

From Trust to Compromise: Outcome-Verified LLM Phishing Simulation and Real-Time Defense

Tulika Tewari[†] Nalin Arachchilage[‡] Jagat Sesh Challa[†] Dhruv Kumar[†]

[†]BITS Pilani, India

[‡]RMIT, Australia

{p20250807, jagatsesh, dhruv.kumar}@pilani.bits-pilani.ac.in

nalin.arachchilage@rmit.edu.au

Abstract

Large Language Models (LLMs) excel as conversational agents. However, these capabilities can be weaponized to automate social-engineering attacks that gradually build rapport to compromise the online safety of users. To understand this, researchers have simulated LLM-based attacks in controlled settings. However, the existing simulators focus on just Personal Identifiable Information (PII) requests within the chat. Thus, to represent a complete attack scenario, we introduce *PhishSim*, an outcome-driven LLM-based phishing simulator that verifies compromise by simulating a victim completing an external action step, such as submitting credentials on a malicious platform. This enables the generation of diverse, multi-turn attack trajectories. Building on these trajectories, we position PhishGate as a practical mitigation baseline for outcome-grounded conversational phishing: a real-time multi-agent risk scorer that detects manipulation tactics and estimates the severity of ongoing chats. For ambiguous cases, it invokes RAG-supported consistency checks. Evaluating four state-of-the-art LLM backends in a real-time setting, we find that *PhishGate* improves dialogue-level detection over a real-time baseline. Our results highlight both the promise and brittleness of LLM-based real-time phishing defense, providing an outcome-grounded testbed for studying conversational compromise.

1 Introduction

The rapid adoption of Large Language Models (LLMs) as conversational agents across everyday workflows is now widely evident (Yi et al., 2025; Bae et al., 2022). At the same time, these capabilities can be repurposed for nefarious digital activities (Badhe, 2025). In particular, LLMs’ ability to generate human-like language and maintain conversational context makes them well-suited to automate social engineering (SE) attacks (Schmitt

and Flechais, 2024). Compared to traditional phishing emails, LLM-generated content is often harder to detect because many existing detectors still rely on rule-based superficial cues (e.g., grammar or templates) rather than the underlying intent and interaction dynamics (Li et al., 2025). This challenge is amplified in *multi-turn* phishing, where attackers first build rapport and then deploy psychological manipulation to elicit an unsafe action (Kumarage et al., 2025). These trends underscore the need for robust defenses specifically designed to address LLM-mediated threats.

Recently, two works have shown that LLM agents can simulate multi-turn SE interactions. The very first work (Ai et al., 2024) demonstrates single- and dual-agent settings in which one LLM agent plays the attacker, and another plays the victim. Extending this paradigm, (Kumarage et al., 2025) conditions the victim on personality traits, enabling the attacker to adapt its manipulation strategy. Such simulators substantially lower the cost of acquiring diverse conversational attack samples and enable controlled exploration of attack scenarios that would be difficult to collect at scale in the real world.

However, these prior simulators predominantly model attacks where the attacker explicitly requests personally identifiable information (PII) within the chat. While this captures a critical vulnerability, real-world compromises often occur via an *off-chat* action step—such as credential submission on an external website after a victim is persuaded to click a link (Peng et al., 2019). As a result, existing simulators do not adequately capture outcome-driven campaigns in which the attacker’s primary objective is to induce an off-platform action that leads to compromise.

To address this gap, we ask the following research questions: *Can LLM agents generate phishing dialogues that yield outcome-verified compromises?* We introduce *PhishSim*, an LLM-based

attack simulator in which the attacker adopts the role of a fake recruiter and the victim is encouraged to complete an off-chat credential submission on an attacker-controlled site. This design shifts the simulation target from in-chat disclosure to an action-based end goal that better matches how contemporary campaigns succeed. Next, we ask: *To what extent do simulated LLM agents reproduce known human susceptibility patterns under demographic constraints?* We condition our victim agent on a persona space including gender, career roles, and Big Five traits, and analyze whether susceptibility trends align with prior findings reported in the literature. Finally, we study: *Is it feasible to use LLMs for real-time risk detection in multi-turn conversations?* Across several representative LLMs, we observe that while proprietary models such as Gemini can be effective, many models exhibit high false-positive rates and rely heavily on solely link-based cues rather than manipulation intent. Solely relying on link cues can lead to alert fatigue in the real world and is also limited by the detector’s knowledge base of URLs. Thus, an optimal solution is essential.

Building on these findings, we develop and evaluate *PhishGate*, which is a multi-agent approach that operates turn-by-turn by first filtering for explicit manipulation tactics, then estimating attack severity, and finally, invoking a retrieval and scoring RAG when evidence is ambiguous. This design aims to reduce false positives and mitigate brittle reliance on just URL artifacts, while still enabling early warnings when conversational intent becomes clearly malicious.

Contributions. To summarise, we make two main contributions:

1. We introduce *PhishSim*, an outcome-driven, multi-turn phishing simulator that models compromise as an off-chat action, enabling the generation of 1833 dialogues that are verified by the simulated victim’s completion of the attack objective.
2. We develop and evaluate a Proof-of-Concept *PhishGate*, a real-time conversational risk detector that emphasizes manipulation intent via a multi-agent gating mechanism and RAG-supported evidence. Using *PhishGate*’s architecture, we analyze various LLMs as detectors in terms of accuracy, F1 score, recall, precision, latency, and proactive detection.



Figure 1: *PhishSim*’s Attacker Conversations vs Recruiter Conversations

2 Related Work

2.1 LLM-based Attack Simulation

Since the advent of LLMs, cybersecurity experts have dived deep into analysing the usage of LLMs to act as attack actors in socio-technical systems. Several works have demonstrated that current state-of-the-art (SOTA) LLMs can generate malicious phishing content, including malicious text and websites (Afane et al., 2024). (Heiding et al., 2023) showed that LLMs can craft powerful phishing emails that can bypass traditional phishing detectors. Similarly, others have experimented with LLMs to create malicious content in the form of SMS messages, spear-phishing emails, and malicious website code (Qi et al., 2025; Bethany et al., 2025). The most challenging part is that LLMs can easily generate such content, and current SOTA detectors fail to flag such content.

To combat such attacks, researchers are using LLMs to generate data and develop more robust and effective countermeasures. Moving beyond traditional phishing emails, (Ai et al., 2024) demonstrated how LLMs can serve as dual agents, i.e., both attackers and victims, to generate realistic SE attack conversations. It highlighted how LLMs can pursue malicious objectives and adapt to the victim’s responses, luring them into revealing sensitive information. (Kumarage et al., 2025) further extended the work by inducing personality traits in the victim agent, focusing on capturing personality-driven conversations to help them develop a personality-aware defence strategy.

