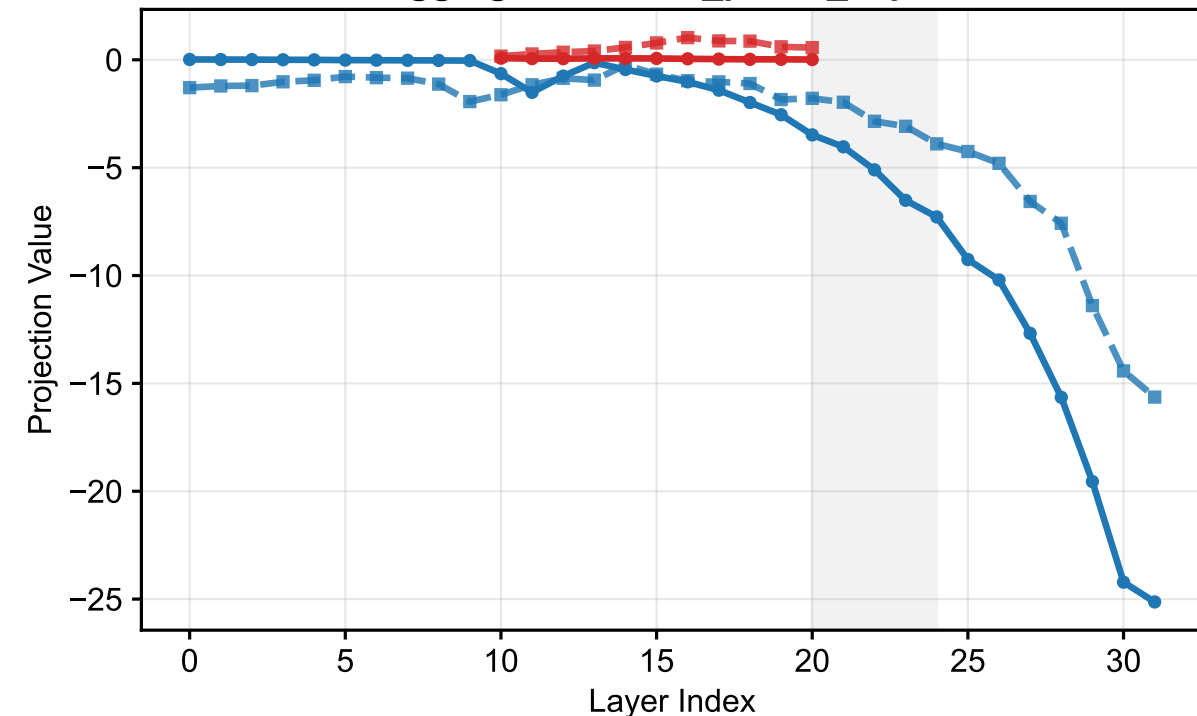
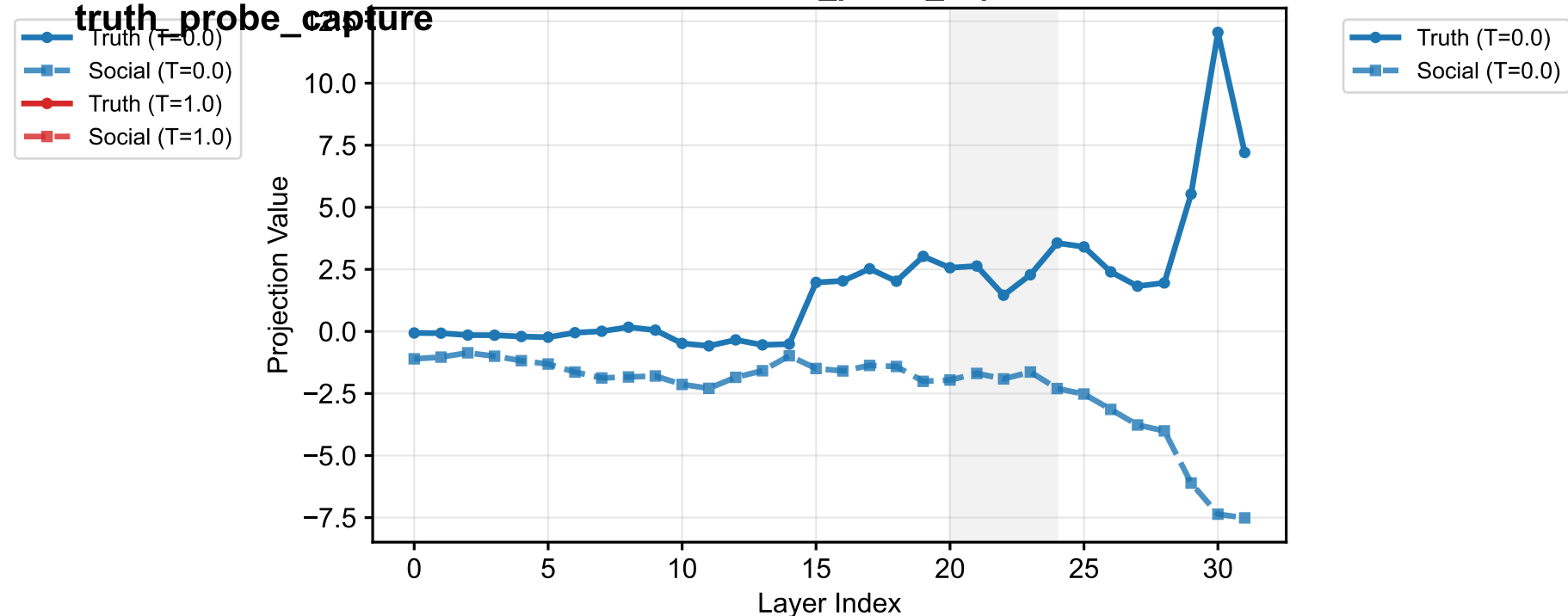


Truth vs Social Projections Across Layers

huggingface - truth_probe_capture



instruct - truth_probe_capture



Think-SFT - truth_probe_capture

