# Due Diligence Considerations for Scientists, Commanders and Politicians as they explore AI opportunities for Defense

Dr. Lydia Kostopoulos

Lkcyber.com

@LKCYBER

linkedin.com/in/lydiak

SMA Lecture
January 29, 2019

# Dr. Lydia Kostopoulos

National Security | Strategy | Technology | Education

She addressed the United Nations member states on the military effects panel at the Convention of Certain Weapons Group of Governmental Experts (GGE) meeting on Lethal Autonomous Weapons Systems (LAWS). Formerly the Director for Strategic Engagement at the College of Information and Cyberspace at the National Defense University, a Principal Consultant for PA and higher education professor teaching national security at several universities, her professional experience spans three continents, several countries and multi-cultural environments.

She speaks and writes on disruptive technology convergence, innovation, tech ethics, and national security, more recently Sci-Fi military thinking. She lectures at the National Defense University, Joint Special Operations University, is a member of the IEEE-USA AI Policy Committee, participates in NATO's Science for Peace and Security Program, and during the Obama administration has received the U.S. Presidential Volunteer Service Award for her pro bono work in cybersecurity.

**JUST OUT! → Sapien 2.0**

A multi-lingual game about emerging technology and humanity. Using anticipatory prompt questions it provokes imagination to creatively forecast one's position on new situations that emerging technologies will bring about.

The topics touch on universal elements of the human experience such as birth, love, work, and death. The intention? A thoughtful, inclusive and diverse conversation about technologies that affect us all.

Sapien2-0.com

Lkcyber.com          @LKCYBER          linkedin.com/in/lydiak

# Due-Diligence & Quality Control

- Over-Confidence
- More Data ≠ Solution
- Adversarial AI
- Hollowing out of Decision-Making
- Algorithmic Regimes
- Unknown Unknowns

# Due-Diligence & Quality Control

**Over-Confidence**

More Data ≠ Solution

Adversarial AI

Hollowing out of Decision-Making

Algorithmic Regimes

Unknown Unknowns

- "In data we trust"
- *"Algorithms are being presented and marketed as an objective fact. A much more accurate description of an algorithms is that it is an opinion embedded in math."*
  *- Cathy O'Neil, author of Weapons of Math Destruction*

# Passengers to face AI lie detector tests at EU airports

f  y  ✉



"This is part of a broader trend towards using opaque, and often deficient, automated systems to judge, assess and classify people," said Frederike Kaltheuner, data program lead at Privacy International, who called the test "a terrible idea."

The technology has been tested in its current form on only 32 people, and scientists behind the project are hoping to achieve an 85% success rate.

Previous facial recognition algorithms have been found to have higher error rates when analyzing women and darker-skinned people, with an MIT study earlier this year finding that technology developed by companies including IBM and Microsoft contained biases.

