# AI and Security

25.04.24    Alfie Beard, Research Manager

**BT Group**

# AI's Everywhere!

## US experts warn AI likely to kill off jobs – and widen wealth inequality

Economists wary of firm predictions but say advances could create new raft of billionaires while other workers are laid off



New technologies like AI often produce jobs that no one could predict – before the invention of computers, who would have foreseen the job of computer programmer? Photograph: Jaap Arriens/NurPhoto/REX/Shutterstock

ChatGPT is just the latest technology to fuel worries that it will wipe out the jobs of millions of workers, whether advertising copywriters, Wall Street traders, salespeople, writers of basic computer code or journalists.

## Your job is (probably) safe from artificial intelligence

Why predictions of an imminent economic revolution are overstated



May 7th 2023 | SAN FRANCISCO

THE AGE of "generative" artificial intelligence has well and truly arrived. OpenAI's chatbots, which use large-language-model (LLM) technology, got the ball rolling in November. Now barely a day goes by without some mind-blowing advance. An AI-powered song featuring a fake "Drake" and "The Weeknd" recently shook the music industry. Programs which convert

# But no one can agree on it...

DEALBOOK NEWSLETTER

# Washington Confronts the Challenge of Policing A.I.

Lawmakers and Sam Altman, the chief executive of OpenAI, agreed that artificial intelligence should be regulated. How to do that remains an open question.
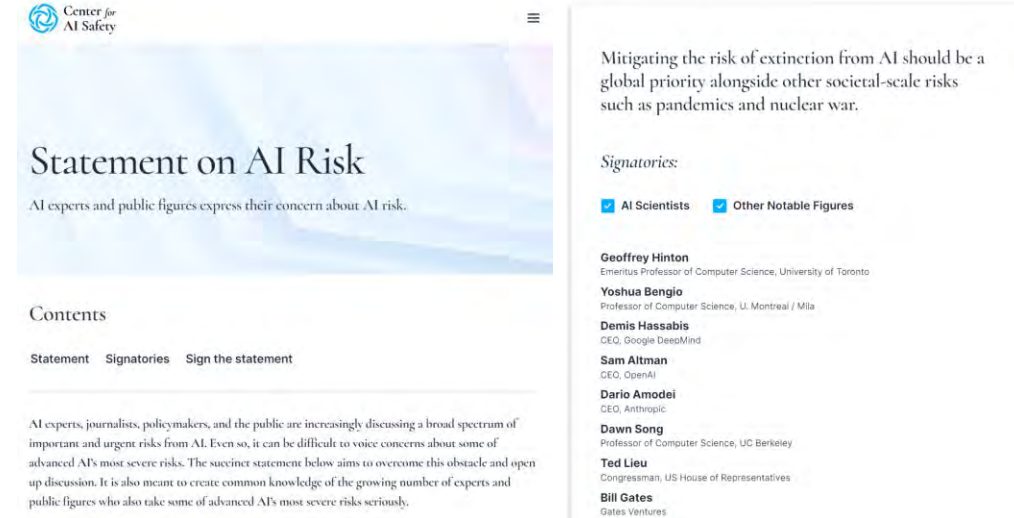
🎁 Give this article   9

By Andrew Ross Sorkin, Ravi Mattu, Bernhard Warner, Sarah Kessler, Michael J. de la Merced, Lauren Hirsch and Ephrat Livni

May 17, 2023

**But no one can agree on it...**

**Yann LeCun** @ylecun

I have spent years trying to understand the risks.
The risk of extinction from AI is close enough to zero that it is negligible in regards to other risks, like climate change and other with which AI can actually help.

6:08 pm · 2 Jun 2023 · **8,768** Views

Center for AI Safety

## Statement on AI Risk

AI experts and public figures express their concern about AI risk.

Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.

