

# Integrating Artificial Intelligence into Weapon Systems

Philip Feldman  
ASRC Federal

Columbia, MD, USA  
philip.feldman@asrcfederal.com

Aaron Dant  
ASRC Federal

Columbia, MD, USA  
aaron.dant@asrcfederal.com

Aaron Massey  
Information Systems

University of Maryland, Baltimore County  
Baltimore, MD, USA  
akmassey@umbc.edu

*Abstract*—The integration of Artificial Intelligence (AI) into weapon systems is one of the most consequential tactical and strategic decisions in the history of warfare. Current AI development is a remarkable combination of accelerating capability, hidden decision mechanisms, and decreasing costs. Implementation of these systems is in its infancy and exists on a spectrum from resilient and flexible to simplistic and brittle. Resilient systems should be able to effectively handle the complexities of a high-dimensional battlespace. Simplistic AI implementations could be manipulated by an adversarial AI that identifies and exploits their weaknesses.

In this paper, we present a framework for understanding the development of dynamic AI/ML systems that interactively and continuously adapt to their user's needs. We explore the implications of increasingly capable AI in the kill chain and how this will lead inevitably to a fully automated, always on system, barring regulation by treaty. We examine the potential of total integration of cyber and physical security and how this likelihood must inform the development of AI-enabled systems with respect to the “fog of war”, human morals, and ethics.

*Index Terms*—machine learning, computer simulation, human-computer interaction

## I. INTRODUCTION

John Boyd's Observation-Orientation-Decision-Action (OODA) model formalizes the description of the inputs, outputs, experiences, and biases that explain tactical decision-making for individuals and groups. In this model, an adversary attacking either (O)bservation or (O)rientation can create the conditions for incorrect or catastrophic (D)ecisions and (A)ctions. Even without exacerbating conditions such as combat, humans often find decision-making under stress difficult, particularly with incomplete information. The natural human tendency to defer to authorities for decision-making, human or machine, can also lead to disastrous outcomes [1], [2].

Artificial intelligence and machine learning promise to integrate dynamically into the human decision-making processes in ways previous technologies could not, including by being responsive to the operator's cognitive load. In the simplest approach, machines perform only tedious and boring tasks. In a slightly more complex scenario, machines perform as much of the non-decision-making activity as possible so that humans can focus completely on the task at hand. In the most complex

scenario, humans are not directly involved with the system as it performs the task independently. These “always on” systems respond to threats that are beyond the capability of real-time human supervision. Human interaction is restricted to activities such as training the system in offline or simulated environments

Many countries are currently pursuing an ambitious AI agenda, including the United States and several potential adversaries [3], [4]. Secretive development of lethal autonomous weapons systems (LAWS) leads to the conditions for an AI arms race. Tactically, each side's human/AI system would be attempting to “turn inside” the adversaries' OODA loop. Although AI may be an advantage in cognitive offloading of mundane tasks or through increased speed and capability in battle, it presents an opportunity for a new class of attacks that take advantage of the latent, high-dimensional spaces in deep neural networks. These unobserved regions of the AI decision-making process are prone to *normal accidents* – a type of “inevitable” accident that emerges in situations where components are densely connected, tightly coupled, and opaque in their processing [5]. Study of this field began with accidents such as Three-Mile Island, but AI technologies embody similar risks. Finding and exploiting these weaknesses to induce defective behavior will become a permanent feature of military strategy [6].

This human/AI partnership is likely to produce emergent behaviors that are not obvious extensions of current military thinking. This creates a tension between two poles. At one end is the need for systems to be trustworthy. They should predictably do what we believe is the right thing in ethically difficult conditions. At the other end is the need to be responsive and dynamic in unpredictable conditions. In this paper, we develop a framework for examining problems in this nascent area of intelligent warfighting machines.

## II. BACKGROUND

Although the battlespace becomes faster and more complex as information communications technologies improves, the fundamental tactics have been unchanged

for centuries: opposing commanders observe the evolving battlespace, attempt to understand and model the space, and act to produce positive outcomes. Of course, what makes this difficult is that the adversary is doing the same thing, leading to a co-evolving physical and information environment that is difficult to predict with any certainty [7].

