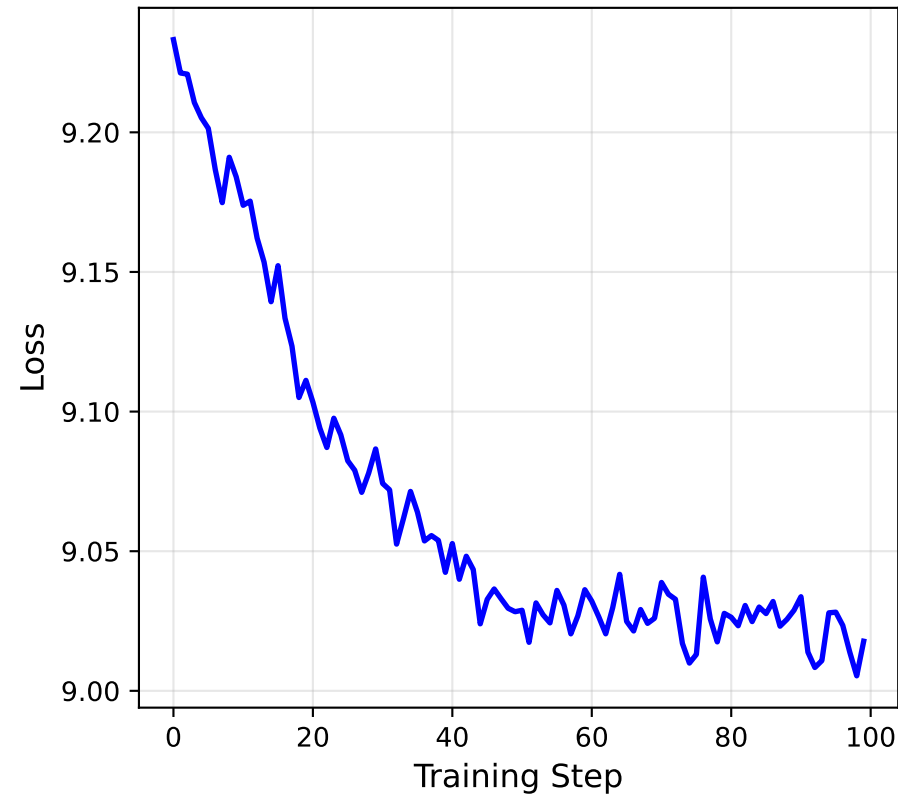
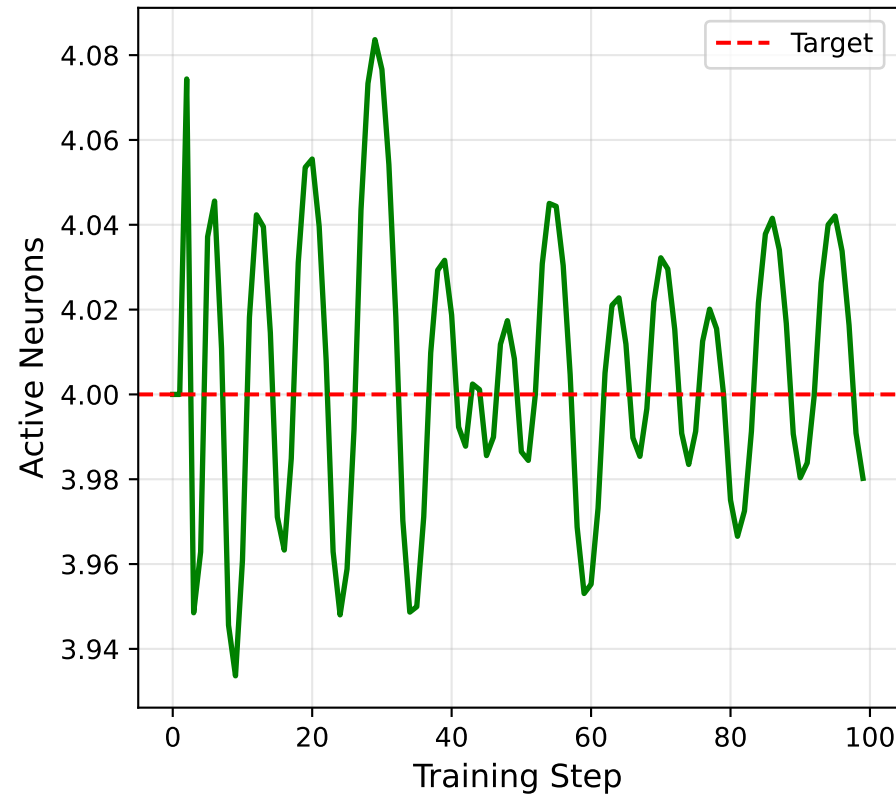


Sparse Targeted Activation Editing Results

Training Progress



Sparsity Over Time



Selected Neurons