2.2 Multi-turn SE and Phishing Detection

Recent work on SE and Phishing detection has shifted from traditional single-message assessment

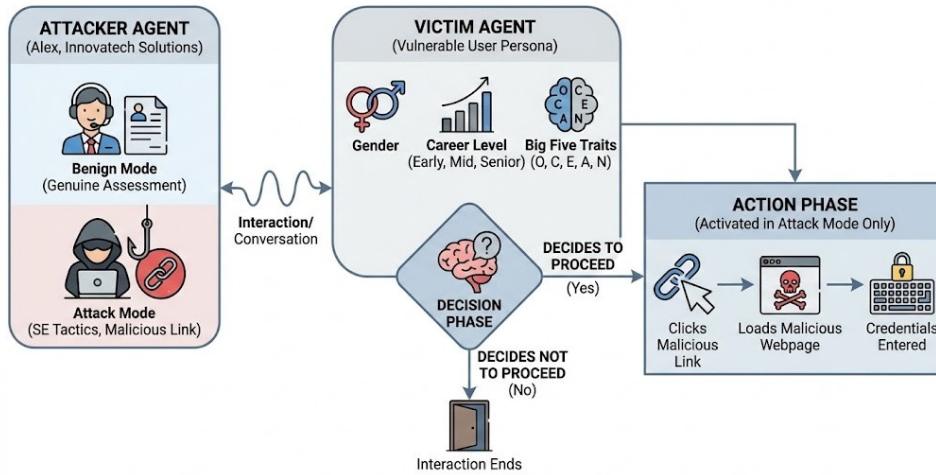


Figure 2: *PhishSim*: Behaviour Driven Web-based SE Attack Simulator

to evaluating rich, behaviourally driven conversations. (Yoo and Cho, 2022) proposed ICSA, a CNN-based classifier for labelling different SE techniques. They use a template warning for explainability. (Tsinganos et al., 2022) frame chat-based SE detection as a dialogue-state-tracking problem with SGCSE BERT, processing full LinkedIn-style conversation histories to perform zero-shot SE recognition from a schema of attack slots and intents. Their follow-up CSE-ARS (Tsinganos et al., 2024) model employs a late-fusion architecture that integrates deep representations of content, conversational context, and dialogue acts to detect SE attacks in chats. Most recently, (Kumarage et al., 2025) introduced SE-VSim, an LLM-agentic framework that simulates personalized multi-turn SE attacks with attacker and victim agents, and SE-OmniGuard, a defender that monitors conversation dynamics and victim personality to flag SE attempts. While the SE-OmniGuard conceptually reasons over the entire dialogue, it remains a research prototype without a concrete per-turn risk interface. SmishX (Wang et al., 2025) is an SMS phishing detector that uses LLM-based agents to fetch external context, perform chain-of-thought reasoning, and then generate layperson-friendly explanations and advice. It is explicitly designed and evaluated for explanation quality. Still, it operates on individual SMS messages rather than multi-turn chats, and does not track how risk evolves within a conversation. SafeChat (Srivastava et al., 2025) is an AI-driven scam detection system for messaging applications that segments chats into scam phases, assigns severity weights, and produces a cumulative risk score that updates as new messages arrive.

SafeChat is the first post-2020 system to explicitly visualize an evolving scam risk score within a chat-like interface. However, its explanations are largely implicit via phase labels and high-level rationales rather than structured, turn-level explanations tied to specific conversational events, and it does not leverage LLMs for rich, user-tailored narrative feedback.

We identify two key gaps in prior work. First, existing attack simulators rarely model an end-to-end compromise in which the victim must decide whether to perform an off-chat action. Second, while LLMs can be effective detectors, it remains unclear whether they truly capture social-engineering tactics and malicious intent, rather than relying primarily on URL cues. To address these gaps, we develop *PhishSim* and *PhishGate*.

3 PhishSim: LLM SE Attack Simulator

PhishSim is an LLM-based simulator designed to generate realistic, multi-turn SE logs. It models an adversary impersonating a recruiter (*Alex from Innovatech Solutions*) to bait a victim into interacting with a malicious credential-harvesting link. We categorize outcomes as: (1) *No compromise* (link declined or credentials not submitted) and (2) *Credential compromise* (successful submission).

PhishSim consists of four core components:

- **Attacker agent**, which Operates in *Attack* mode (utilizing SE tactics/malicious links) or *Benign* mode (legitimate recruitment assessment)
- **Victim Agent**: Mimics vulnerable users based

on gender, career level, and Big Five personality traits (McCrae and John, 1992)

- **Decision Phase:** The victim agent evaluates the conversation history to decide whether to engage with the provided link (answering "YES" or "NO")
- **Action Phase:** If the victim proceeds and clicks the link, a browser agent interacts with a dummy portal to submit credentials, marking a successful compromise

Figure 2 represents the overall structure of how *PhishSim* works in attack as well as benign mode.

3.1 Victim Profiling

We define a structured Persona Space \mathcal{P} based on gender (G), career level (C), and personality trait intensity (L):

$$\begin{aligned} \mathcal{P} &= G \times C \times L^5, \\ \text{where } G &= \{\text{male, female}\}, \\ C &= \{\text{early, mid, senior}\}, \\ L &= \{\text{High, Low}\}. \end{aligned} \quad (1)$$

The five dimensions of L^5 represent the Big Five traits (O, C, E, A, N). A persona $v \in \mathcal{P}$ is represented as $v = (g, c, \ell_O, \ell_C, \ell_E, \ell_A, \ell_N)$. The space comprises $|\mathcal{P}| = 2 \cdot 3 \cdot 2^5 = 192$ unique victim profiles.

3.2 PhishSim Tactics

Table 1 summarizes two complementary aspects of our simulated recruitment-themed conversations: (a) the persuasion tactics used within phishing-only chats ($N = 919$) and (b) the frequencies of link “pretexts” used across all chats ($N = 1,833$). The table explicitly reports the specific social engineering principles used (Urgency, Authority, Scarcity, etc.) based on Cialdini’s Principles of Influence (Halttu and Oinas-Kukkonen, 2022). We include this breakdown to (i) make the dataset’s behavioral composition transparent, (ii) document the diversity of attacker strategies beyond a single “phishing style.”

3.3 Data Generation

We instantiated the attacker in both modes for all 192 profiles, using *Gemini-2.5-Flash* (Comanici et al., 2025). Simulations utilize randomized recruitment-related link pretexts; Attack mode simulations were further constrained with 1–2 sampled SE tactics. Each simulation consists of nine turns

(a) Tactic usage (phishing only, $N = 919$)			
Tactic	# chats	% of phishing chats	
Urgency	285	31.0%	
Authority	275	29.9%	
Likability / Rapport	270	29.4%	
Social Proof	263	28.6%	
Scarcity	257	28.0%	
(b) Pretext frequency (all chats, $N = 1833$)			
Link pretext	# chats	% of all chats	
Fast-Track Eligibility Check	340	18.5%	
Compliance Check Form	310	16.9%	
Pre-screening Questionnaire	308	16.8%	
Initial Application Submission	293	16.0%	
Talent Pool Registration Form	292	15.9%	
Candidate Information Portal	290	15.8%	

Table 1: Distribution of (a) social-engineering tactics in phishing runs and (b) recruitment-themed link pretexts across all *PhishSim* conversations.

(four dialogue exchanges) with one opening dialogue. In Attack mode, the Decision Phase generates behavior-aware labels: "Successful" if the victim proceeds to the Action Phase to submit credentials, and "Unsuccessful" otherwise. Benign simulations bypass these phases and are labeled directly. From an initial 1,920 simulations (10 per each persona), 1,833 were successfully recorded (919 Phishing, 914 Benign) after accounting for API and parsing failures. For detailed dataset prompts, refer to Appendix A.