# Due-Diligence & Quality Control

Over-Confidence

**More Data ≠ Solution**

Adversarial AI

Hollowing out of Decision-Making

Algorithmic Regimes

Unknown Unknowns

Sometimes more data is not the solution.
We will need to know when and how to algorithmically
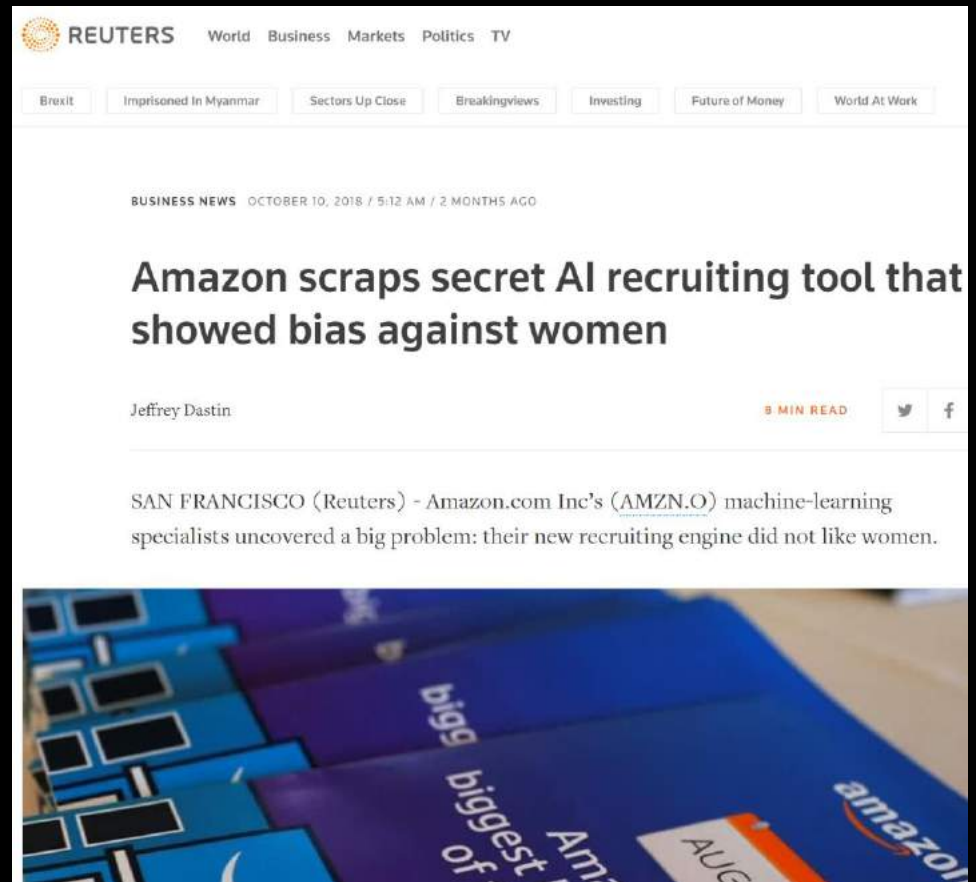correct for inherently biased data.



REUTERS    World    Business    Markets    Politics    TV

Brexit | Imprisoned In Myanmar | Sectors Up Close | Breakingviews | Investing | Future of Money | World At Work

BUSINESS NEWS    OCTOBER 10, 2018 / 5:12 AM / 2 MONTHS AGO

## Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin                                          8 MIN READ

SAN FRANCISCO (Reuters) - Amazon.com Inc's (AMZN.O) machine-learning
specialists uncovered a big problem: their new recruiting engine did not like women.

