**CoT Locality: Intervention Effects by Mode and Layer**

*Baseline Accuracy: 65.0%*
*Δ = Intervention Accuracy − Baseline Accuracy*
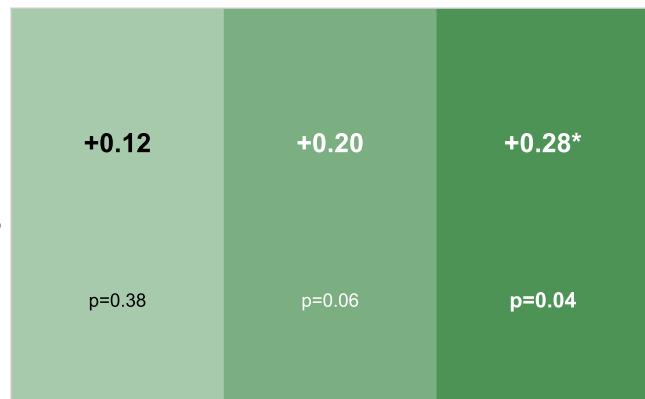*p-values: McNemar's Test | \* p<0.05, \*\* p<0.01, \*\*\* p<0.001 | Bold = Significant*

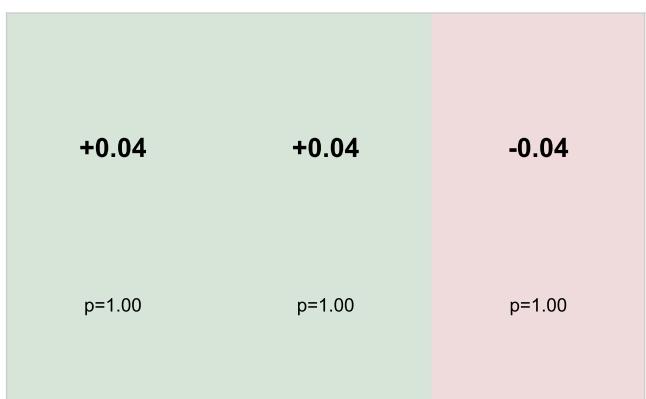|  | **Add** | | | **Lesion** | | | **Rescue** | | |
|---|---|---|---|---|---|---|---|---|---|
| **Layer 25** | +0.12 (p=0.38) | +0.20 (p=0.06) | +0.28* (p=0.04) | −0.16 (p=0.12) | −0.28* (p=0.02) | −0.24 (p=0.11) | +0.08 (p=0.50) | +0.20 (p=0.06) | +0.32** (p=0.008) |
| **Layer 26** | +0.00 (p=1.00) | +0.04 (p=1.00) | −0.08 (p=0.62) | +0.04 (p=1.00) | +0.04 (p=1.00) | −0.04 (p=1.00) | −0.04 (p=1.00) | +0.00 (p=1.00) | −0.08 (p=0.50) |
| **Layer 27** | −0.08 (p=0.62) | +0.00 (p=1.00) | +0.04 (p=1.00) | +0.00 (p=1.00) | +0.04 (p=1.00) | −0.04 (p=1.00) | +0.04 (p=1.00) | −0.04 (p=1.00) | −0.04 (p=1.00) |

α / γ / β: 0.5, 1.0, 2.0

ΔAnswer Accuracy (Intervention − Baseline)