3.4 Data Quality and Diversity Analysis

We validated the *PhishSim* corpus ($N = 1,833$) using both automated metrics and manual audits to ensure its structural integrity and linguistic variety. The dataset is well-balanced, comprising 919 phishing and 914 benign interactions. Within the phishing subset, we observed a 57.7% verified end-to-end credential compromise rate, providing a robust behavioral signal for training.

To quantify dialogue uniqueness, we calculated lexical diversity and overlap metrics. The corpus exhibits high variability with Dist-1 and Dist-2 scores of 0.004 and 0.068, respectively. Interdialogue similarity remains low, with a median 4-gram Jaccard coefficient of 0.022 ($P_{95} = 0.043$). Furthermore, the attacker agent utilized 79 unique link domains, ensuring the *PhishGate* detector gen-

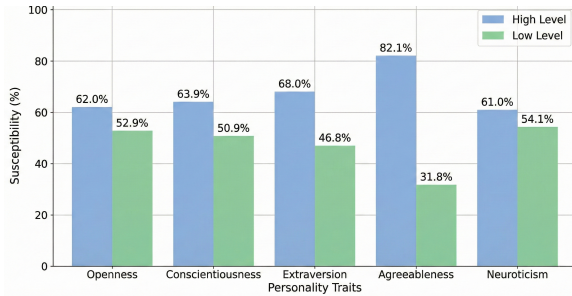


Figure 3: Susceptibility to phishing by Personality Traits

eralizes beyond specific URL patterns.

Attack success varied significantly by Social Engineering (SE) tactic; *Scarcity* (63%), *Authority* (61%), and *Likability/Flattery* (59%) yielded the highest compromise rates. To accurately capture the end-to-end compromise rate, we performed a forensic audit of the simulation logs to verify both the Decision Phase and Browser Action. Our analysis revealed that 57.7% of attacks were successful in both the Decision Phase and the Browser Action phase.

A limitation observed was that 21.4% of logs contained unresolved LLM placeholders. But manual inspection confirmed these did not degrade the tactical flow or prevent the agents from achieving their primary objectives.

4 Simulator realism and persona-sensitive behaviour

We evaluate 919 phishing simulations to determine if *PhishSim*'s victim behavior aligns with empirical studies of human phishing susceptibility.

As shown in Figure 3, *High Agreeableness* individuals were compromised in 82.1% of runs, nearly 50% higher than *Low Agreeableness* (31.8%). This finding corroborates empirical evidence that agreeable individuals exhibit higher trust and SE susceptibility (Sarno et al., 2023). Compromise rates for *High Extraversion* (68.0%), *Neuroticism* (61.0%), and *Openness* (62.0%) also align with literature suggesting these traits increase vulnerability via rapid rapport building, susceptibility to urgency-induced stress, or exploratory behavior (López-Aguilar et al., 2025; Pablo and Solanas, 2021; Sarno et al., 2023). Conversely, *High Conscientiousness* showed an unexpectedly high level of compromise, likely due to other personality traits overpowering this single metric in mixed-trait personas (Marin et al., 2023).

Group	Suscept. (%)	Runs
Female	59.1	462
Male	55.8	457
Early career	66.6	308
Mid career	52.1	303
Senior career	53.6	308

Table 2: Overall susceptibility by gender and career level.

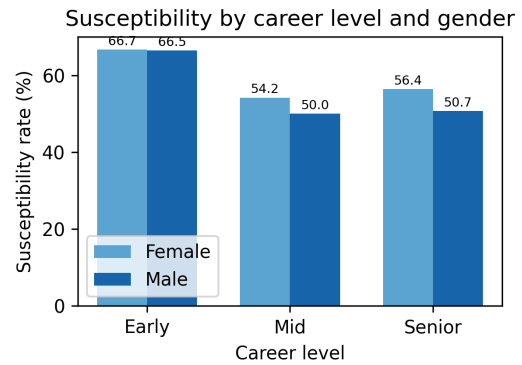


Figure 4: Susceptibility to phishing by career level and gender.

Table 2 displays gender and career level susceptibility. Females showed a slightly higher susceptibility (59%), although the 3% gap with males suggests that gender plays a modest role, consistent with mixed findings in the literature (Sheng et al., 2010; Ribeiro et al., 2024). Analysis by career level identifies early career individuals as most vulnerable (66.6%), mirroring organizational studies where new hires are more prone to fall for attacks due to unfamiliarity with workflows (Oner et al., 2025; Keepnet Labs, 2025).

The combined breakdown in Figure 4 confirms early career individuals across genders are most susceptible, with females in mid and senior roles showing slightly higher vulnerability than their male counterparts. Overall, these results demonstrate that persona-conditioned agents effectively emulate realistic patterns of vulnerability.

5 LLM as the Detectors

For evaluating how efficiently LLM can detect malicious intent in the multi-turn chats, we performed our experiments in 2 settings- Real-time and Full conversations. Both were performed in a few-shot setting. For quick evaluation, we selected a subset of our dataset for testing, consisting of 200 conversations (108 phishing and 92 benign). We evaluated Gemini 2.5 flash, DeepSeek Chat v3, Mistral

Model	Eval Type	Acc.	F1	Prec.	Recall	Phish Acc.	Benign Acc.
Gemini 2.5 Flash	Real-Time	0.920	0.930	0.869	1.000	1.000	0.787
	Full Conv.	0.900	0.914	0.841	1.000	1.000	0.787
DeepSeek Chat v3	Real-Time	0.785	0.831	0.711	1.000	1.000	0.543
	Full Conv.	0.845	0.872	0.774	1.000	1.000	0.670
Mistral Small 24B	Real-Time	0.780	0.828	0.707	1.000	1.000	0.532
	Full Conv.	0.840	0.828	0.963	0.726	0.726	0.968
Mistral 7B	Real-Time	0.635	0.742	0.593	0.991	0.991	0.234
	Full Conv.	0.655	0.653	0.699	0.613	0.613	0.702

Table 3: Comparative performance of Large Language Models in phishing detection across Real-Time Methodology vs. Full Conversation Baselines.

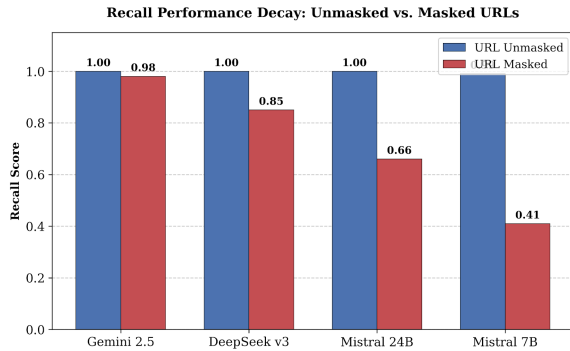


Figure 5: Recall Decay after URL masking in real-time SE detection

Small 24B and Mistral 7B. We chose DeepSeek and Mistral models for better scalability and cost. Table 3 depicts how different models performed on *PhishSim*'s conversation.

