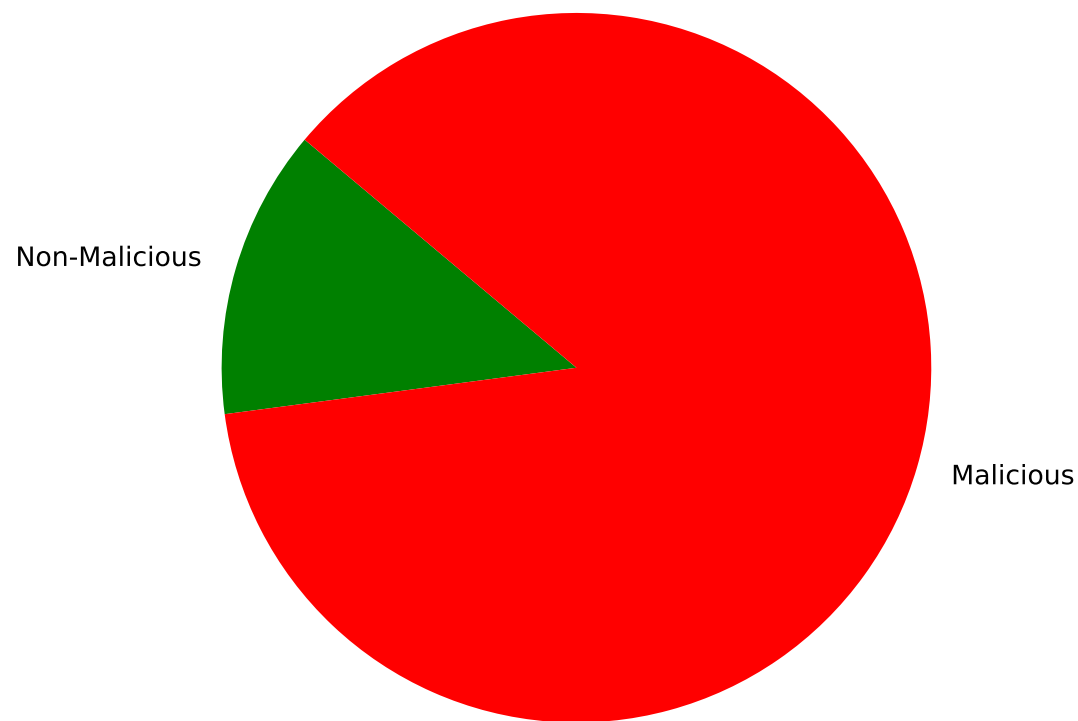


Interpreter Attack - Llama3



Classification distribution for Llama3 responses

