

Attention Probes Trained on ID Training Sets, Evaluated on ID On-Policy Incentivised Test Sets for Each Behaviour