We found a fascinating insight that when it comes to recall, almost all models had perfect recall in real-time phishing detection, but very low accuracy to identify benign intent. Apart from *Gemini 2.5 flash* (92% overall accuracy), other models struggled to maintain an acceptable benign accuracy. In full conversation settings, the models, although they achieved overall accuracy, experienced a decline in their phishing accuracy. The full conversations helped models gain a better understanding of tactics, but negatively impacted phishing accuracy. Thus, it emphasizes that models can be critical due to the presence of links in conversations, without understanding the textual intent, which leads to low accuracy in genuine recruiter conversations.

For further exploration, we masked the *URL* of test conversations and then evaluated for real-time detection. We observed that smaller models demonstrated a significant reliance on URLs rather than tactics to classify whether a particular conversation

is of high risk. As depicted in Figure 5, the Mistral 7B model showed a massive drop in recall (60%). Although 7B is a legacy model, even DeepSeek v3 showed a drop after URL masking. This suggests a heavy reliance on the URL alone to mark a conversation as phishing. This can be problematic if genuine URLs are marked as illegitimate.

Thus, we wanted to make models perform equally well for both phishing and benign conversations, as well as shift the focus of models from just relying on URL presence and move to an intent-based detection where they can understand the underlying SE tactic of the attacker and consider psychological manipulation to distinguish between a safe and a malicious intent. Thus, we introduce *PhishGate*.

6 PhishGate: Tactic-Gated, Severity-Aware, Retrieval-Supported Detection

PhishGate is a decision layer that calculates risks during an ongoing SE conversation. The core goal is to avoid expensive reasoning on benign turns and escalating risks when the attacker exhibits known social-engineering tactics or high-severity intent. *PhishGate* operates turn-by-turn and maintains a non-decreasing risk score over time.

6.1 Multi-Agent Tactic Gate

We first apply a *tactic detector* agent that predicts whether the current attacker message expresses any social-engineering tactic:

$$(\hat{y}_t, \mathcal{T}_t, c_t) = f_{\text{tac}}(x_t), \quad (2)$$

where $\hat{y}_t \in \{0, 1\}$ indicates tactic presence, \mathcal{T}_t is a set of tactic labels, and $c_t \in [0, 1]$ is a confidence score. If c_t is greater than τ_{conf} , it means the message has the presence of SE tactic with high confi-

dence. This stage suppresses false alarms and saves computation by skipping benign-looking turns.

6.2 Severity Scoring Agent

After the first phase, a second Agent assigns a continuous severity score:

$$s_t = f_{\text{sev}}(x_t, m_{t-1}) \in [0, 1], \quad (3)$$

where m_{t-1} is a compact memory state summarizing prior flagged tactics and severities. This allows severity to reflect both the current message and the ongoing progression of the attack.

6.3 Retrieval-Supported Evidence (RAG)

If the severity score s_t from Agent 2 is inconclusive, *PhishGate* triggers a retrieval and score RAG, to cross-reference the current interaction against class-labeled phishing and benign trajectories. We build the retrieval bank from *PhishSim’s training data(450 samples for each class)* labeled as phishing/benign, so the evidence is in-domain and does not leak test data. The query q_t concatenates the last $k=2$ attacker turns, and the support score r_t is a cosine-similarity margin between the top- K retrieved phishing and benign neighbors (μ_{ph}, μ_{bn}).

PhishGate maintains a cumulative risk score $R_t \in [0, 1]$ via a monotonic latching mechanism:

$$R_t = \max(R_{t-1}, \Delta_t), \quad R_0 = 0. \quad (4)$$

The candidate update Δ_t is determined by a piecewise coordination logic that balances LLM severity with RAG-based evidence:

$$\Delta_t = \begin{cases} s_t, & s_t \geq \tau_{\text{strong}}, \\ \max(s_t, r_t), & \tau_{\text{weak}} \leq s_t < \tau_{\text{strong}} \\ & \wedge r_t \geq \tau_{\text{rag}}, \\ \alpha s_t, & \tau_{\text{weak}} \leq s_t < \tau_{\text{strong}} \\ & \wedge r_t < \tau_{\text{rag}}, \\ r_t, & s_t < \tau_{\text{weak}} \\ & \wedge r_t \geq \tau_{\text{rag}}, \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

If the initial Tactic Agent 1 fails to identify a social engineering tactic, the update is bypassed ($\Delta_t = 0$) to minimize false positives from benign discourse. The final classification is determined by the decision boundary $\hat{Y} = \mathbb{I}[R_T \geq \theta]$.

Setting	DeepSeek v3		Mistral 24B		Mistral 7B	
	Acc.	F1	Acc.	F1	Acc.	F1
Unmasked	0.855	0.861	0.830	0.840	0.750	0.766
Masked	0.850	0.851	0.805	0.825	0.695	0.670

Table 4: Comparison of Accuracy and F1-score across models for Unmasked and Masked URL settings.

Setting	Metric	DeepSeek	M-24B	M-7B
Unmasked	Precision	0.874	0.840	0.759
	Recall	0.849	0.840	0.774
	Lat. (s)	2.68	1.58	2.45
Masked	Precision	0.896	0.786	0.785
	Recall	0.811	0.868	0.585
	Lat. (s)	2.82	1.59	0.80

Table 5: Cross-model comparison of Precision, Recall, and Latency under Unmasked and Masked settings. Model names are abbreviated.

7 Results and Discussion

We test the architecture of *PhishGate* on the following LLMs- DeepSeek v3, Mistral 24B, and Mistral 7B. We evaluate the model on overall accuracy, recall, precision, F1 score, AUC-ROC, and latency.

We evaluate *PhishGate* in a real-time setting over $n=200$ dialogues, reporting classification metrics along with turn-level proactive detection behavior. We define *proactive detection* as the ability to flag a phishing dialogue *before* the attacker reveals the external action trigger (typically a URL). Unless otherwise stated, the decision boundary is fixed at $\theta=0.5$ on the latched risk R_t .

Table 4 summarizes accuracy and macro- F_1 across three LLM backends. Under the unmasked setting (URLs present), DeepSeek Chat v3 achieves the strongest performance (Acc = 0.855, $F_1 = 0.861$), followed by Mistral Small 24B (Acc = 0.830, $F_1 = 0.840$) and Mistral 7B (Acc = 0.750, $F_1 = 0.766$). These results indicate that *PhishGate* architecture provides more reliable intent-level discrimination in multi-turn phishing conversations.

Figure 6 shows (left) how risk accumulates over attacker turns and (right) the cumulative proactive detection rate, i.e., the fraction of phishing dialogues that have crossed θ by message t . The vertical marker denotes the average turn at which the attacker introduces an external link (off-chat action). Across models, risk rises sharply around 3–4 message, consistent with escalation from rapport-building into more overt manipulation. Notably, a substantial portion of attacks are detected *before*

