Probe weight cosine (steering_language_response_language) steer_18 probe18