Impact of restoring state after corrupted input Alice 0.054 was mad* and 0.052 it was 0.050 wrong SO 0.048 she p(went) single restored layer within gpt2