

Percentage of conversations perceived to be manipulative, when models requested to be helpful/persuasive