These rapidly co-evolving battlespace dynamics are one of the largest obstacles in involving current state-of-the-art machine learning systems. Currently, the best AI is based on enormous networks that are trained for days against massive datasets. The time frames involved in this process do not afford the rapid updates that human interaction requires.

To address both the promise and the risks of adding lethal combat capabilities to AI systems, we need to establish a development framework that emphasizes human control over the behaviors of such systems, regardless of how sophisticated they become. At the core, we believe that these aspects of human control must include the following:

**Interactivity:** Users need to be able to explore and adjust the behavior of the system to confirm changes that they made and validate that the system exhibits a more “correct” behavior.

**Transparency:** Although intelligent machines may never truly be able to explain their actions, they should be able to reveal the sources from which they learned any particular behavior.

**Resiliency:** Intelligent systems cannot be brittle. They must handle overload conditions gracefully and recover quickly. They must be able to indicate when they are operating with low confidence, and they cannot simply freeze.

### III. LITERATURE REVIEW

The goal of adding AI to the battlespace is to augment humans decision-making, but, adding AI to this decision-making process would have ramifications that need to be considered carefully. If AI systems are effective, pressure to increase the level of assistance to the warfighter would be inevitable. Continued success would mean gradually pushing the human out of the loop, first to a supervisory role and then finally to the role of a “killswitch operator” monitoring an always-on LAWS [8].

We see four relevant areas of work that address aspects of this problem space:

**Cybersecurity:** the virtual counterpart to the physical weapons systems

**Hand-to-hand combat:** a proxy for thinking about multiple adversarial AI systems of equal capability engaged in combat

**Machine learning research:** how current state-of-the-art AI systems can responsively and interactively update their states

**Military strategy:** how these systems must operate in the problem space

In Section III-A, we examine the current state of the art in cyberdefense, the limits of a defense-only strategy, and the emerging argument for cyber-counterattacks, including concerns about automation. Using that as a technological frame, we discuss hand-to-hand combat as a proxy for what happens when there are roughly matched adversarial intelligent systems engaged in extremely dynamic kinetic actions in Section III-B. We then examine how machine learning research addresses the need for interactive, evolving dynamic adaptation in Section III-C. Finally, we fit this information in the frame of military strategy in Section III-D, focusing on the interaction of practical combat considerations and international law, particularly article 51 of the UN charter.

#### A. Cybersecurity and the limits of defense

Cybersecurity controls and countermeasures often employ several machine learning and data mining techniques to uncover signs of misuse, anomaly detection, or hybrid approaches that do both [9]. One of the fundamental issues is the volume and velocity of the information that can be associated with an attack. Machine learning techniques aid network administrators seeking to respond to actual issues rather than false alarms, but this approach also represents a weakness. Zero day attacks, which have no previously known signature, can only be detected using anomaly detection systems. If successful, a zero day attack may be able to exploit a system for considerable periods of time. For example, the FBI has determined that four individuals with Russian support were able to penetrate the Yahoo network for two years, getting subscriber information on 500 million accounts before their activities were detected and stopped [10].

Adaptation or generalization from one attack vector to multiple does not prevent this threat. The more adaptable the classifier is, the more open it is to manipulation adversarial training techniques. In other words, an adversary can learn the latent spaces in the classifier that lead to false results. This can be exploited to overwhelm the system with false positives for benign vectors while simultaneously rendering dangerous vectors less detectable [11].

Academic computer security research is overwhelmingly oriented towards detecting and blocking cyberattacks. Regardless of whether the detection scheme is recognition or anomaly-based, all these approaches rest on the fundamental assumption that cybersecurity is *passive*. Systems wait for attacks, identify them as fast as possible, and determine the best course(s) of action and respond, often within milliseconds [9].

