

Refusal vs. Harmful Rates by Model × Alignment