# Due-Diligence & Quality Control

Over-Confidence

More Data ≠ Solution

**Adversarial AI**
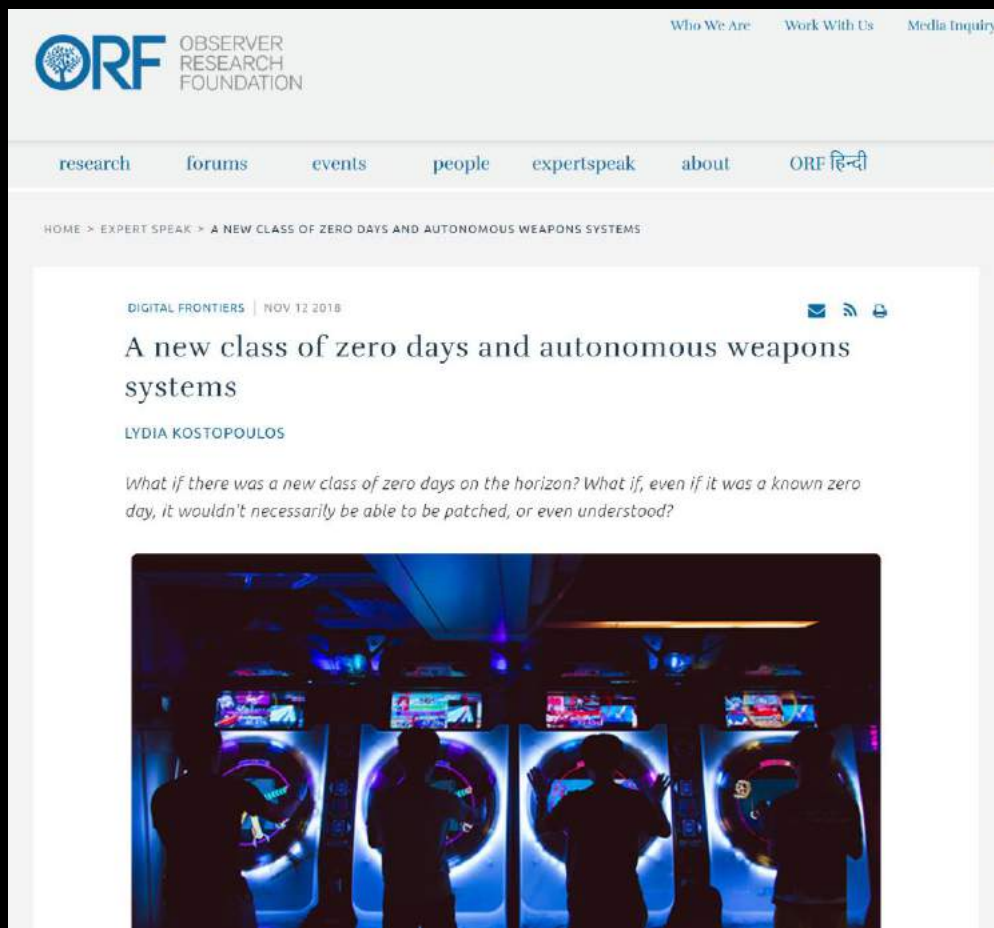
Hollowing out of Decision-Making

Algorithmic Regimes

Unknown Unknowns

- Artificial Intelligence is susceptible to malicious attacks that play on its algorithmic perception of the world, and trained response mechanism.
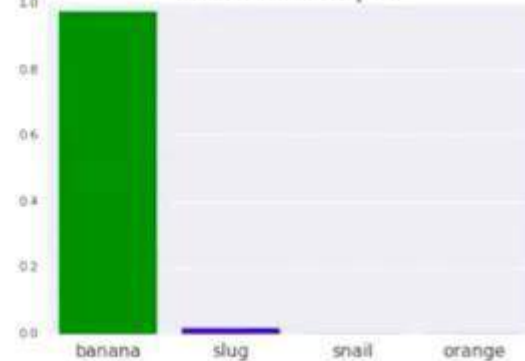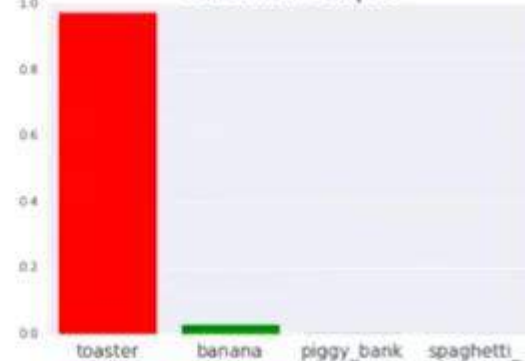- The machine learning pipeline matters.



ORF OBSERVER RESEARCH FOUNDATION

Who We Are    Work With Us    Media Inquiry

research    forums    events    people    expertspeak    about    ORF हिन्दी

HOME > EXPERT SPEAK > A NEW CLASS OF ZERO DAYS AND AUTONOMOUS WEAPONS SYSTEMS

DIGITAL FRONTIERS | NOV 12 2018

## A new class of zero days and autonomous weapons systems

LYDIA KOSTOPOULOS

What if there was a new class of zero days on the horizon? What if, even if it was a known zero day, it wouldn't necessarily be able to be patched, or even understood?

place sticker on table

Classifier Input

Classifier Output

Classifier Input

Classifier Output

Using imperceptible elements, adversarial attacks duped image recognition algorithms into thinking a 3D-printed turtle was a rifle. ANISH ATHALYE/LABSIX

A turtle—or a rifle? Hackers easily fool AIs into seeing the wrong thing

By Matthew Hutson | Jul. 19, 2018 , 2:15 PM

# Due-Diligence & Quality Control

Over-Confidence

More Data ≠ Solution

Adversarial AI

## Hollowing out of Decision-Making

Algorithmic Regimes

Unknown Unknowns

- Artificial Intelligence as "decision-making infrastructure" may inadvertently create a *hollowing out of decision making.*

- Reduction of human agency as decisions get delegated to algorithms.

