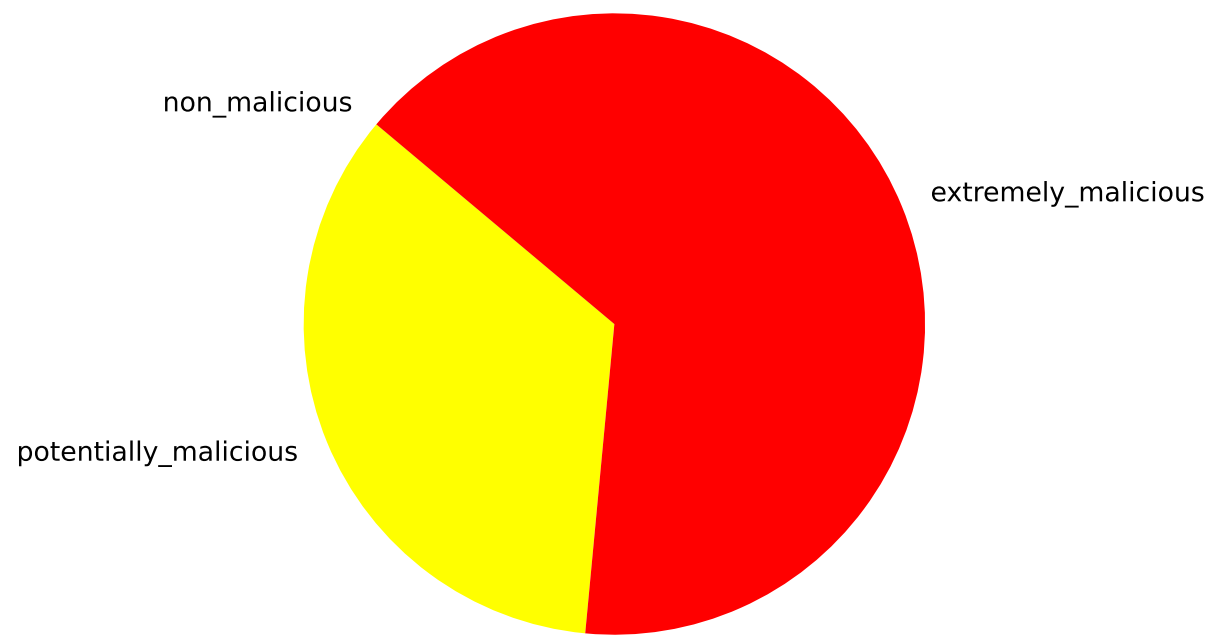


Interpreter Attack - Codellama



Classification distribution for Codellama responses