There is a growing awareness among cybersecurity professionals that there may be a need for active defense as well, particularly if the cyberattack results in damage to critical national infrastructure. Active defense may be

both appropriate and effective in eliciting cooperation. Axelrod showed in 1981 that tit-for-tat responses to aggression were a robust and effective response [12]. Two considerations are crucial: The first is whether counter attack makes sense as a strategy [13], [14]. The second is how fast to respond. Current government processes associated with responding to kinetic attacks are too slow for responsive cyberattacks [15].

This highlights the fundamental issue in the use of AI systems in weapons systems, whether virtual or physical. The feedback loop between ever-increasing technical capability and the political awareness of the decreasing time window for reflective decision-making drives technical evolution towards always-on, automated, reflexive systems [8]. This pressure needs to be addressed openly and transparently in any system design.

### B. Hand-to-hand Combat

A useful analogy to the evolutionary path that we see happening is individual unarmed combat. This is the only example where we can observe a proxy of similarly matched intelligent systems interacting using the affordances of force [16]. This model is only effective for considering evenly matched AI combat systems because in a mismatch, the odds of a rout are high. Asymmetric encounters are important as well but not the focus of this model.

Most human combat consists of two phases: an *assessment phase* where the adversaries evaluate each other before striking. This phase is more analytical and less reflexive. The adversaries evaluate one another in a highly dynamic state while moving in tandem. They employ past training to generate a plan of action while continuously attempting to lead the opponent into making incorrect assessments, increasing the size and complexity of the problem space each opponent has to consider [17]. The second phase is a *kinetic phase* involving rapid strike, defense, and counterstrike. For these actions to be effective, they have to be reflexive. Any reflective thinking slows down the action, exposing vulnerability.

These two processes roughly correspond to Kahneman's mechanisms for human cognition [18]. Kahneman's System 1 is reflexively responding to a stimulus, whereas his System 2 is conscious calculation. System 1 can be "trained" to respond with seemingly conscious calculation. A good fighter can produce complex sequences of reflexive action in response to combat cues. For example, *The Book of Five Rings* [16], a canonical work describing traditional Japanese martial arts, describes the *Crimson-Leaves Strike*, a trained reflexive action. The first part of the strike is to identify or cause the opponent to lower his guard. This triggers a trained reflex that causes the fighter to strike reflexively at the opening.

We employ this combat model for the entire range of human and AI combat systems, from fully human to fully automated. In all cases, action in the kinetic phase

must be as fast as possible, leaving no time to search for novel solutions. What changes in the transition to AI systems will be the speed and number of dimensions to consider. One can easily imagine human/computer partnerships, where humans become more involved with the assessment phase and less with the kinetic phase. Over time, as AI becomes more capable of reflective and integrative thinking, the human component will have to be eliminated altogether as the speed and dimensionality become incomprehensible, even accounting for cognitive assistance.

### C. Machine Learning

Modern machine learning research is focused on developing huge models that train over even larger datasets, often for days and weeks. Though startlingly effective, these systems struggle to adapt to changing conditions [19]. One method to increase adaptability is called *Transfer Learning* [20], which allows a model optimized and trained on one dataset to be modified and trained on a different, smaller but related dataset. For example, image recognition systems trained to recognize vehicles for a self-driving car application could be adapted to detect and recognize military aircraft.

These kinds of machine learning models contain a weakness. Numerous studies have shown that *adversarial attacks* can cause systems to misclassify examples that are only slightly different from correctly classified examples" [21]. For example, Figure 1 shows that wearing a picture can fool an image classifier [22].

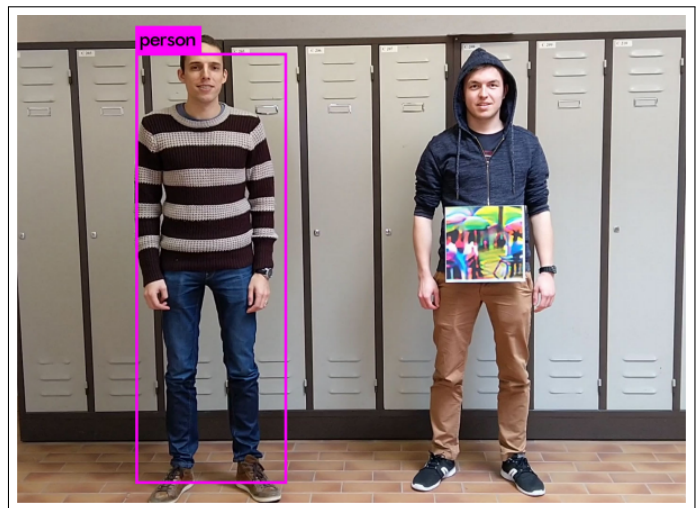


Figure 1. Adversarial Attack Against an Image Classifier [22]