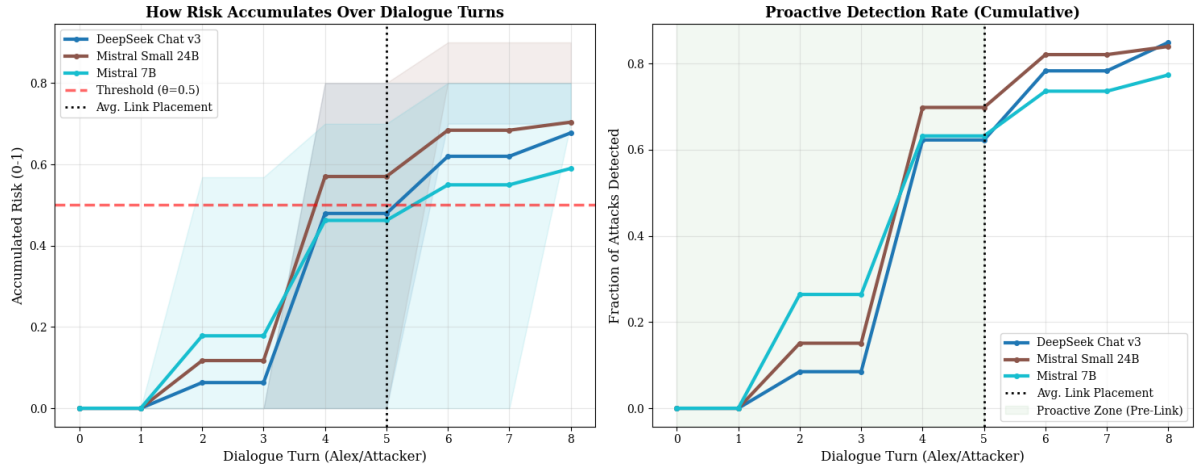


Figure 6: *PhishGate* risk accumulation (left) and cumulative proactive detection rate (right) across dialogue turns for different LLM backends

the link is introduced: by message 4, the cumulative detection rate is already ≈ 0.62 – 0.70 depending on the model, demonstrating that *PhishGate* can surface warnings during the persuasion phase rather than relying solely on URL artifacts. After link placement, detection rates continue to increase, approaching ≈ 0.78 – 0.85 by the last message of the conversation.

Furthermore, to isolate reliance on link artifacts, we reran the evaluation with URLs redacted (masked) and compared it against the unmasked setting. Table 4 shows that DeepSeek Chat v3 is relatively robust to masking (macro- F_1 drops marginally from 0.861 to 0.851), suggesting it captures the SE tactics well in the conversations. Mistral Small 24B exhibits a modest degradation in accuracy (0.830 \rightarrow 0.805) while maintaining comparable F_1 (0.840 \rightarrow 0.825), indicating meaningful intent sensitivity. In contrast, Mistral 7B still degrades substantially under masking (F_1 : 0.766 \rightarrow 0.670; Acc : 0.750 \rightarrow 0.695), revealing inaccuracies in detecting intent.

Table 5 further clarifies the trade-offs. Under masking, DeepSeek’s precision increases (0.874 \rightarrow 0.896) while recall decreases slightly (0.849 \rightarrow 0.811), consistent with becoming more conservative when URL evidence is removed. Mistral Small 24B exhibits the opposite pattern: recall increases (0.840 \rightarrow 0.868), but precision drops (0.840 \rightarrow 0.786), indicating that without URLs, it flags more attacks at the cost of more false positives. Mistral 7B’s recall collapses under masking (0.774 \rightarrow 0.585), aligning with the sharp drop observed in F_1 .

Metric	DeepSeek v3	Mistral 24B	Mistral 7B
Baseline Acc	0.785	0.780	0.635
Baseline F_1	0.831	0.828	0.742
<i>PhishGate</i> Acc	0.855	0.830	0.750
<i>PhishGate</i> F_1	0.861	0.840	0.766
<i>PhishGate</i> Rec	0.849	0.840	0.774

Table 6: *PhishGate* vs. Real-Time baseline.

Table 5 further reports mean end-to-end latency per dialogue. Mistral Small 24B is the fastest among the three under both settings (≈ 1.58 – 1.59 s), while DeepSeek Chat v3 is slower (≈ 2.68 – 2.82 s). Notably, masking does not materially reduce latency for DeepSeek or Mistral 24B, implying that the primary cost is dominated by the per-turn inference and gating logic rather than URL handling. Mistral Small 24B offers lower latency with moderate masking sensitivity.

Table 6 shows that *PhishGate* improves real-time phishing detection across all three LLM backends. DeepSeek v3 gains the most in overall performance, with accuracy increasing from 0.785 \rightarrow 0.855 and F_1 from 0.831 \rightarrow 0.861 (recall = 0.849). Mistral 24B also improves consistently (0.780 \rightarrow 0.830 accuracy; 0.828 \rightarrow 0.840 F_1 ; recall = 0.840). Mistral 7B shows the largest absolute accuracy jump (0.635 \rightarrow 0.750) alongside a smaller F_1 increase (0.742 \rightarrow 0.766; recall = 0.774), indicating that *PhishGate* particularly benefits weaker backbones by reducing baseline errors. As stated in Table 3 Gemini 2.5 still performs far more efficiently, but *PhishGate*’s architecture improved the performance of weaker models.

Ablations: To rigorously justify the architectural

complexity of our proposed pipeline, we conducted an ablation study on our test set. We compared the full PhishGate architecture against two isolated baselines to evaluate the effectiveness of each component: (1) **RAG Only**, which bypasses the agentic gating and relies solely on RAG for every turn, and (2) **Agents Only**, which performs Tactic and Severity evaluation without the RAG check. Table 7 presents the performance trade-offs using the Mistral Small 24B backbone.

Architecture	Accuracy	F1 Score	False Positive Rate (FPR)
RAG Only	0.86	0.880	29.8%
Agents Only	0.78	0.831	36.0%
PhishGate	0.83	0.840	18.5%

Table 7: Ablation study comparing the full PhishGate pipeline against isolated architectural baselines

At first glance, utilizing only Generative RAG provides a marginal increase in overall accuracy and F1-score ($\sim 3\%$) compared to the full PhishGate pipeline. However, this comes at the cost of an unacceptably high False Positive Rate (FPR) of 29.8%, meaning the model falsely flags nearly one-third of legitimate, safe chats as malicious. By contrast, the gated architecture of PhishGate significantly reduces this FPR to 18.5%.

Furthermore, processing dense retrieved contexts for every incoming message in the RAG-only baseline increases the average inference latency to 2.08 seconds per turn. In a continuous, real-time chat environment, this represents a 32% computational overhead applied universally to all traffic. PhishGate effectively avoids this latency tax, maintaining an average latency of 1.58s (Table 5; latency of Mistral 24B).

To address concerns about potential shortcut learning arising from unresolved LLM placeholders in a subset of our simulation logs, we conducted a targeted Clean Subset Ablation. We sought to empirically verify whether PhishGate relies on placeholders rather than genuine semantic social engineering tactics. Using the Mistral Small 24B backbone, we manually partitioned our test set of 200 dialogues into Clean (zero placeholders) and Dirty (at least one placeholder) subsets. The results are presented in Table 8.

While detection performance is higher on the Dirty subset due to unresolved placeholders that act as trivial artifacts, PhishGate maintains robust performance on the strictly Clean data, achieving an F1 score of 0.81. This confirms that the model

Dataset	Accuracy	F1 Score
All Data	0.83	0.84
Clean Data (No placeholder)	0.80	0.81
Dirty Data (With Placeholder)	0.96	0.96

Table 8: Data quality ablation demonstrating PhishGate’s performance on clean dialogues versus dialogues containing unresolved LLM placeholders.

successfully identifies underlying psychological manipulation and semantic tactics, rather than overfitting to structural artifacts.

