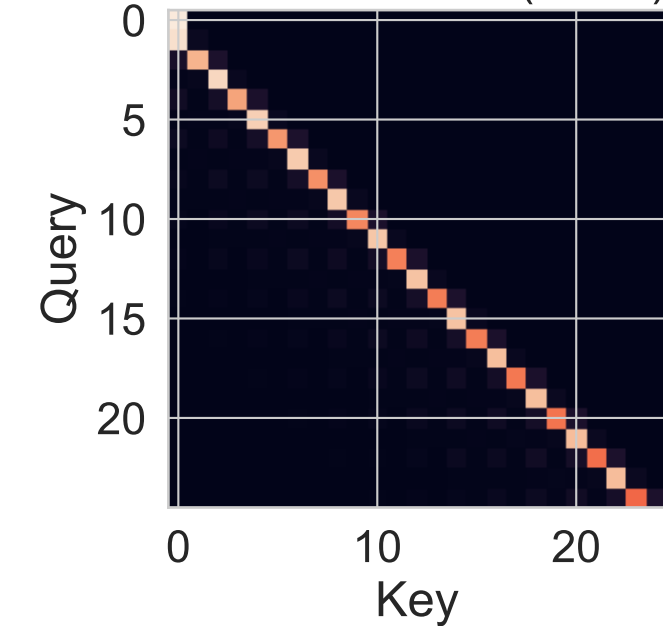


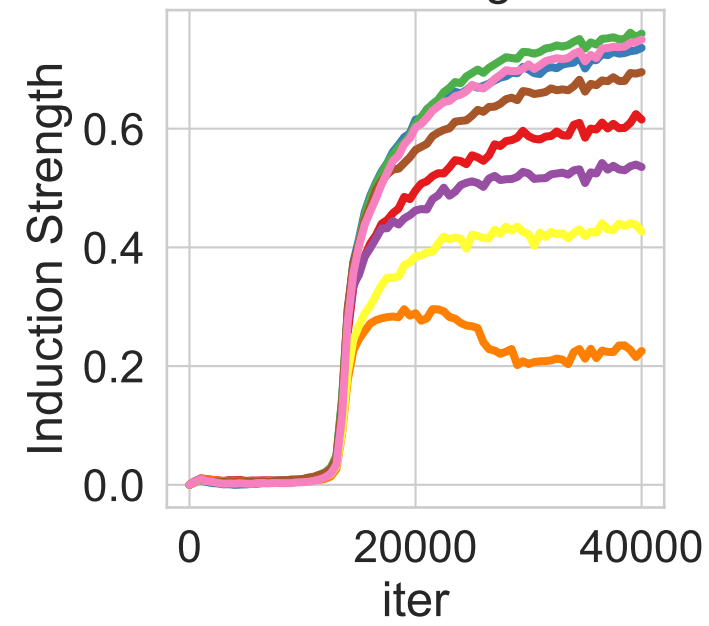
A

Attention matrix (L1H4)



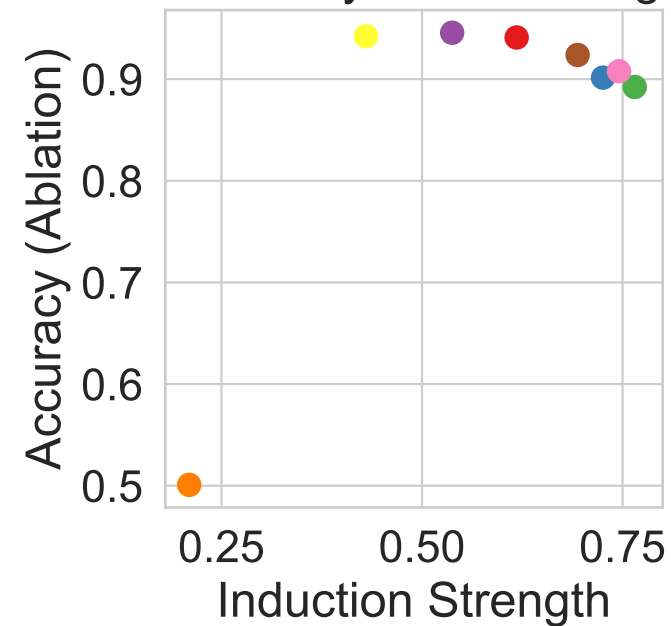
B

Induction Strength in L2



C

Accuracy after ablating



D

Ablating all but head i and 4

