Gradient-Ranked CausalDetox vs ITI Selected Heads