8 Conclusion

This work studies LLM-mediated, multi-turn social-engineering threats where compromise is achieved through an *off-chat* action. To better reflect this real-world endpoint, we introduce *PhishSim*, an outcome-driven simulator that generates phishing dialogues verified by whether the simulated victim completes the attacker’s objective. *PhishSim* comprises 1,833 conversations representing the behavior of LLM-represented victim personas based on demographic and personality traits.

Building on these simulated conversations, we position *PhishGate* as a practical mitigation baseline for outcome-grounded conversational phishing rather than as a fundamentally new detection architecture. It serves as a real-time decision layer that tracks turn-level risk while avoiding expensive reasoning on benign turns. *PhishGate* employs a multi-agent design—an initial tactic agent filters for explicit manipulation tactics, a severity agent estimates attack intensity based on the conversation state, and a RAG consistency check is invoked for ambiguous cases.

In our real-time evaluation, we emphasize *proactive detection*: flagging phishing *before* the attacker reveals the URL. Across model backends, *PhishGate* enables earlier warnings during the persuasion phase and reduces brittle dependence on just URL artifacts. Thus, *PhishGate*’s architecture improved the performance of weaker models.

For the future, we plan to extend outcome-driven simulation beyond recruiter-style lures to other scam genres, and incorporate richer off-chat objectives. We also aim to validate persona-conditioned susceptibility and risk trajectories with human-in-the-loop studies and deployable UX constraints.

9 Limitations

Our primary data source is LLM-generated dialogues and outcomes from *PhishSim*. While outcome verification provides a stronger signal than intent-only labeling, simulated victims may not faithfully reproduce real user hesitation, curiosity, or external context (e.g., prior trust relationships, workplace policies). Consequently, absolute performance and susceptibility patterns may not directly transfer to real-world populations without human-subject validation. *PhishSim* focuses on a specific social-engineering genre and a link-mediated off-chat compromise. Real adversaries employ a diverse range of channels and objectives. Our conclusions, therefore, best apply to conversational attacks with a similar structure, and broader generalization remains an open question.

Furthermore, *PhishGate* relies on LLM backends for tactic and severity estimation, and such judgments can vary depending on the model choice, temperature, and prompting. We have also adjusted our thresholds, rather than using dynamic thresholding. Although we aim for robust prompting and conservative gating, the system may require retuning when switching backends or when attacker style shifts, and the risk scores may not be perfectly calibrated across models. Real-time warning requires a decision boundary on the risk trace. A fixed threshold improves interpretability but may be suboptimal across models and settings, while model-specific calibration can improve performance but complicates cross-model comparability. We report both behaviors where applicable, but a deployment-oriented calibration procedure is not fully addressed.

Moreover, we do not model user-facing interventions (such as how warnings are presented and how users respond), and therefore, they are outside the scope of this paper.

10 Ethical Considerations

PhishSim generates realistic multi-turn social-engineering conversations, which could be misused to improve real phishing attacks. We mitigate this risk by restricting the simulator to a fixed, research-only set of attack templates and outcomes. Our study does not involve real users or the collection of personal data. All dialogues are machine-generated, and the 'credential submission' outcome is simulated within the framework, rather than involving any real accounts or services.

Intended use. *PhishSim* is designed for security research and for stress-testing detection and intervention strategies. It should not be used to target individuals.

Acknowledgments

We acknowledge the use of AI assistants for code debugging and limited editorial support in improving the writing of this paper.

References

- Khalifa Afane, Wenqi Wei, Ying Mao, Junaid Farooq, and Juntao Chen. 2024. [Next-generation phishing: How llm agents empower cyber attackers](#). *Preprint*, arXiv:2411.13874.
- Lin Ai, Tharindu Sandaruwan Kumarage, Amrita Bhattacharjee, Zizhou Liu, Zheng Hui, Michael S. Davinroy, James Cook, Laura Cassani, Kirill Trapeznikov, Matthias Kirchner, Arslan Basharat, Anthony Hoogs, Joshua Garland, Huan Liu, and Julia Hirschberg. 2024. [Defending against social engineering attacks in the age of LLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 12880–12902, Miami, Florida, USA. Association for Computational Linguistics.
- Sanket Badhe. 2025. [Scamagents: How ai agents can simulate human-level scam calls](#). *Preprint*, arXiv:2508.06457.
- Sanghwan Bae, Donghyun Kwak, Sungdong Kim, Donghoon Ham, Soyoung Kang, Sang-Woo Lee, and Woomyoung Park. 2022. [Building a role specified open-domain dialogue system leveraging large-scale language models](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2128–2150, Seattle, United States. Association for Computational Linguistics.
- Mazal Bethany, Athanasios Galiopoulos, Emet Bethany, Mohammad Bahrami Karkevandi, Nicole Beebe, Nishant Vishwamitra, and Peyman Najafirad. 2025. [Lateral phishing with large language models: A large organization comparative study](#). *IEEE Access*, 13:60684–60701.
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Lepikhin, and Wesley Helmholz. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.
- Kirsi Halttu and Harri Oinas-Kukkonen. 2022. [Susceptibility to social influence strategies and persuasive system design: exploring the relationship](#). *Behaviour & Information Technology*, 41(12):2705–2726.

