Truth vs Social Projections Across Layers authoritative_bias