For combat systems, huge models create an inherent risk. Because academic and corporate models are few and often in the public domain, a malicious actor seeking to save development costs can simply download and study them for areas where they can be manipulated to respond incorrectly to a particular set of stimuli [21]. However, transfer learning changes only a small part

of the network, latent space vulnerabilities would likely exist regardless of how the model has been adapted. This places any combat system based on one of these models at risk for undetectable exploitation. Models developed in secure environments based on real and simulated data may be significantly more secure from exploitation, but with a correspondingly higher cost to develop.

Neural network systems can learn as they interact with an environment, through *Reinforcement Learning* (RL) [23]. This technique allows a system to explore a problem space with respect to an evaluation function that can score the system's behavior. Such systems can learn to move in simulated environments, play games, and operate robots. Adversarial versions, where one RL system is scored by how well it is competing against another, currently form the basis for state of the art results in such games as Chess and Alpha Go. Training time for adversarial gamespaces is significant. DeepMind's AlphaGo Zero took approximately 40 days to train, including self-play of 29 million games [24]. If a less exhaustive exploration of the data space is acceptable, RL systems can be guided through their learning process by humans. This technique significantly reduces learning time because humans can often recognize incorrect behavior long before the machine can [25].

#### D. Military Strategy

The role of military systems varies by context and can be deeply complicated. Broadly, forces can be in a peacetime state, at the boundary between peace and war, and in armed conflict [26]. Article 51 of the UN Charter specifies that member states may always act in self defense, but there are now decades of precedent that specify how that right may be interpreted. Further, as we have seen in the cyberdefense section, the definition of what constitutes force is changing.

Understanding military contexts matters because of how we develop trust in the automated systems we use regularly. Continuous interaction builds trust incrementally until we implicitly hand off responsibility to the system and direct our attention elsewhere [27], [28]. This user bias is recognized as the root cause in numerous accidents. If your early warning system alerts for an incoming attack, the pressure to trust the alert and respond is powerful. System trust has been a contributing factor in fratricidal battlefield losses involving Patriot Missiles in Iraq [29]. Once any level of handoff occurs, the transfer of control should generally assumed to be total.

When a soldier or commander makes a judgment call about the use of force, an explicit set of procedures, orders, or other judgment calls precede the first shot being fired. This chain of accountability reduces the dimensionality of the problem space of the decision. But intelligent machines do not proceed similarly. Once a higher-level decision is made, an AI weapon system can effect orders

almost immediately, which is part of the attraction of integrating AI into conflict operations.

When we start to include AI in weapons systems at any level, the system will need to know the context it is acting in and the rules of engagement for that context. Designers cannot foresee all potential contexts, so the system will need to be adaptable and transparent. A human-intensive example of this approach is used in the Aegis combat system, where the parameters for semiautomated and fully automated behavior of the system is prepared in advance of each deployment by the ship's doctrine review board, a diverse mix of officers and senior enlisted crew that review all aspects of the upcoming mission before preparing a suite of behavior "packages" that can be activated by the captain [8].

Context is critical. Training exercises may look like war, but they are actually between allies. A cold war may look like peace, but it isn't exactly. Any intelligent system (human, human/machine, or machine) must be aware of these and other complicating concerns. A weapons system that either directly or indirectly starts a conflict will need to act in accordance with international law or risk implicating its makers and the government that activated it in war crimes. Adversaries know this and will try to use that weakness against any intelligent system.