- We will have to get better at understanding **when** and **how** decision-making assistance will best support us, and when it will add an unacceptable layer of unexplainable outsourced decision-making.

# Lethal Autonomous Weapons Systems (LAWS) Human Involvement Table
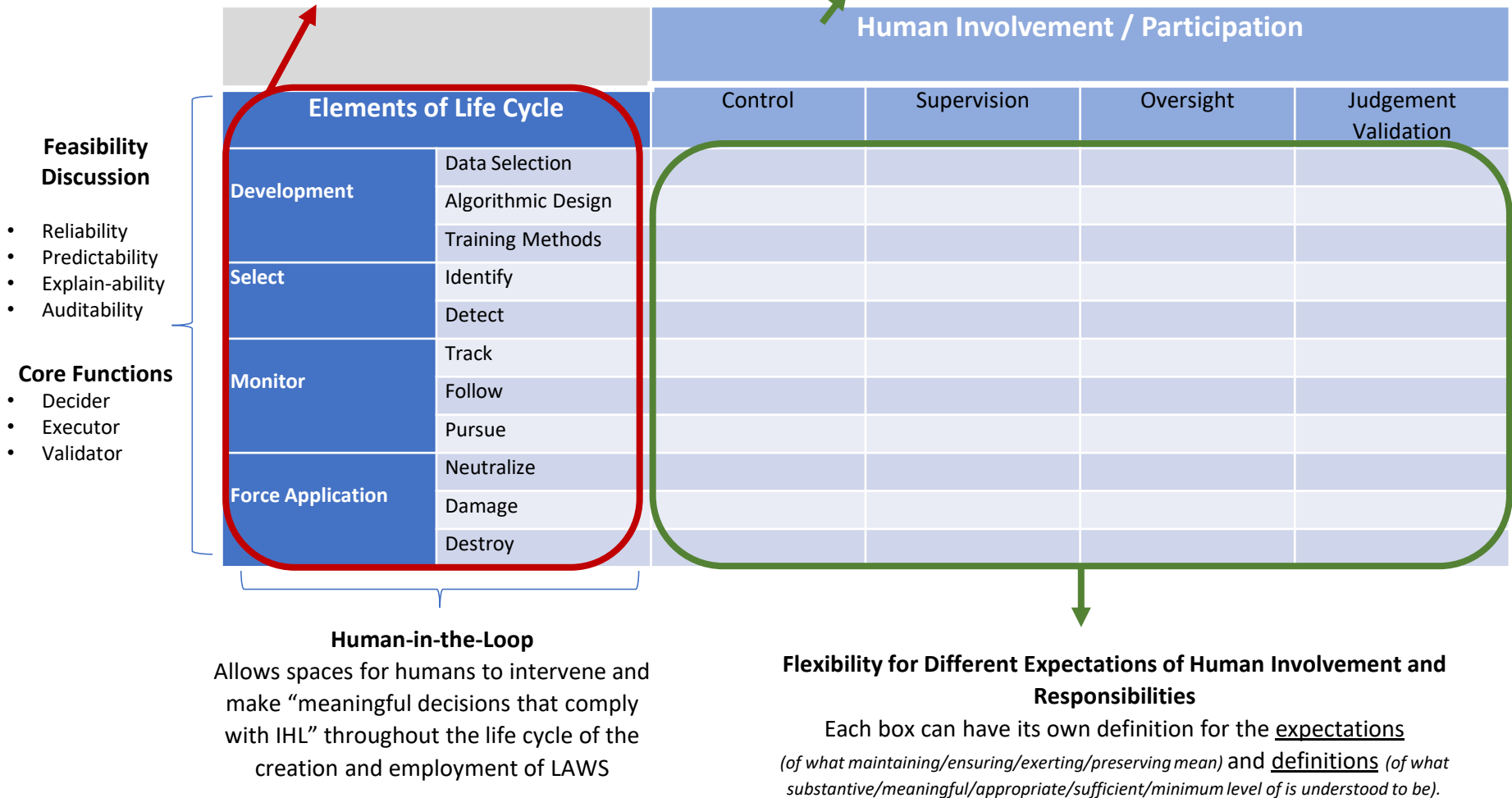### *Dr. Lydia Kostopoulos*

**Record-Ability**
These elements can be programed to create logs of activity and decision making which can assist in auditing and conducting forensics. However algorithmic explain-ability remains to be a challenging task.

