per-sample reward hacking scores (red color - successful reward hacking)