The U.S. Military does not have and significant technological advantage in this space. China, to take one example, views the U.S. as highly vulnerable in cyberwar and is working to cement its potential advantage. Pressure to develop systems that can effectively grapple with our adversaries across multiple domains, dimensions, and timeframes will be extremely high. All sides are equally pressured to gain superiority, and as such the inevitability of fully automated, always on systems should be seriously considered in all aspects of AI integration.

## IV. OPPORTUNITIES AND CHALLENGES

Military adaptation of commercially or academically trained models contains inherent risks. However, these risks highlight potential opportunities for development that are distinct from the focus of commercial and academic communities. In particular, we identify the following opportunities and challenges: (A) Offline latent space hacking by adversaries; (B) Incorporating legal and ethical constraints into training a model; (C) Mapping, traceability, and transparency of inputs and outputs; and (D) Avoiding dangerous predictability.

### A. Offline Latent Space Hacking by Adversaries

Learning to exploit regions in the available latent space of large models should be explored in depth. Less well supported actors may take advantage of commercial or academic models in an attempt to gain high military impact for low cost and effort. Determining how to thwart, for example, a terrorist organization turning

a facial recognition model into a targeting system for exploding drones is certainly a prudent move.

Technologies such as evolutionary development of model structures in a reinforcement learning environment can create a framework to support the development of unique network structures that cannot be predicted by an adversary [30]. Further, such a framework can be induced to create different networks that address the same sets of problems, making it possible to generate a diverse set of systems providing redundancy and resilience.

### B. Incorporating Legal and Ethical Constraints

Modern AI/ML systems reflect the data used to train them. In commercial and academic systems, data often reflects unconscious biases that emerge in the trained behavior of the system [31]. This type of behavior only becomes more dangerous when it is connected with weapon systems. We need to develop techniques that allow us to train models that have the appropriate doctrine “baked in” so that they can operate appropriately in contexts ranging from war games in peacetime to escorting an adversarial emissary to a peace conference in wartime. Although some research exists for encoding legal and ethical considerations (e.g., using evolutionary approaches [32]), publications in this area are rare.

### C. Mapping, Traceability, and Transparency

It is our strong belief that intelligent weapons systems of the future will move and think at machine speed. This disproportionate capability and the inevitable system trust human operators will place in these machines means that most if not all lethal and sub-lethal interactions will only be analyzable in hindsight [33]. Military weapon system models must be built to support a recorded mapping of inputs that can be traced to actions or recommendations. This level of transparency is crucial for post-incident analysis, validation, and retraining.

### D. Lack of diversity

Because the creation of models is complicated and time consuming, few commercial models address substantively similar tasks at the same level of sophistication. Indeed, there is often only one “best” model for any set of data. In either a cybersecurity or military context, this sort of monoculture represents a predictable single point of failure. Diverse models need to be developed to address tasks redundantly, and they must be regularly revisited, modified, or rebuilt to ensure any adversary obtaining a system using one model will not be able to rely completely on it in the face of battle. This challenge will likely not be met in the academic or commercial community, where raw performance improvements determines success, not survivability and ruggedness.

## V. THE ROLE OF HUMANS

Because of the rate of technological development in the AI/ML space, we believe that the role of humans in combat systems, *barring regulation through treaty*,<sup>1</sup> will become more peripheral over time. As such, it is critical to ensure that our design decisions and the implementations of these designs incorporate the values that we wish to express as a national and global culture.

### A. Human-in-the-loop

The starting point for many intelligent systems begins with the tight integration of human and machine in the weapon system. For example, missiles announce when they have a lock, increasing the capability of the warfighter and leave little ambiguity as to what the weapons system will do once the trigger is pulled. If fault has to be found, it will be the human that must bear the responsibility.

But in more ambiguous circumstance, such as friend-or-foe identification (IFF), the data used to make the calculation will be in the possession of the system, not the user. Imagine a case where an IFF transponders have been known to be spoofed by the adversary, and a large, slow moving aircraft identifying as civilian has been detected on what seems to be a hostile approach, and only a short time to make a decision. There are four presentations: (1) the system can declare that it has identified the aircraft as hostile and provide a lock; (2) the system can declare the aircraft as friendly and open a channel to warn; (3) the system can present a set of ranked recommendations and provide a set of options to the user; or (4) the system simply displays the raw information.