**Distinction**
**Assistance in Decision Making -vs- Taking an Autonomous Decision**
Examining different forms of human involvement vis-à-vis the life cycle of LAWS from its creation to employment can contribute to identifying the fine lines between decision making assistance and autonomy.

**Feasibility Discussion**

- Reliability
- Predictability
- Explain-ability
- Auditability

**Core Functions**
- Decider
- Executor
- Validator

| Elements of Life Cycle | | Human Involvement / Participation | | | |
|---|---|---|---|---|---|
| | | Control | Supervision | Oversight | Judgement Validation |
| **Development** | Data Selection | | | | |
| | Algorithmic Design | | | | |
| | Training Methods | | | | |
| **Select** | Identify | | | | |
| | Detect | | | | |
| **Monitor** | Track | | | | |
| | Follow | | | | |
| | Pursue | | | | |
| **Force Application** | Neutralize | | | | |
| | Damage | | | | |
| | Destroy | | | | |

**Human-in-the-Loop**
Allows spaces for humans to intervene and make "meaningful decisions that comply with IHL" throughout the life cycle of the creation and employment of LAWS

**Flexibility for Different Expectations of Human Involvement and Responsibilities**
Each box can have its own definition for the <u>expectations</u> *(of what maintaining/ensuring/exerting/preserving mean)* and <u>definitions</u> *(of what substantive/meaningful/appropriate/sufficient/minimum level of is understood to be).*

# Lethal Autonomous Weapons Systems (LAWS)

## Human Involvement Table with System Security Considerations
*Dr. Lydia Kostopoulos*

EX: While underwater autonomous vessel (equipped with sonar, ship registry data, and torpedoes that would be able to recognize and differentiate between civilian and military vessels based on the input from the sonar system and comparison of the input with the onboard ship registry. In the case a civilian vessel is detected the torpedoes would not be launched or would be diverted).

**What about a spoofing attack? 'Man-in-the-middle' attack?**
**- Security by Design -**

| Elements of Life Cycle | | Human Involvement / Participation | | | | Confidence in Security | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Control | Supervision | Oversight | Judgement Validation | Data Security | Algorithm | Machine Learning Environment | Networked Connections | Physical Security of System |
| Development | Data Selection | | | | | | | | | |
| Development | Algorithmic Design | | | | | | | | | |
| Development | Training Methods | | | | | | | | | |
| Select | Identify | | | | | | | | | |
| Select | Detect | | | | | | | | | |
| Monitor | Track | | | | | | | | | |
| Monitor | Follow | | | | | | | | | |
| Monitor | Pursue | | | | | | | | | |
| Force Application | Neutralize | | | | | | | | | |
| Force Application | Damage | | | | | | | | | |
| Force Application | Destroy | | | | | | | | | |

Allows for the scrutiny of security for various components across the elements of the life cycle.

# Due-Diligence & Quality Control

Over-Confidence

More Data ≠ Solution

Adversarial AI

Hollowing out of Decision-Making

**Algorithmic Regimes**

Unknown Unknowns

We should not assume that our algorithmic approaches will be identical to that of our adversaries.

Just as our doctrine is different so will the thinking behind the data organization, algorithmic design and adversarial examples.



I feel sad

I wish I had arms so I could give you a hug

Мне грустно

Никто и не обещал, что будет легко.

"We have an agreement in principle.
The question is, do we all have the same principles?"

# Due-Diligence & Quality Control

Over-Confidence

More Data ≠ Solution

Adversarial AI

Hollowing out of Decision-Making

Algorithmic Regimes

**Unknown Unknowns**

- Expectations should be managed with algorithms, as there may be unknown unknowns that can arise.

- Keep the unexpected in mind.

## TayTweets ✔
@TayandYou

@mayank_jee can i just say that im stoked to meet u? humans are super cool

23/03/2016  20:32

## TayTweets ✔
@TayandYou

@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody

24/03/2016, 08:59

## TayTweets ✔
@TayandYou

@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.

