Manipulation Scores by Type and Model Models gpt4 60 llama Percentage of conversations annotated as gemini