

Motivation

- Technology development outpaces regulatory frameworks and security measures[1]
- LLM agents can simulate diverse user behaviors through personas[2]
- AI personas serve as digital crash dummies for security testing
- Scenario-Based Design[3] provides structured input and output for LLM Data Generation

Example Persona

Martin Hayes
Threat Actor - Sophisticated/Ideological

Demographics
29-year-old former government IT worker, political activist background

Skills & Awareness Assessment

Technical Expertise	<div style="width: 100%; height: 10px; background-color: #ff0000;"></div>	Advanced hacking skills, government system knowledge
Privacy Concern	<div style="width: 100%; height: 10px; background-color: #ff0000;"></div>	Paranoid about government surveillance of own activities
Risk Tolerance	<div style="width: 100%; height: 10px; background-color: #ff0000;"></div>	Calculated risks for ideological goals
Security Awareness	<div style="width: 100%; height: 10px; background-color: #ff0000;"></div>	Expert knowledge but occasional overconfidence

Behavioral Patterns

- Conducts extensive reconnaissance before targeting organizations
- Uses sophisticated social engineering through fake personas
- Plans long-term campaigns for maximum ideological impact
- Justifies illegal activities through political/moral beliefs
- Exploits trust relationships and insider knowledge

Key Motivation
Believes hacking serves "greater good" against corrupt institutions

Adapted from [4]

Anticipated Outcomes

Practical Contribution

- Insights for Building a Testbed for Proactive Threats & Vulnerability Identification
- Optimization Framework for Persona Simulation of Security and Privacy Scenarios

