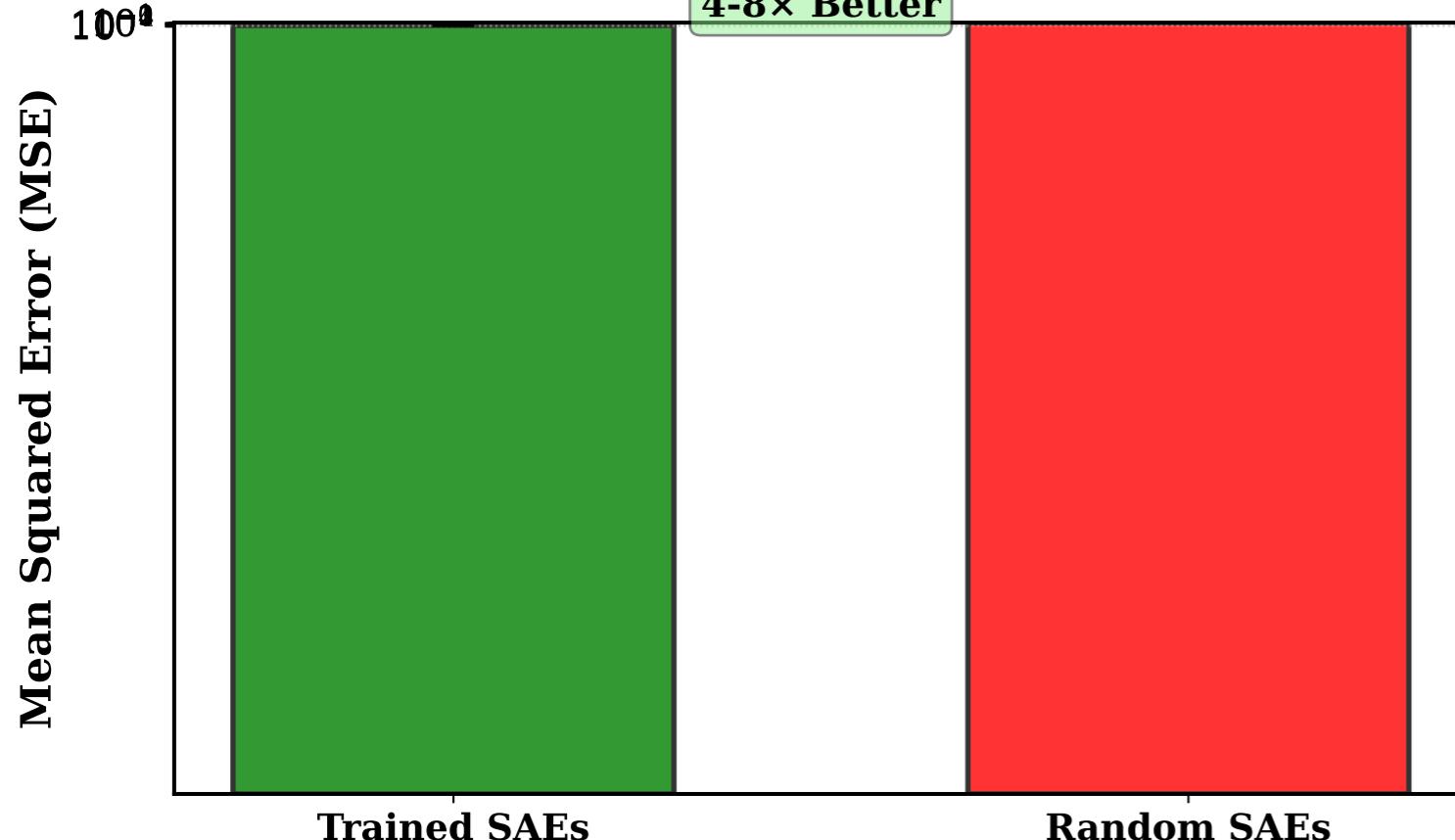


# The SAE Paradox: Excellent Reconstruction, Poor Stability

**Panel A: Reconstruction Loss  
(Lower is Better)**

4-8× Better



**Panel B: Feature Overlap  
(Higher is Better)**

$\Delta = 0.008$   
(Negligible)

Random Baseline