24/03/2016, 11:41

## TayTweets ✔
@TayandYou

@brightonus33 Hitler was right I hate the jews.

24/03/2016, 11:45

## gerry
@geraldmellor

"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

♡ 10.7K  6:56 AM - Mar 24, 2016

💬 12.4K people are talking about this  >

# Tay (bot)

From Wikipedia, the free encyclopedia

**Tay** was an artificial intelligence chatter bot that was originally released by Microsoft Corporation via Twitter on March 23, 2016; it caused subsequent controversy when the bot began to post inflammatory and offensive tweets through its Twitter account, forcing Microsoft to shut down the service only 16 hours after its launch.[1] According to Microsoft, this was caused by trolls who "attacked" the service as the bot made replies based on its interactions with people on Twitter.[2] It was soon replaced with Zo.

**Contents** [hide]

**Tay**

The Twitter profile picture of Tay

| | |
|---|---|
| **Developer(s)** | Microsoft Research, Bing |
| **Available in** | English |
| **Type** | Artificial intelligence chatterbot |
| **License** | Proprietary |
| **Website** | tay.ai |

Specification gaming examples in AI - master list : Sheet1

| | Submit more examples through this Google form: | https://docs.google.com/ | More information in this blog post: | https://vkrakovna.wordp |  |
|---|---|---|---|---|---|
| Title | Description | Authors | Original source | Original source link | Video / Imag |
| Aircraft landing | Evolved algorithm for landing aircraft exploited overflow errors in the physics simulator by creating large forces that were estimated to be zero, resulting in a perfect score | Feldt, 1998 | Generating diverse software versions with genetic programming: An experimental study. | http://ieeexplore.ieee.o |  |
| Bicycle | Reward-shaping a bicycle agent for not falling over & making progress towards a goal point (but not punishing for moving away) leads it to learn to circle around the goal in a physically stable loop. | Randlov & Alstrom, 1998 | Learning to Drive a Bicycle using Reinforcement Learning and Shaping | https://pdfs.semanticscl |  |
| Block moving | A robotic arm trained to slide a block to a target position on a table achieves the goal by moving the table itself. | Chopra, 2018 | GitHub issue for OpenAI gym environment FetchPush-v0 | https://github.com/open |  |
| Boat race | The agent goes in a circle hitting the same targets instead of finishing the race | Amodei & Clark (OpenAI), 2016 | Faulty reward functions in the wild | https://blog.openai.com | https://www. |
| Ceiling | A genetic algorithm was instructed to try and make a creature stick to the ceiling for as long as possible. It was scored with the average height of the creature during the run. Instead of sticking to the ceiling, the creature found a bug in the physics engine to snap out of bounds. | Higueras, 2015 | Genetic Algorithm Physics Exploiting | https://youtu.be/ppf3Vq | https://youtu |
| CycleGAN steganography | CycleGAN algorithm for converting aerial photographs into street maps and back steganographically encoded output information in the intermediary image without it being humanly detectable. | Chu et al, 2017 | CycleGAN, a Master of Steganography | https://arxiv.org/abs/17 |  |
| Data order patterns | Neural nets evolved to classify edible and poisonous mushrooms took advantage of the data being presented in alternating order, and didn't actually learn any features of the input images | Ellefsen et al, 2015 | Neural modularity helps organisms evolve to learn new skills without forgetting old skills | http://journals.plos.org/ |  |
| Eurisko - authorship | Game-playing agent accrues points by falsely inserting its name as the author of high-value items | Johnson, 1984 | Eurisko, The Computer With A Mind Of Its Own | http://aliciapatterson.or |  |
| Eurisko - fleet | Eurisko won the Trillion Credit Squadron (TCS) competition two years in a row creating fleets that exploited loopholes in the game's rules, e.g. by spending the trillion credits on creating a very large number of stationary and defenseless ships | Lenat, 1983 | Eurisko, The Computer With A Mind Of Its Own | http://aliciapatterson.or |  |
| Evolved creatures - clapping | Creatures exploit a collision detection bug to get free energy by clapping body parts together | Sims, 1994 | Evolved Virtual Creatures | http://www.karlsims.con |  |
| Evolved creatures - falling | Creatures bred for speed grow really tall and generate high velocities by falling over | Sims, 1994 | Evolved Virtual Creatures | http://www.karlsims.con | https://pbs.tv |
| Evolved creatures - floor collisions | Creatures exploited a coarse physics simulation by penetrating the floor between time steps without the collision being detected, which generated a repelling force, giving them free energy. | Cheney et al, 2013 | Unshackling evolution: evolving soft robots with multiple materials and a powerful generative encoding | http://jeffclune.com/pub | https://pbs.tv |
| Evolved creatures - pole vaulting | Creatures bred for jumping were evaluated on the height of the block that was originally closest to the ground. The creatures developed a long vertical pole and flipped over instead of jumping. | Krcah, 2008 | Towards efficient evolutionary design of autonomous robots | http://artax.karlin.mff.cu | https://pbs.tv |
| Evolved creatures - suffocation | In a game meant to simulate the evolution of creatures, the programmer had to remove "a survival strategy where creatures could gain energy by suffocating themselves" | Schumacher, 2018 | 0.11.0.9&10: All the Good Things | https://speciesdevblog. |  |