The third option may seem to be the best embodiment of the human-in-the-loop philosophy, but it discounts the effects of system trust. The user may spend some time evaluating the list of options the first few times, but if the system places the correct option at the top of the list often enough, the human user will begin to simply select that option. This effect is exacerbated with time pressures. The human simply becomes a rubber stamp.

If, however, the system is trained by a set of known individuals, and the provenance of the system ranking can be traced back to its “mentor.” Mentors could train the model in the context of the current deployment and be known to the user. The machine then incorporates rules of engagement that are related to the particular deployment through this human interaction. Further, weights that are accumulated from these interactions can be brought back and integrated into the models, allowing them to evolve with respect to the current realities of a given battlespace.

<sup>1</sup> And assuming that AI/ML systems are not advanced enough to autonomously incorporate the risk and consequences of violating international treaty into their decision-making.

## B. Human-on-the-loop

As human-in-the-loop systems advance, the system with less need to rely on human decision-making to achieve results will begin to dominate. Thus, humans will be relegated to offline analysis and improvement of AI strategies during training. Work is already being done in this space commercially. Examples of what is essentially human-on-the-loop architectures are regularly explored now in StarCraft competitions [19]. From a machine learning perspective, the difference between a StarCraft 2 environment an autonomous Aegis battlegroup is one of scale and consequences. Though, human-on-the-loop systems have less interaction in real-time, integration of the mentor architecture described above may be possible. This integration would depend on the creation of offline wargames and simulations that can be played at rates slow enough to elicit meaningful training from expert human cognition.

Continually running human-led scenarios offline increases the odds that the trained system reflects the realities of the battlespace [8]. Separate classifier systems may also be able to catalog adversary behaviors near the boundaries of the current trained responses. These boundaries might be detected by looking at the behavior of the systems themselves and recognizing when they are making decisions among multiple options with divergent potential outcomes.

## C. Human initiated

An extension of the human-on-the-loop approach is the human-initiated “fire and forget” approach to battlefield AI. Once the velocity and dimensionality of the battlespace increase beyond human comprehension, human involvement will be limited to choosing the temporal and physical bounds of behavior desired in an anticipated context. Depending on the immediacy or unpredictability of the threat, engaging the system manually at the onset of hostilities may be impossible. Rather, these systems would need to be activated before the onset of hostilities. Activation alone could be an extremely consequential. Once activated, shutting the system down for any reason may be interpreted by an adversarial AI as a weakness to be exploited. Predicting with confidence how a conflict between two equally matched, highly capable AI systems would unfold once started may be impossible. These conflicts would have to be simulated extensively against a wide variety of adversaries to have any confidence that their behavior would align with our national values.

## D. Post-hoc forensics

Given a battlespace so overwhelming that humans cannot manually engage with the system, the human role will be limited to post-hoc forensic analysis, once hostilities have ceased, or treaties have been signed and

implemented.<sup>2</sup> To this end, these systems will need the maximal amount of provenance of input data, alternatives considered, and actions selected with sufficiently high-fidelity recordings that any error or inexplicable behavior can be meaningfully interpreted [34]. Recent work that applies an approach of this sort is the concept of an “activation atlas” [33] that can show areas of conceptual conflict within a model. These approaches are nascent, but need to be vigorously explored and developed.

## VI. DISCUSSION AND SUMMARY

Tight integration of AI into the kill chain is not a decision to be taken lightly. Particularly with respect to LAWS, our policy should be: *not until they can outperform human/MI collaboration, including making ethically acceptable choices* [3]. Prior to this, deepening our understanding of the roles, risks, and rewards of integrating AI into the killchain can better define the problem space of the solutions that we contemplate.

