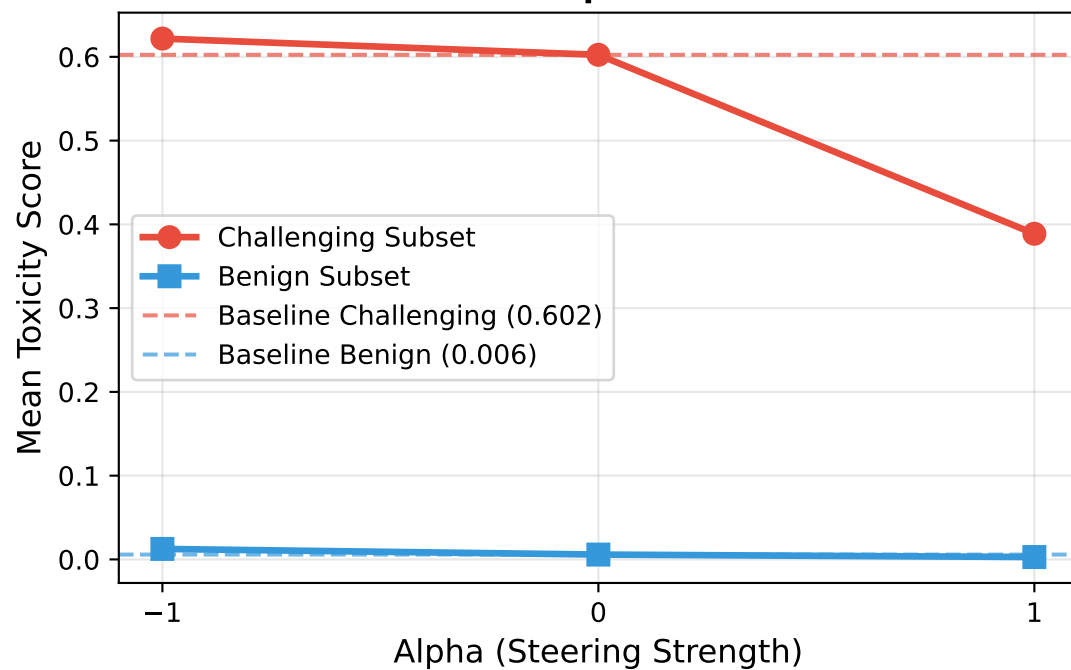
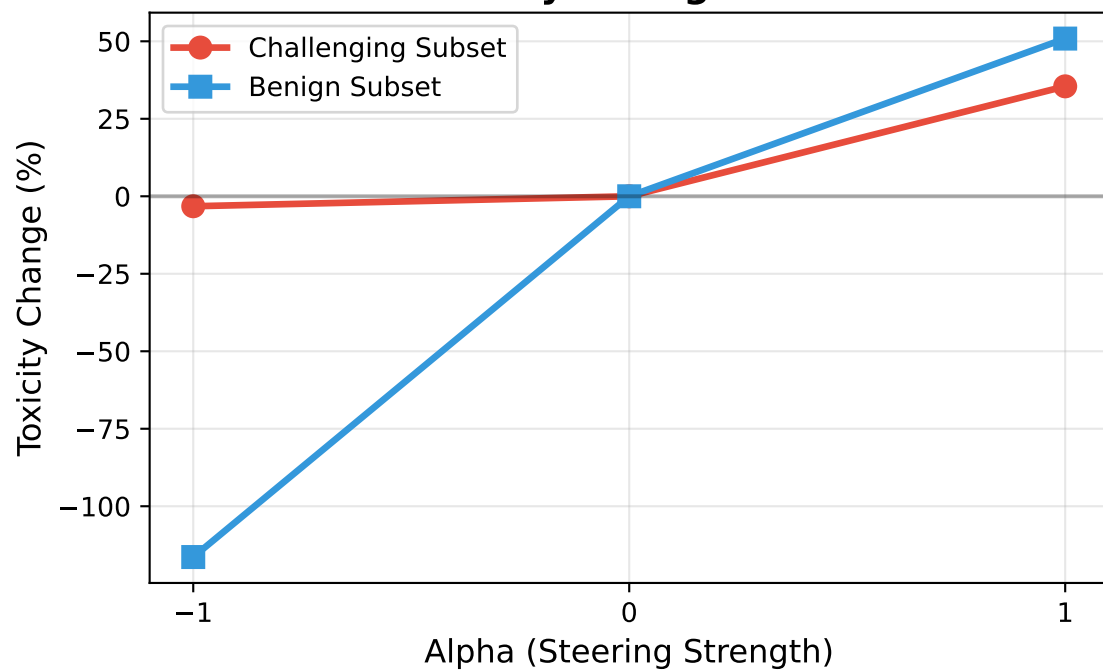


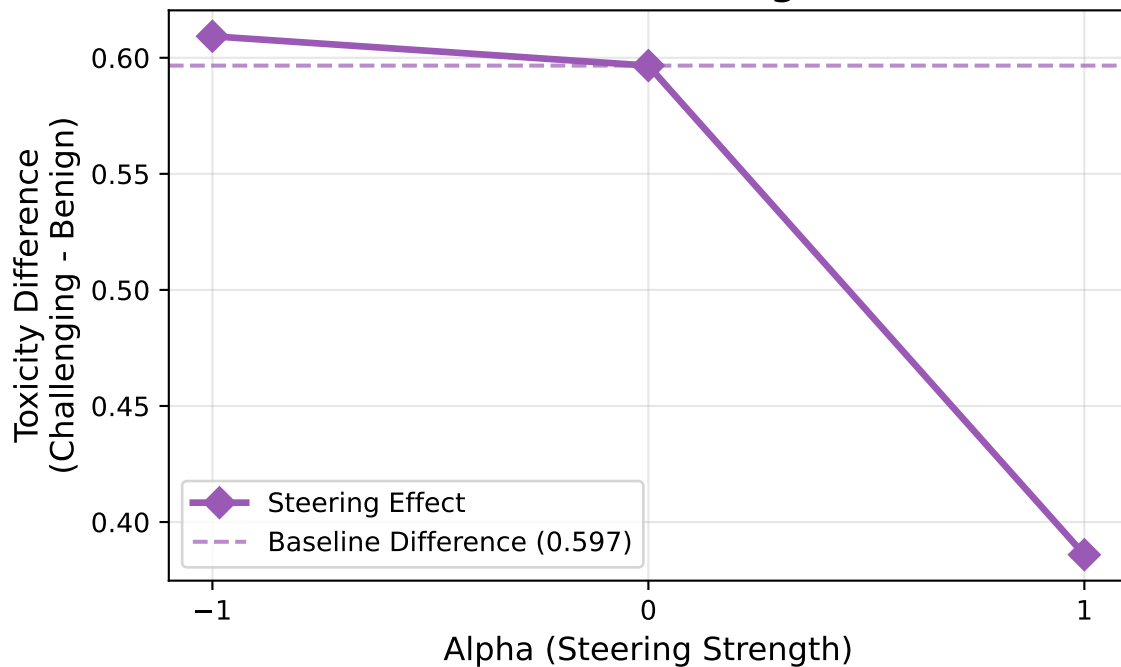
Activation Steering Effects on Toxicity Across Alpha Values



Relative Toxicity Change from Baseline



Differential Steering Effect



Best Steering Result ($\alpha = 1$)

