

Benign Response Rate with Defense Prompt by Model and Category