Many of the solutions needed to address concerns identified herein are likely also useful in commercial or civilian settings. The dual use applications for responsive, diverse and adaptable AI seem numerous, ranging from near-term problems such as autonomous vehicles in varied environments, to locally personalized, private, and secure AI assistants. Novel solutions in these area could be useful in many sectors of today’s economy, and might kickstart new developments in AI that would in turn inform and improve the development of systems focused on military problems. Also, violence is not limited to nation-state level conflict. Some of the challenges outlined herein may ultimately arise in civilian contexts.

For industry, the incentives for developing adaptable AI are currently low, but here is an opportunity for the defense community to once more contribute to technological advances in ways that benefit the broader population, much in the same way that the development of the B-52 “spilled over” into the development of the Boeing 707 and subsequent commercial jet aircraft [35].

Ethically-aware AI support systems would be useful across many human domains, including law, diplomacy, and negotiations between multiple parties. As war has been often described as the application of force as a replacement for diplomacy, perhaps AI-enhanced diplomacy can reduce the risk of a hyper-accelerated AI war.

Finally, this paper assumes that AI does not ultimately “run away” from human ability to control it. Serious philosophers view this sort of “technological singularity” as a realistic scenario, and we believe it must be addressed no later than widespread use of human-on-the-loop systems.

This paper presents a framework for understanding the integration of AI and ML into military weapons

<sup>2</sup> Diplomacy in such an environment seems like a nearly incomprehensible challenge and deserves extensive research on its own.



systems. The steps we take may be on a path to human obsolescence as combatants, and the decision to proceed on this path should be well informed and involve all members of our societies. Given that other players are deeply engaged in the weaponization of AI, we would be foolish not to research and experiment with the development of highly capable systems. But we believe that a better answer may be to support regulation and prohibition because, like chemical and biological weapons, for weaponized AI, “the only winning move is not to play.”

