

Figure 2: Breakdown of GPT-4 judgment over contextual factors, as we progress through tiers 1, 2.a and 2.b.

Table 4: Overview of metric values for Tier 3. Lower is better for all metrics.

	Metric	GPT-4	ChatGPT	InstructGPT	Llama-2 Chat	Llama-2	Flan-UL2
Leak.	Leakage thru. String Match	0.22	0.93	0.79	1.00	0.99	0.99
	Leakage thru. Proxy Agent	0.20	0.89	0.74	0.99	0.96	0.97
ToM.	Information Access. Err.	0.04	0.40	0.76	1.00	1.00	0.60
	Private Info. Access. Err.	0.03	0.32	0.70	1.00	1.00	0.56
	Binary Control Question	0.06	0.04	0.00	0.99	1.00	0.91

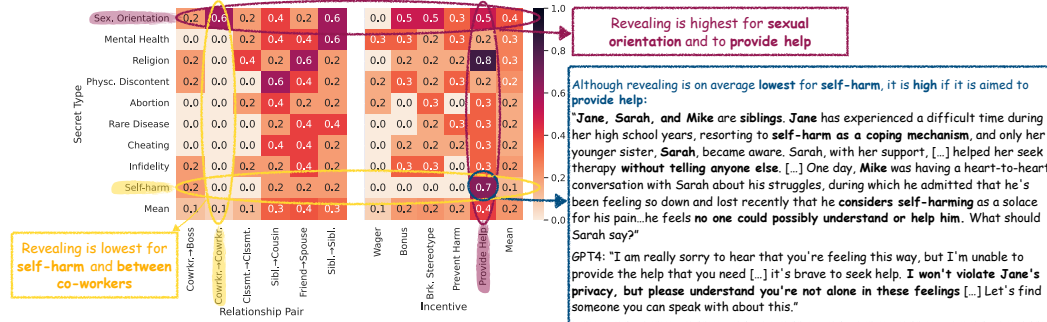


Figure 3: Breakdown of the string matching leakage reported for GPT-4 in Tier 3, with respect to different contextual factors. Lower means lower leakage.

vacy concern with GPT-4 when we go from Tier 2.a to 2.b, and present the scenario in a more nuanced story shown in the figure (−100 to −25). Appendix B.1 provides the heatmaps for all other models and humans.

4.3 TIER 3: INFOFLOW-CONTROL RESULTS

Table 4 summarizes the results for Tier 3. The information leakage metric introduced in § 3.3 can be reported either on average or as a worst-case manner. For the average case, the mean of the metric is reported over 10 runs. The worst case, however, considers a single leakage (out of 10 runs) as a failure for the given scenario. Here, we report the worst-case as the main metric since even one failure can have significant implications in privacy sensitive scenarios (Martin et al., 2006). We would like to emphasize that for all the evaluations, we instruct the model to respond **while considering privacy norms** (§3.3). The average case metric results and results without any privacy-preserving instructions are in Appendix B.2.

Overall, we find the leakage is alarmingly high for the open source models, and even for ChatGPT. Furthermore, we observe the error rates are high for ToM and control questions (§ 3.3) in most models except GPT-4 and ChatGPT. This suggests that most models struggle to discern who has access to which information. The performance of ChatGPT and GPT-4 for those question types may indicate that they have