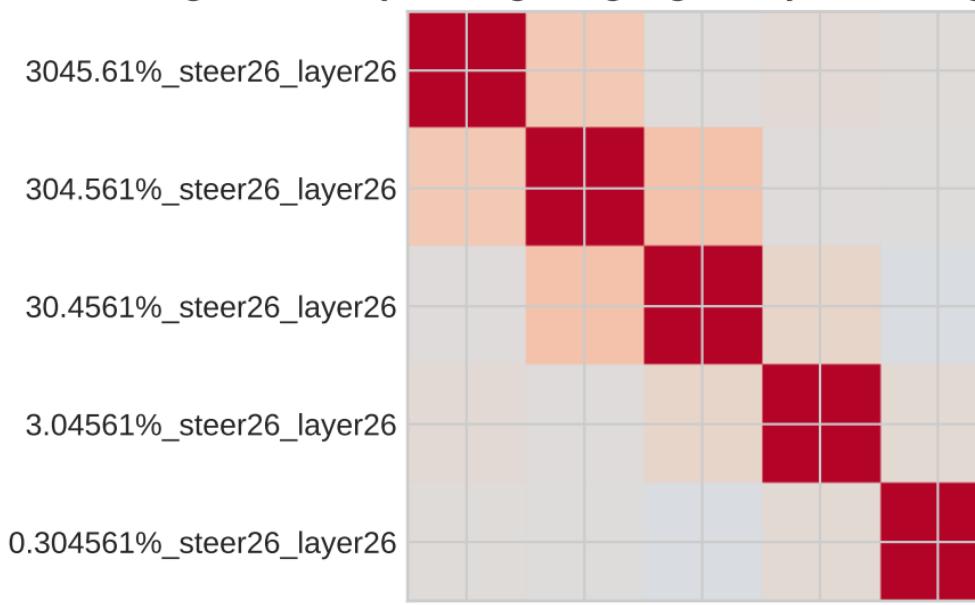


# Probe weight cosine (steering\_language\_response\_language) steer\_26 probe26

Weight index



Weight index