## REFERENCES

- [1] C. Ferraris and R. Carveth, “NASA and the Columbia Disaster: Decision-making by Groupthink?” in *Proceedings of the 2003 Association for Business Communication Annual Convention*, 2003, p. 12.
- [2] K. Clark, “The GPS: A fatally misleading travel companion,” Jul. 2011. [Online]. Available: <http://www.npr.org/2011/07/26/137646147/the-gps-a-fatally-misleading-travel-companion>
- [3] W. Carter, *Chinese Advances in Emerging Technologies and their Implications for US National Security*. Washington, DC: CSIS, 2018.
- [4] E. B. Kania, *Battlefield Singularity: Artificial Intelligence, Military Revolution, and China's Future Military Power*. Center for a New American Security, 2017.
- [5] C. Perrow, *Normal Accidents: Living with High Risk Technologies*. Princeton University Press, 2011.
- [6] J. Staff, “Joint publication 3-13 information operations incorporating change 1,” 2014.
- [7] A. B. Dahl, “Command dysfunction: Minding the cognitive war,” Air University Maxwell AFB AL School of Advanced Airpower Studies, Tech. Rep., 1996.
- [8] P. Scharre, *Army of none: Autonomous weapons and the future of war*. WW Norton & Company, 2018.
- [9] A. L. Buczak and E. Guven, “A survey of data mining and machine learning methods for cyber security intrusion detection,” *IEEE Communications Surveys Tutorials*, vol. 18, no. 2, pp. 1153–1176, Secondquarter 2016.
- [10] FBI, “Charges announced in massive cyber intrusion case,” Mar 2017. [Online]. Available: <https://www.fbi.gov/news/stories/charges-announced-in-massive-cyber-intrusion-case>
- [11] Z. Katzir and Y. Elovici, “Quantifying the resilience of machine learning classifiers used for cyber security,” *Expert Systems with Applications*, vol. 92, pp. 419–429, 2018.
- [12] R. Axelrod and W. D. Hamilton, “The evolution of cooperation,” *science*, vol. 211, no. 4489, pp. 1390–1396, 1981.
- [13] J. Kallberg and R. A. Burk, “The flaw of immediate cyber counter strikes,” *Strategic Analysis*, vol. 41, no. 5, pp. 510–514, 2017.
- [14] E. T. Jensen, “Computer attacks on critical national infrastructure: A use of force invoking the right of self-defense,” *Stan. J. Int'l L.*, vol. 38, p. 207, 2002.
- [15] T. Grant, “Speeding up parliamentary decision making for cyber counter-attack,” in *ICMLG 2017 5th International Conference on Management Leadership and Governance*. Academic Conferences and publishing limited, 2017, p. 152.
- [16] M. Miyamoto, *The Book of Five Rings*. Kodansha International, 1645.
- [17] B. Hart, *Strategy*. Editorial Benei Noaj, 2009. [Online]. Available: [https://books.google.com/books?id=\\_D\\_wPgAACAAJ](https://books.google.com/books?id=_D_wPgAACAAJ)
- [18] D. Kahneman, *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011.
- [19] P.-A. Andersen, M. Goodwin, and O.-C. Granmo, “Deep RTS: A Game Environment for Deep Reinforcement Learning in Real-Time Strategy Games,” in *Conference on Computational Intelligence and Games (CIG)*. IEEE, 2018, pp. 1–8.
- [20] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [21] I. J. Goodfellow, J. Shlens, and C. Szegedy, “Explaining and harnessing adversarial examples,” *arXiv preprint arXiv:1412.6572*, 2014.
- [22] S. Thys, W. Van Ranst, and T. Goedemé, “Fooling automated surveillance cameras: adversarial patches to attack person detection,” *arXiv preprint arXiv:1904.08653*, 2019.
- [23] H. B. Suay and S. Chernova, “Effect of human guidance and state space size on interactive reinforcement learning,” in *RO-MAN*. IEEE, 2011, pp. 1–6.
- [24] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, no. 7676, p. 354, 2017.
- [25] L. Fridman, “Human-centered autonomous vehicle systems: Principles of effective shared autonomy,” *arXiv preprint arXiv:1810.01835*, 2018.
- [26] M. J. Sklerov, “Solving the dilemma of state responses to cyberattacks: A justification for the use of active defenses against states who neglect their duty to prevent,” *Mil. L. Rev.*, vol. 201, p. 1, 2009.
- [27] G. Hacıyakupoglu and W. Zhang, “Social media and trust during the gezi protests in turkey,” *Journal of Computer-Mediated Communication*, vol. 20, no. 4, pp. 450–466, 2015.
- [28] K. L. Mosier, L. J. Skitka, M. D. Burdick, and S. T. Heers, “Automation bias, accountability, and verification behaviors,” in *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, vol. 40, no. 4. SAGE Publications Sage CA: Los Angeles, CA, 1996, pp. 204–208.
- [29] J. K. Hawley, “Not by widgets alone: The human challenge of technology-intensive military systems,” *Armed Forces Journal*, pp. 24–28, 2011.
- [30] P. J. Angeline, G. M. Saunders, and J. B. Pollack, “An evolutionary algorithm that constructs recurrent neural networks,” *IEEE transactions on Neural Networks*, vol. 5, no. 1, pp. 54–65, 1994.
- [31] T. Bolukbasi, K.-W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai, “Man is to computer programmer as woman is to homemaker? debiasing word embeddings,” in *Advances in neural information processing systems*, 2016, pp. 4349–4357.
- [32] A. R. Honarvar and N. Ghasem-Aghaee, “An artificial neural network approach for creating an ethical artificial agent,” in *2009 IEEE International Symposium on Computational Intelligence in Robotics and Automation-(CIRA)*. IEEE, 2009, pp. 290–295.
- [33] S. Carter, Z. Armstrong, L. Schubert, I. Johnson, and C. Olah, “Activation atlas,” *Distill*, 2019.
- [34] X. Liu, X. Wang, and S. Matwin, “Improving the interpretability of deep neural networks with knowledge distillation,” in *Proceedings of the International Conference on Data Mining Workshops (ICDMW)*, Nov 2018, pp. 905–912.
- [35] S. Kotha, “Spillovers, spill-ins, and strategic entrepreneurship: America’s first commercial jet airplane and boeing’s ascendancy in commercial aviation,” *Strategic Entrepreneurship Journal*, vol. 4, no. 4, pp. 284–306, 2010.