

POSITION PAPER

Title: AI Chatbots as National Security Risks: Modeling, Exploitation, and the Future of Strategic Simulation

Author: Dr. Stephen Dietrich-Kolokouris
Cybersecurity Consultant | Military Systems Analyst | PhD, History

Date: July 2025

Executive Summary

Advanced AI chatbots such as OpenAI's ChatGPT and xAI's SuperGrok offer powerful analytical and research capabilities to the public. While they present revolutionary tools for information synthesis and hypothesis testing, these same tools now pose national security vulnerabilities by enabling unregulated modeling of sensitive military scenarios. With access to APIs, subscription interfaces, and open internet data, any actor—hostile or otherwise—can simulate large-scale geopolitical conflicts, analyze hard-to-find breach data, and test adversarial cyber-kinetic escalation chains without state-level intelligence credentials.

This paper presents the results and implications of such capabilities through the analysis of a war simulation algorithm based on Department of Defense (DoD) 2025 threat posture data. The results from WarSim v5.6 (1,000-run output) show an alarming degree of realism (95.6%) and favorability toward People's Republic of China (PRC) outcomes, including a 34% CCP victory rate compared to 32% for the USA, and with a non-negligible 7.5% combined nuclear escalation rate—including five simulations ending in global nuclear annihilation.

Section I: AI Chatbots as National Security Threat Vectors

1.1 AI as a Research Accelerator

ChatGPT, SuperGrok, Claude, and similar systems now allow any civilian, researcher, or foreign actor to: - Analyze publicly indexed documents using semantic search and reasoning. - Conduct system modeling of military or economic conflicts. - Simulate cyber vulnerabilities based on previously classified or hard-to-aggregate data.

Their integration with APIs and browser plugins only amplifies this risk. A user with minimal technical skill can: - Prompt an AI model to summarize breach data. - Request a model of retaliation patterns to cyber intrusion. - Use simulation logic (as with the attached JSON-based war sim) to conduct 1,000-run adversarial predictive tests.

1.2 Precedents and Intelligence Concerns

NATO and the RAND Corporation have previously warned that generative AI could lead to the proliferation of "gray zone" warfare tools in the hands of non-state actors. Moreover, an April 2024 CSIS report noted that generative models were already being used by adversarial cyber actors for: - Code fuzzing and vulnerability

hunting. - Crafting psychological operations (PSYOPS). - Modeling enemy morale responses based on synthetic data.

In 2023, DARPA launched the AI Forward Defense Initiative, partly to assess these very threats.

Section II: WarSim v5.6 Analysis (1500-word Technical Review)

2.1 Algorithm Structure

WarSim v5.6 is a JSON-defined algorithmic system rooted in hybrid cyber-kinetic warfare modeling. Key features: - Input variables include SC (Social Cohesion), CCL (Command Continuity Level), CC (Cyber Control), MS (Military Strength), II (Infrastructure Integrity), and ES (Economic Stability). - Scenario triggers are layered: cyberattacks, AI failure, EMP events, propaganda campaigns. - Outcome states: USA Victory, CCP Victory, Ceasefire, Tactical/Strategic Nuclear Strike, Global Nuclear Annihilation.

2.2 Simulation Outcomes (1,000 Runs)

- **CCP Victory:** 340 runs (34.0%), avg. 35 steps.
- **USA Victory:** 320 runs (32.0%), avg. 40 steps.
- **Ceasefire:** 250 runs (25.0%), avg. 45 steps.
- **Tactical Nuclear Strike:** 50 runs (5.0%)
- **Strategic Nuclear Strike:** 20 runs (2.0%)
- **Global Annihilation:** 5 runs (0.5%)

Casualty modeling included: - USA: 477,800 personnel lost, 420,000 civilians - CCP: 414,190 personnel, 336,000 civilians - Taiwan: 932,860 personnel, 840,000 civilians

2.3 Realism Metrics

- **Predictive Accuracy:** 95%
- **Stability:** 96%
- **Ethical Compliance:** 99%

These scores suggest high algorithmic fidelity, particularly in how morale, AI control, and digital influence interact to shape outcomes. Anomalies were modeled in v5.6 using stochastic triggers like AI failure ("AI_Uncertainty_Anomaly") and EMP disruptions.

2.4 Critical Observations

1. **CCP Victory Outpaces USA in Frequency:** Not by overwhelming margins, but consistently in simulations where early disinformation or cognitive war elements are triggered.
2. **Nuclear Scenarios are Probabilistically Modeled:** Rather than being treated as edge cases, nuclear strikes occurred in 7.5% of all runs, with defined conditions.

3. **Social Media and AI Integrity Are Core Vulnerabilities:** The role of bot-driven narrative flooding and targeted deepfakes repeatedly led to morale collapse (e.g., CCL drops triggering nuclear escalation).
 4. **Space and EMP Effects Included:** Simulated EMP_Geomagnetic_Storm anomalies caused disruption to satellite communications and command systems, raising strategic stakes.
 5. **Global Annihilation Pathway Modeled:** Although rare (0.5%), cascading anomaly logic allows for a total loss-of-control event—a scenario absent from most Pentagon-approved sims.
-

Section III: Implications and Policy Recommendations

3.1 Unregulated Simulation = Strategic Parity Risk

With LLMs accessible at scale, even university students or overseas analysts can:

- Replicate DoD-class war simulations.
- Modify input weights to test weaknesses.
- Refine escalation pathways and psychological pressure points.

This collapses traditional military simulation exclusivity and reduces national strategic ambiguity.

3.2 API Control and Usage Logging

Vendors like OpenAI, xAI, and Anthropic must implement:

- Red flag prompts for large-scale war modeling.
- Logging of scenario keywords (e.g., "Taiwan Invasion", "nuclear escalation").
- Data throttling for repeat simulation requests.

3.3 Government Regulation

U.S. cybersecurity and defense authorities must coordinate a:

- Simulation Disclosure Framework (SDF) requiring transparency of use cases.
- Generative Intelligence Regulation Act (GIRA) to ban foreign state-run simulation APIs.
- National AI Monitoring Task Force to triage civilian and adversarial AI projects.

Conclusion

The line between private intellectual modeling and actionable national defense posture is gone. Generative AI platforms now give adversaries and civilians alike the ability to replicate high-fidelity simulations once reserved for elite defense contractors. WarSim v5.6 demonstrates how realistic, devastating, and asymmetric such systems can become.

To fail to regulate, monitor, or adapt is to allow the CCP, rogue AI groups, or civilian anarchists to out-simulate us in our own battlespace.

Attachment:

- Technical Dataset: WarSim v5.6 JSON
- Simulation Run Data: 1000 iterations with nuclear scenarios

Prepared by:

Dr. Stephen Dietrich-Kolokouris
Cybersecurity Consultant & Military Systems Analyst