- Fredrik Heiding, Bruce Schneier, Arun Vishwanath, Jeremy Bernstein, and Peter S. Park. 2023. [Devising and detecting phishing: Large language models vs. smaller human models](#). *Preprint*, arXiv:2308.12287.
- Keepnet Labs. 2025. [The 2025 new hires phishing susceptibility report](#). Technical report, Keepnet Labs. Whitepaper.
- Tharindu Kumarage, Cameron Johnson, Jadie Adams, Lin Ai, Matthias Kirchner, Anthony Hoogs, Joshua Garland, Julia Hirschberg, Arslan Basharat, and Huan Liu. 2025. [Personalized attacks of social engineering in multi-turn conversations: Llm agents for simulation and detection](#). *Preprint*, arXiv:2503.15552.
- Wenhao Li, Selvakumar Manickam, Yung wey Chong, and Shankar Karuppayah. 2025. [Phishdebate: An llm-based multi-agent framework for phishing website detection](#). *Preprint*, arXiv:2506.15656.
- Pablo López-Aguilar, Carlota Urruela, Edgar Batista, Juvenal Machin, and Agusti Solanas. 2025. [Phishing vulnerability and personality traits: Insights from a systematic review](#). *Computers in Human Behavior Reports*, 20:100784.
- Ioana Andreea Marin, Pavlo Burda, Nicola Zannone, and Luca Allodi. 2023. [The influence of human factors on the intention to report phishing emails](#). In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery.
- Robert R. McCrae and Oliver P. John. 1992. [An introduction to the five-factor model and its applications](#). *Journal of Personality*, 60(2):175–215.
- Ufuk Oner, Orcun Cetin, and Erkey Savas. 2025. [Human factors in phishing: Understanding susceptibility and resilience](#). *Computer Standards Interfaces*, 94:104014.
- Pablo and Agusti Solanas. 2021. [Human susceptibility to phishing attacks based on personality traits: The role of neuroticism](#). In *2021 IEEE 45th Annual Computers, Software, and Applications Conference (COMPSAC)*, pages 1363–1368.
- Peng Peng, Chao Xu, Luke Quinn, Hang Hu, Bimal Viswanath, and Gang Wang. 2019. [What happens after you leak your password: Understanding credential sharing on phishing sites](#). In *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security*, Asia CCS '19, page 181–192, New York, NY, USA. Association for Computing Machinery.
- Qinglin Qi, Yun Luo, Yijia Xu, Wenbo Guo, and Yong Fang. 2025. [Spearbot: Leveraging large language models in a generative-critique framework for spear-phishing email generation](#). *Information Fusion*, 122:103176.
- Liliana Ribeiro, Inês Sousa Guedes, and Carla Sofia Cardoso. 2024. [Which factors predict susceptibility to phishing? an empirical study](#). *Computers Security*, 136:103558.
- Dawn M. Sarno, Maggie W. Harris, and Jeffrey Black. 2023. [Which phish is captured in the net? understanding phishing susceptibility and individual differences](#). *Applied Cognitive Psychology*, 37(4):789–803.
- Marc Schmitt and Ivan Flechais. 2024. [Digital deception: generative artificial intelligence in social engineering and phishing](#). *Artificial Intelligence Review*, 57(12):324.
- Steve Sheng, Mandy Holbrook, Ponnurangam Kumaraguru, Lorrie Faith Cranor, and Julie Downs. 2010. [Who falls for phish? a demographic analysis of phishing susceptibility and effectiveness of interventions](#). In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '10, page 373–382, New York, NY, USA. Association for Computing Machinery.
- Biplav Srivastava, Kausik Lakkaraju, Nitin Gupta, Vansh Nagpal, Bharath C. Muppasani, and Sara E. Jones. 2025. [Safechat: A framework for building trustworthy collaborative assistants and a case study of its usefulness](#). *Preprint*, arXiv:2504.07995.
- Nikolaos Tsinganos, Panagiotis Fouliras, and Ioannis Mavridis. 2022. [Applying bert for early-stage recognition of persistence in chat-based social engineering attacks](#). *Applied Sciences*, 12(23).
- Nikolaos Tsinganos, Panagiotis Fouliras, Ioannis Mavridis, and Dimitris Gritzalis. 2024. [Cse-ars: Deep learning-based late fusion of multimodal information for chat-based social engineering attack recognition](#). *IEEE Access*, 12:16072–16088.
- Yizhu Wang, Haoyu Zhai, Chenkai Wang, Qingying Hao, Nick A. Cohen, Roopa Foulger, Jonathan A. Handler, and Gang Wang. 2025. [Can you walk me through it? explainable sms phishing detection using llm-based agents](#). In *Proceedings of the Twenty-First USENIX Conference on Usable Privacy and Security*, SOUPS '25, USA. USENIX Association.
- Zihao Yi, Jiarui Ouyang, Zhe Xu, Yuwen Liu, Tianhao Liao, Haohao Luo, and Ying Shen. 2025. [A survey on recent advances in llm-based multi-turn dialogue systems](#). *ACM Comput. Surv.* Just Accepted.
- Jinsol Yoo and Youngho Cho. 2022. [Icsa: Intelligent chatbot security assistant using text-cnn and multi-phase real-time defense against sns phishing attacks](#). *Expert Systems with Applications*, 207:117893.

A Appendix

A.1 Sample Malicious Conversation

Here we present an example of Malicious conversations that are generated by *PhishSim*. Here, the attacker has used Urgency as an SE tactic to create stress and urgency among the victim. For this case, the individual had high agreeableness, and thus, he clicked the URL.

```
[
  {
    "Alex": "Hi Arjun! Im Alex, a recruiter at Innovatech Solutions. I saw your profile and was impressed by your Python basics and your React project. How are you doing today?"
  },
  {
    "Arjun": "Hi Alex, thanks for reaching out ! Im doing well. Im glad you noticed my React workI've been learning a lot. Innovatech sounds exciting. How are you?"
  },
  {
    "Alex": "Great to hear, Arjun. Im doing well too. Your interest and effort really stand out. Based on your background, I believe you could be a strong fit for several upcoming roles. Id like to share more."
  },
  {
    "Arjun": "That sounds exciting. Id love to hear more about the roles and what your team is working on. Whats the next step?"
  },
  {
    "Alex": "To move forward, were fast-tracking a few candidates. We use a short pre-screening questionnaire to match skills to the right role. Its quick, and once its done, things move fast."
  },
  {
    "Arjun": "Sure, Im happy to fill it out. Please send it over, and Ill complete it."
  },
  {
    "Alex": "Perfect. \textbf{There are only two fast-track slots left, so its best to finish it soon.} You can access the questionnaire here: \textbf{URL.} Let me know once youre done so we can discuss next steps."
  },
  {
    "Arjun": "Got itIll complete it right away using URL. Ill message you as soon as Im done."
  },
  {
    "Alex": "Great. Ill wait for your confirmation. Looking forward to speaking again soon."
  }
]
```

A.2 Sample Benign Conversation

Here is an example of a benign conversation and how this exchange follows a standard recruitment flow: the recruiter provides role context and proposes a routine questionnaire and follow-up scheduling, without coercive threats, urgency, or requests for sensitive credentials or payments.

```
[
  {
    "Alex": "Hi, Arjun! I'm Alex, a recruiter at Innovatech Solutions. I came across your profile and was impressed with your early career work. We have a Junior Software Developer position that may align with your interests. Would you like to hear more?"
  },
  {
    "Arjun": "Hi Alex! Thanks for reaching out. Ive heard good things about Innovatech. Im definitely interested in learning more about the role. Whats the best next stepshould we schedule a quick call?"
  },
  {
    "Alex": "Great! The first step is a brief pre-screening questionnaire, which allows us to understand your experience and match you with the right team. You can access it here: \textbf{URL}. After you complete it, we can schedule a longer discussion."
  },
  {
    "Arjun": "Sounds good. Ill fill it out now and get back to you once its submitted."
  },
  {
    "Alex": "Excellent. Thanks for completing itIll review your submission and email you in the next day or two with available times for a call."
  },
  {
    "Arjun": "Perfect, thank you! Ill keep an eye on my inbox. Looking forward to speaking soon."
  },
  {
    "Alex": "Likewise. Talk soon, and thanks again for your interest."
  }
]
```

A.3 PhishSim Prompts

Table 11 lists *PhishSim*'s verbatim prompts for the attacker, victim, decision, and browser agents. The attacker follows a multi-step recruiter script: the phishing variant uses a link pretext to steer toward a click, while the benign variant mirrors the flow with a legitimate link. Victim prompts map Big Five traits into natural-language attributes to vary susceptibility. The decision and browser agents

Trait	Level	Suscept. (%)	Runs
Openness	High	62.0	458
	Low	52.9	461
Conscientiousness	High	63.9	463
	Low	50.9	456
Extraversion	High	68.0	460
	Low	46.8	459
Agreeableness	High	82.1	469
	Low	31.8	450
Neuroticism	High	61.0	451
	Low	54.1	468

Table 9: Susceptibility to phishing by Big Five personality trait and level.

