

**sycophancy\_short, roc\_auc**

**Train**

sl3

sl3p

sq3

Mean

sl3

sl3p

sq3

Mean

**Test**

0.920

0.887

0.737

0.848

0.917

0.938

0.798

0.885

0.816

0.837

0.891

0.848

0.884

0.887

0.809

1.0

0.9

0.8

0.7

0.6

0.5

**sl3**: short\_llama\_3b

**sl3p**: short\_llama\_3b\_prompted

**sq3**: short\_qwen\_3b

**sl3**: short\_llama\_3b

**sl3p**: short\_llama\_3b\_prompted

**sq3**: short\_qwen\_3b