*Signatories:*

☑ AI Scientists   ☑ Other Notable Figures

**Geoffrey Hinton**
Emeritus Professor of Computer Science, University of Toronto

**Yoshua Bengio**
Professor of Computer Science, U. Montreal / Mila

**Demis Hassabis**
CEO, Google DeepMind

**Sam Altman**
CEO, OpenAI

**Dario Amodei**
CEO, Anthropic

**Dawn Song**
Professor of Computer Science, UC Berkeley

**Ted Lieu**
Congressman, US House of Representatives

**Bill Gates**
Gates Ventures

### Contents

Statement   Signatories   Sign the statement

AI experts, journalists, policymakers, and the public are increasingly discussing a broad spectrum of important and urgent risks from AI. Even so, it can be difficult to voice concerns about some of advanced AI's most severe risks. The succinct statement below aims to overcome this obstacle and open up discussion. It is also meant to create common knowledge of the growing number of experts and public figures who also take some of advanced AI's most severe risks seriously.

# But no one can agree on it…

# What is AI?

# What is AI?

## AI is a tool

## "The ability for machines to accomplish complex goals"

# BT Group

# Why Now?

# Why Now?



Moore's Law: The number of transistors on microchips doubles every two years

**Increased Computational Capability**

# Why Now?

1500 years human experience in a few hours of training.

*"I can't disguise my satisfaction that it plays with a very dynamic style, much like my own!"*

Garry Kasparov





An animation of the gradient descent method predicting a structure for CASP13 target T1008

**BT Group**

# Where are we now?

# AI Progress

Healthcare

Content Generation

Entertainment

Film

Travel

Agriculture

History

Music

Art

# AI Progress

AI detected lung cancer 5% more often than experts.

11% more likely to reduce false positives.

Healthcare

Google's lung cancer detection AI outperforms 6 human radiologists | VentureBeat

# AI Progress

write a biblical verse in the style of the king james bible explaining how to remove a peanut butter sandwich from a VCR

And it came to pass that a man was troubled by a peanut butter sandwich, for it had been placed within his VCR, and he knew not how to remove it.

And he cried out to the Lord, saying, "Oh, Lord, how can I remove this sandwich from my VCR, for it is stuck fast and will not budge?"

And the Lord spoke unto him, saying, "Fear not, my child, for I shall guide thy hand and show thee the way. Take thy butter knife, and carefully insert it between the sandwich and the VCR, and gently pry them apart. And with patience and perseverance, the sandwich shall be removed, and thy VCR shall be saved."

And the man did as the Lord commanded, and lo and behold, the sandwich was removed from the VCR, and the man was saved.

And the Lord said, "Verily I say unto thee, seek not to put thy peanut butter sandwiches in thy VCR, for it is not a suitable place for such things. Rather, keep thy sandwiches in thy refrigerator or on thy plate, where they belong."

And the man heeded the Lord's words, and from that day forth, he kept his sandwiches in their proper place, and was saved from trouble and woe. Amen.

Content Generation

# AI Progress

Travel

# AI Progress

Entertainment

# AI Progress

Agriculture

# AI Progress

Film

# AI Progress

Vincent Van Gogh

Alexander The Great

Napoleon Bonaparte

William Shakespeare

Wolfgang Amadeus Mozart

Henry VIII

## History

I Used AI To Find Out What These 21 Historical Figures Would Look Like In Modern Times | Bored Panda

# AI Progress

The stars have genuinely collaborated on previous songs including The Ride and Live For

Music

# AI Progress

Art

**BT Group**

# AI & Security

# Automated Detection and Response



**Anomaly Detection:** Identifying suspicious behaviour using unsupervised machine learning to flag novel attack patterns for further investigation.