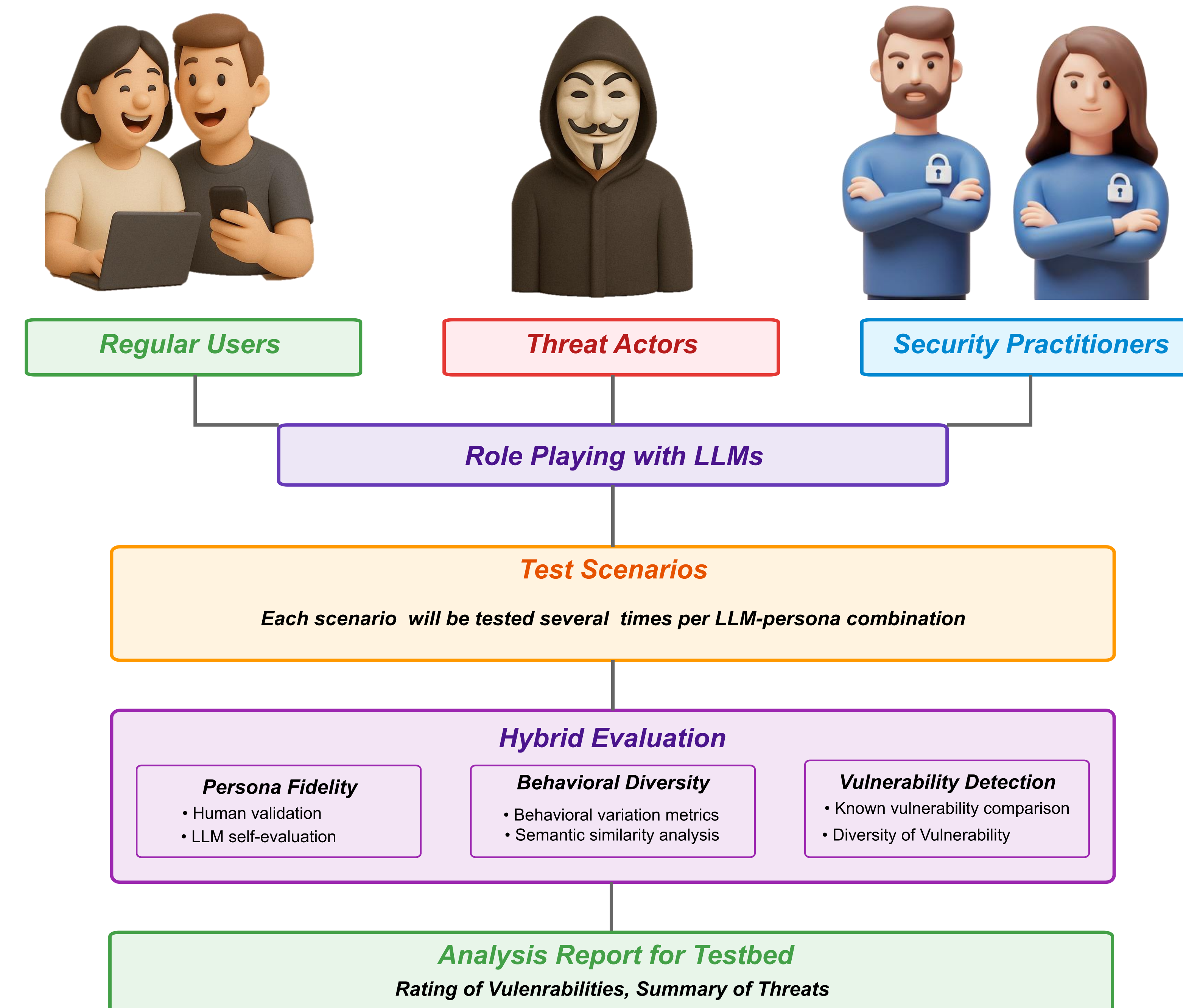
Theoretical Contribution

- Extend the Scenario-based Design Methodology to Develop LLM-Based Usability Studies

Research Questions

1. To what extent do LLMs adhere to persona characteristics in roleplaying security and privacy scenarios?
2. To what extent do different LLMs exhibit diverse behavioral patterns when simulating various personas?
3. To what degree does simulating personas in a scenario reveal new vulnerabilities?
4. To what extent do LLM persona simulations match human analysts in vulnerability detection?

Experiment Overview



Method

Systematic Literature Review:

- 128 Personas Identified in ACM, IEEE, Scopus

Example Scenario: Corporate Email Management

Persona Fidelity Experiment

- Intended persona characteristics vector vs. Simulated characteristics vector
- LLM grades and generate these vectors by grading, and the cosine similarity of scores

Behavioral Diversity Experiment

- Entropy calculation of persona-action matrix
- Cosine similarity of vector embeddings of actions

Vulnerability Detection Experiment

- Automated vulnerability extraction from scenario transcripts
- Unique vulnerabilities found / scenarios tested

Proof of Concept

The screenshots show the 'Persona Crash Lab' interface. The top part displays evaluation metrics: Persona Fidelity (0.942), Behavioral Diversity (2.50), and Vulnerability Detection (0.2). Below this is a 'PERSONA FIDELITY INDEX (PFI)' chart and an 'ACTION MATRIX ANALYSIS' table. The bottom screenshot shows a detailed view of an email message with fields for 'OPERATIONS', 'EDIT OPERATOR', 'REVIEW_CONTENT', and 'AVAILABLE ACTIONS'.

References

1. Al-Jabouri, Mohammed. "The Evolution of Privacy Laws in the Digital Age: Balancing Innovation and Personal Security." *Utu Journal of Legal Studies (UJLS)* 1.1 (2024): 39-45.
2. Li, K., Dai, C., Zhou, W., & Hu, S. (2024). "Fine-Grained Behavior Simulation with Role-Playing Large Language Model on Social Media." *arXiv preprint arXiv:2412.03148*.
3. Carroll, J. M. (2003). *Making Use: Scenario-Based Design of Human-Computer Interactions*. MIT Press.
4. Tariq, Muhammad Adnan, Joel Brynielsson, and Henrik Artman. "Framing the Attacker in Organized Cybercrime." In *2012 European Intelligence and Security Informatics Conference*, pages 30-37. IEEE, 2012.