Token ranks in QK circuit (inc. bias) 100 Percentage of Model Vocabulary Query and Key Inputs: 80 $Q = MLP_0(W_E)$ $K = MLP_0(W_E)$ $Q = W_U$ 60 $K = MLP_0(W_F)$ 40 20

1

5

