region of low sparsity (red) at low attention values and a sparse region of high sparsity (blue) at high attention values.

normalized density	sparsity distribution
0.04	1.0
0.03	0.9
0.02	0.8
0.02	0.7
0.03	0.6
0.03	0.5
0.03	0.4
0.04	0.3
0.06	0.2
0.71	0.1
	0.0

Figure 1 is a plot showing the normalized density of attention values for different sparsity distributions. The x-axis represents attention values on a logarithmic scale from 10^{-7} to 10^{-1} . The y-axis represents the normalized density from 0.00 to 0.20. A color bar on the right indicates the sparsity distribution from 0.0 (blue) to 1.0 (red). The plot shows that as sparsity increases, the distribution of attention values shifts towards higher values (right).

normalized density	sparsity
0.00	0.00
0.02	0.01
0.04	0.00
0.06	0.00
0.08	0.00
0.10	0.00
0.12	0.00
0.14	0.00
0.16	0.00
0.18	0.00
0.20	0.98

normalized density	sparsity distribution
0.20	0.00
0.18	0.00
0.16	0.00
0.14	0.00
0.12	0.00
0.10	0.00
0.08	0.00
0.06	0.00
0.04	0.00
0.02	0.01
0.00	0.97

Figure 1 is a heatmap illustrating the normalized density of attention values for different sparsity distributions. The x-axis represents attention values on a logarithmic scale, ranging from 10^{-7} to 10^{-1} . The y-axis represents the normalized density, ranging from 0.00 to 0.20. The color scale on the right indicates the sparsity distribution, ranging from 0.0 (blue) to 1.0 (red). The plot shows that as the sparsity distribution increases, the density of attention values shifts towards higher values (right side of the plot).