Freq. of Diff. in Per Prompt Acc. per Model with Instr. Reordering