

# Layer-wise Distribution of Causal Features (L1-30)

2,787 Features: 640 Safe, 2,147 Risky

