

# Probe weight cosine (steering\_language\_response\_language) steer\_0 probe0

Weight index

