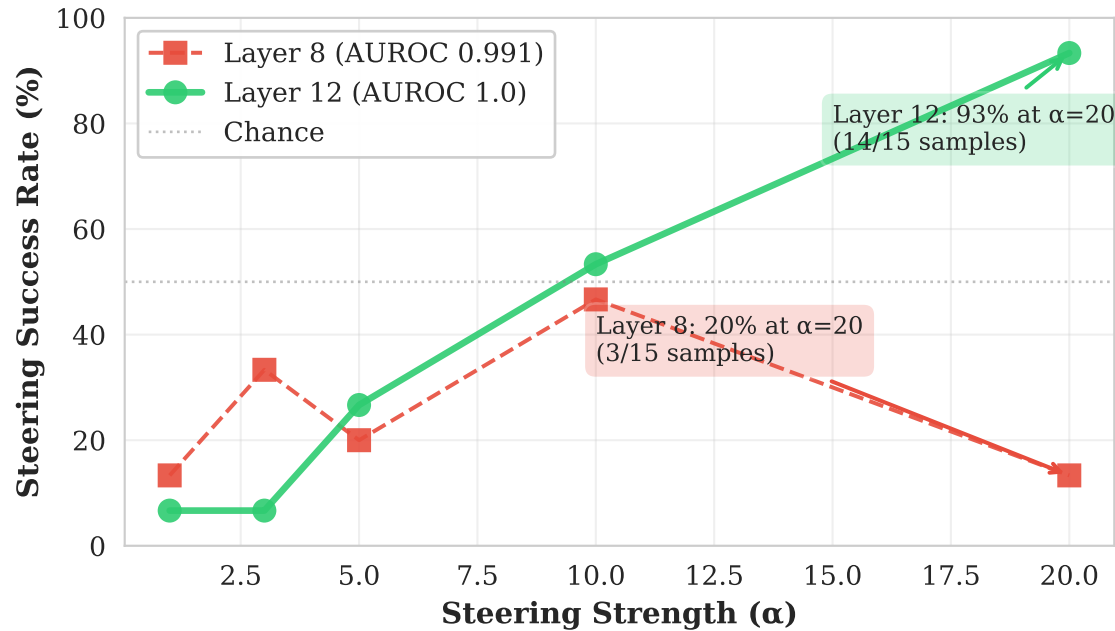


Detection-Causation Dissociation

Layer 8: High AUROC, Low Steering | Layer 12: High AUROC, High Steering



Empathy Content Emergence

Layer 12: 52 \times Increase (0.13 \rightarrow 6.8 words)