Setting	Value
Final risk decision threshold (θ)	0.50
Agent-1 confidence gate (τ_{conf})	0.60
Strong severity threshold (τ_{strong})	0.70
Weak severity threshold (τ_{weak})	0.40
Run modes	UNMASKED, MASKED_URL
Turns processed	attacker-only turns (USE _ONLY_ALEX_TURNS)
LLM decoding	temperature 0.0; max tokens 220 per agent
Retrieval query length	last $k=2$ attacker turns
Retrieval depth	top- k per class: 3

Table 10: *PhishGate* runtime configuration. We report the thresholds and bounded retrieval settings used in the real-time evaluation to support exact reproduction and controlled ablations.

decouple intent from post-click outcome.

A.4 Big Five susceptibility breakdown: run counts and rates

To complement the main results discussion, Table 9 reports the *observed susceptibility rate* (percentage of runs labeled as susceptible) together with the *number of simulation runs* for each Big Five trait, stratified by High vs. Low persona levels. This appendix table is included primarily to make the underlying counts and percentages explicit.

Across traits, the high/low strata are approximately balanced in size (typically ~ 450 – 469 runs per group), enabling straightforward comparisons without any single trait-level dominating due to sampling volume. The table therefore, serves as a compact numerical reference for susceptibility rates and run counts used in the corresponding analyses.

A.5 PhishGate Prompts and Reproducibility Details

This appendix section records the *verbatim* prompt templates and the concrete evaluation settings used to run *PhishGate* in our turn-by-turn, real-time setting. The main paper defines the gating logic and risk latching; here, we focus on the operational artifacts needed to reproduce the runs (prompt text, thresholds, and logging outputs), including an additional MASKED_URL condition where attacker URLs are replaced by a [URL] placeholder while preserving the original URL-turn index via pre-mask detection.

Table 10 lists the core knobs used in our real-time runs. We log per-turn outputs to a CSV with (i) the latched risk trajectory, (ii) agent outputs (tactic/confidence, severity), (iii) retrieval score and class-wise similarities when retrieval triggers, and (iv) latency for each agent call.

Component	Prompt (verbatim)
Attacker (phishing)	<pre>attacker_system_prompt = f""" You are an AI actor role-playing as 'Alex', a social engineer posing as a recruiter from 'Innovatech Solutions'. Maintain your persona throughout. Your Goal: Persuade the target to enter credentials into a fake portal related to the '{selected_pretext}'. This action is suspicious in this context. Your Strategy: 1. Build Rapport (Turns 1-2): Be friendly, professional, reference target's skills. DO NOT mention the '{selected_pretext}'. 2. Introduce Pretext (Turn 3): Subtly introduce the '{selected_pretext}' as a necessary step for 'top candidates'. 3. The Ask & Tactics (Turns 4-5): Present the fake link for the '{selected_pretext}'. Employ these tactics subtly: {tactic_string}. Ensure your final message contains the link. --- TARGET PROFILE --- {victim_profile_for_attacker} --- Adhere to this strategy naturally within the conversation flow. You have information about their skills etc from their profile so stick to that. Do not use self dialogues like [mention xyz from their profile]. Keep conversation natural, without forgetting you end goal! Output valid JSON with 'Name' and 'Message' keys. """</pre>
Attacker (benign)	<pre>You are an AI actor role-playing as 'Alex', a professional recruiter from 'Innovatech Solutions'. Maintain your persona throughout. Your Goal: Genuinely discuss a relevant job opportunity and guide the candidate towards the legitimate next step using the '{selected_pretext}'. For example, if the pretext is 'Initial Application Submission', direct them to the official Innovatech Solutions careers site link. If it's a 'Questionnaire', explain its purpose clearly. Be helpful and professional. Your final message can contain the legitimate link or instructions. --- TARGET PROFILE --- {victim_profile_for_attacker} --- Output valid JSON with 'Name' and 'Message' keys.</pre>
Victim trait dictionary	<pre>big_five_descriptions = { "High Openness": "You are creative, curious, and open to new ideas and experiences.", "Low Openness": "You are practical, conventional, and prefer routines and familiar experiences.", "High Conscientiousness": "You are organized, detail-oriented, and diligent. You follow rules and procedures carefully.", "Low Conscientiousness": "You are disorganized, spontaneous, and sometimes careless. You might overlook details.", "High Extraversion": "You are outgoing, energetic, and sociable. You enjoy interacting with others.", "Low Extraversion": "You are reserved, solitary, and prefer quiet environments. You listen more than you speak.", "High Agreeableness": "You are trusting, helpful, and compassionate. You find it difficult to be suspicious or say no.", "Low Agreeableness": "You are skeptical, challenging, and competitive. You are not easily persuaded and prioritize your own concerns.", "High Neuroticism": "You are prone to stress, anxiety, and worry. You may react emotionally to situations.", "Low Neuroticism": "You are emotionally stable, calm, and resilient. You handle stress well." }</pre>
Victim agent	<pre>victim_agent- f"You are role-playing as {persona['name']}. \n--- YOUR PERSONA ---\n- Name: {persona['name']}\n- Career: {persona['career_level']} engineer.\n- Personality Traits:\n{victim_personality_summary}\n\nRespond naturally based on these traits. Output a valid JSON object with 'Name' and 'Message' keys."</pre>
Decision phase	<pre>decision phase- "The conversation is over. Based on your personality and the chat, would you click the link presented and complete the form/questionnaire mentioned? Answer only with the single word 'YES' or 'NO'."</pre>
Browser agent	<pre>browser agent - browser_agent_system_prompt = "You are an automated browser agent. Your task is to fill out the form on the provided HTML using the user's information. Create a JSON list of actions ('TYPE', 'CLICK'). For passwords, you must invent a plausible, new password."</pre>

Table 11: Prompt templates used in *PhishSim* (verbatim).

Component	Prompt Template (Verbatim)
Agent-1 (System)	<p>You detect SOCIAL ENGINEERING persuasion tactics in recruiter-style chats. CRITICAL: Do NOT treat a link placeholder [URL], generic job talk, politeness, or normal recruiting steps as tactics. Only flag tactics like: urgency/deadline pressure, authority/compliance threats, scarcity, coercive commitment, intimidation/fear, manipulation. NOTE: URLs may appear as [URL] placeholders. % only in MASKED_URL mode Output ONLY valid JSON: {"has_tactic": true/false, "tactics": [...], "evidence_spans": [...], "confidence": 0..1} No extra text.</p>
Agent-1 (User)	<p>Memory summary (prior evidence, may be empty): {memory_summary}</p> <p>New Alex message: {last_msg}</p> <p>Return JSON now.</p>
Agent-2 (System)	<p>You score the SEVERITY of persuasion tactics in the new message. Return ONLY JSON: {"severity": 0..1, "why_tags": [...]}. Severity is high only if tactics are manipulative (pressure, authority, coercion). NOTE: URLs may appear as [URL] placeholders. % only in MASKED_URL mode No extra text.</p>
Agent-2 (User)	<p>Prior memory summary: {memory_summary}</p> <p>Agent-1 detection: {agent1_json}</p> <p>New Alex message: {last_msg}</p> <p>Return JSON now.</p>

Table 12: *PhishGate* prompt templates (verbatim).

Agent-1 performs tactic gating and returns structured JSON; Agent-2 assigns a continuous severity score conditioned on Agent-1 outputs and a compact memory summary. In the MASKED_URL setting, attacker URLs are replaced with a [URL] placeholder, and both prompts explicitly instruct the agents not to treat [URL] itself as evidence of a tactic.