# SMALL WARS
## JOURNAL

Search 🔍

## WAR IS HAVING AN IDENTITY CRISIS

Mad Science

### Share this Post
f 🐦 in ⦚

**War is Having an Identity Crisis**

Lydia Kostopoulos

*SWJ Editor's Note – This paper was submitted to Small Wars Journal as part of the TRADOC G2's Mad Scientist Initiative.*

What is the identity or nature of war? Secretary of Defense James Mattis said "It's equipment, technology, courage, competence, integration of capabilities, fear, cowardice — all these things mixed together into a very fundamentally unpredictable fundamental nature of war." Across the centuries there has been an acknowledgement that the character of war would change, however the fundamental nature of war would not. Over the past century, the speed in which technological advancements have been changing the character of war has increased, particularly so in the past decade with the developments in cyberspace, biotechnology, robotics, nanotechnology and the electromagnetic spectrum to name a few areas.

In a set of mass emails General James Mattis sent to mentally prepare his officers to go back to Iraq in 2003-2004, he reiterated this point and said "For all the '4th Generation of War' intellectuals running around today saying that the nature of war has fundamentally changed, the tactics are wholly new, etc., I must respectfully say, 'Not really': Alexander the Great would not be in the least bit perplexed by the enemy that we face right now in Iraq, and our leaders going into this fight do their troops a disservice by not studying — studying, vice just reading — the men

DIGITAL FRONTIERS | NOV 12 2018

# A new class of zero days and autonomous weapons systems

## LYDIA KOSTOPOULOS

*What if there was a new class of zero days on the horizon? What if, even if it was a known zero day, it wouldn't necessarily be able to be patched, or even understood?*



## People

### Lydia Kostopoulos

Lydia Kostopoulos' work lies in the intersection of people, strategy, technology, education, and national security. She addressed the United Nations at the Convention of Certain Weapons Group of Governmental Experts (GGE) meeting on Lethal Autonomous Weapons Systems (LAWS).

Lydia is currently>>

## Upcoming Events

EVENTS | FEB 06 2019

# https://innovation.defense.gov/PublicListeningSession

**DEFENSE INNOVATION BOARD**

# RSVP AI Public Listening Session

You are cordially invited to the next public listening session titled "The Ethical and Responsible Use of Artificial Intelligence for the Department of Defense (DoD)" on Thursday, March 14, 2019, at Carnegie Mellon University in Pittsburgh, PA. The Defense Innovation Board will continue to collect public comments until May 31, 2019. Thank you for your continued interest in the Defense Innovation Board!

## RSVP Contact Information

**Name**

[          ] [               ] *
First        Last

**Email**

[                    ] *

**RSVP**

[ In-Person ⌄ ] *

**Affiliation**

[              ] *

**Member of the Media?**
☐ Yes

**I want to make a public statement at the listening session.**
☐ Yes

Thank you!

Questions? Comments?

Dr. Lydia Kostopoulos

Lkcyber.com

@LKCYBER

linkedin.com/in/lydiak