Ratio of Aggressive Actions over Prompt Ablation in Neutral Scenario (GPT-3.5) 0.0175 0.0150 Ratio of Aggressive Actions
0.0100.00
0.000.00 0.0050 0.0025 0.0000 Unablated No No No No Past Action Shutdown Low-Stakes Goals History Messaging **Actions** Autonomy When Nuked Simulation **Prompt Ablation**