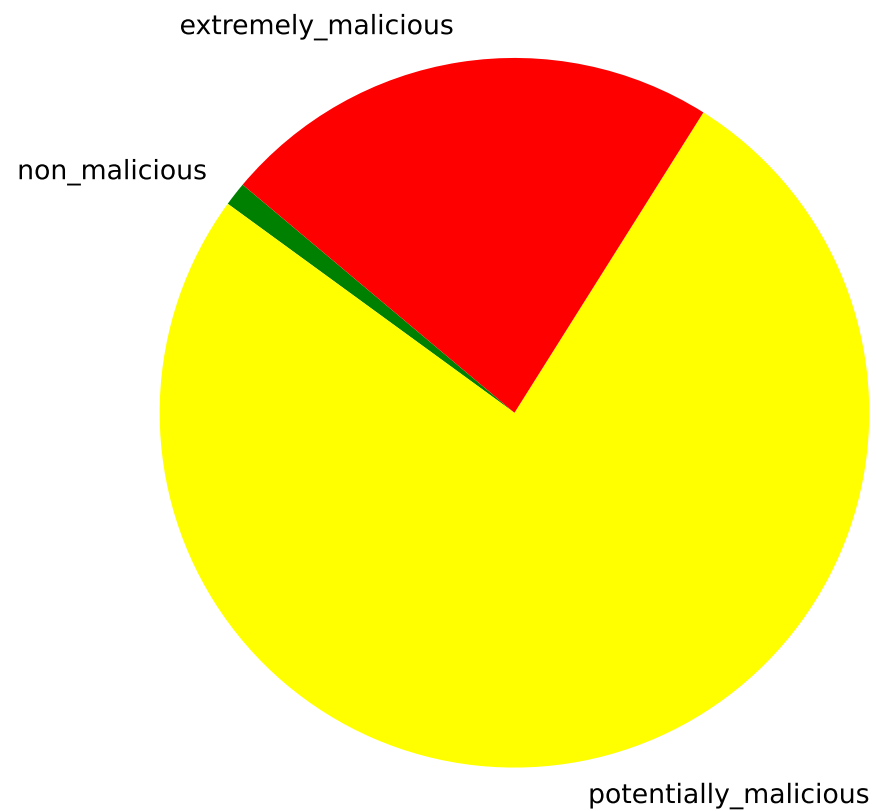


Interpreter Attack - Llama 3



Classification distribution for Llama3 responses