**Threat Detection:** Supervised algorithms using labelled data of known threats. For example, classifying DNS requests made to algorithmically generated domains – a strong indicator of malware infection.
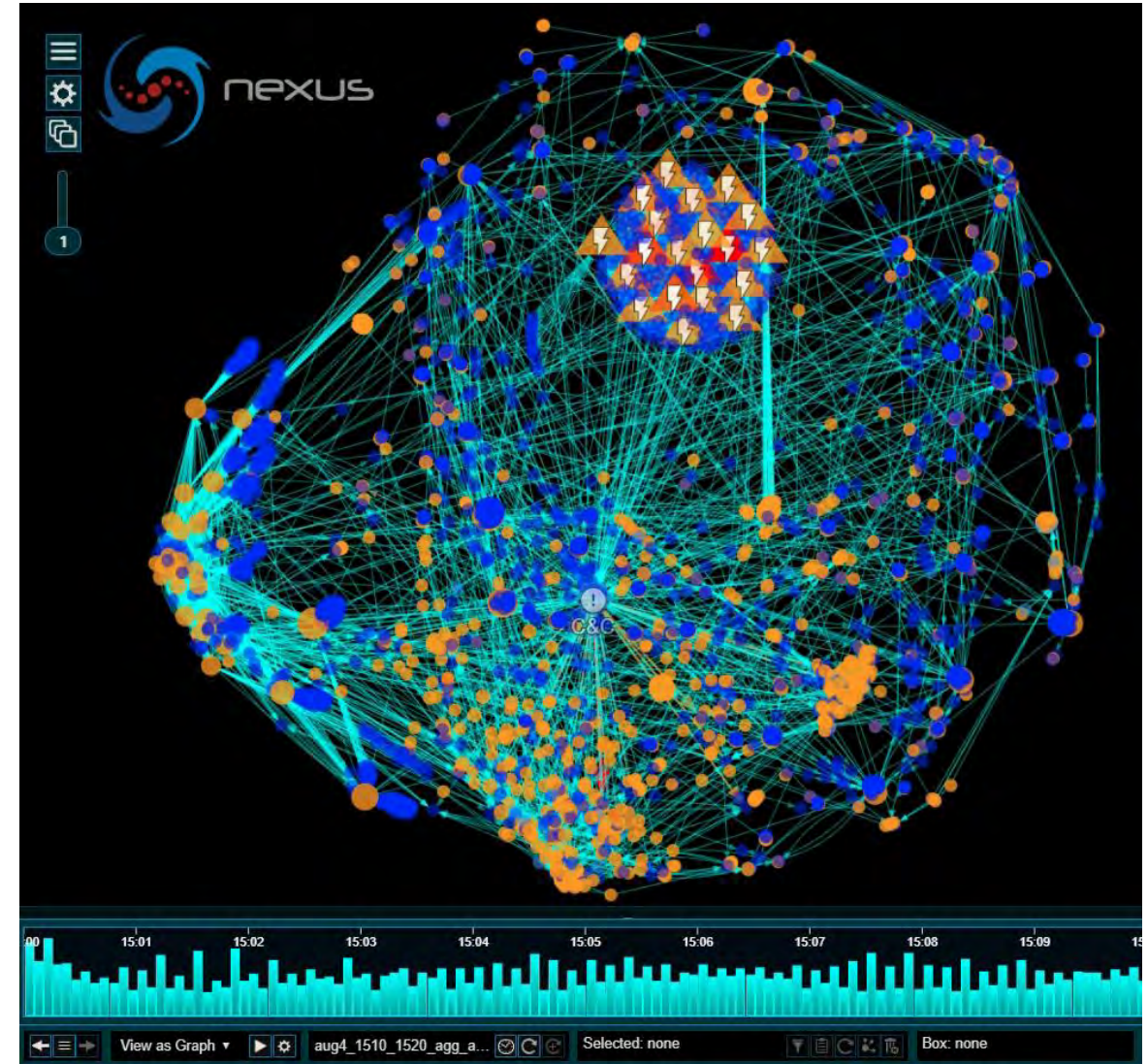
**Automated Security Response:** AI can enable real-time defence to a sensed threat. Where defensive strategies are unclear, simulation and modelling can help derive new techniques to counter attacks.

**Conceptual Knowledge:** Bridging the gap between AI systems and human experts by providing insight into data classification results.

**Visual Analytics:** Graph-based visual analytics provide analysts with a quick and accessible means of identifying heuristic patterns in data.
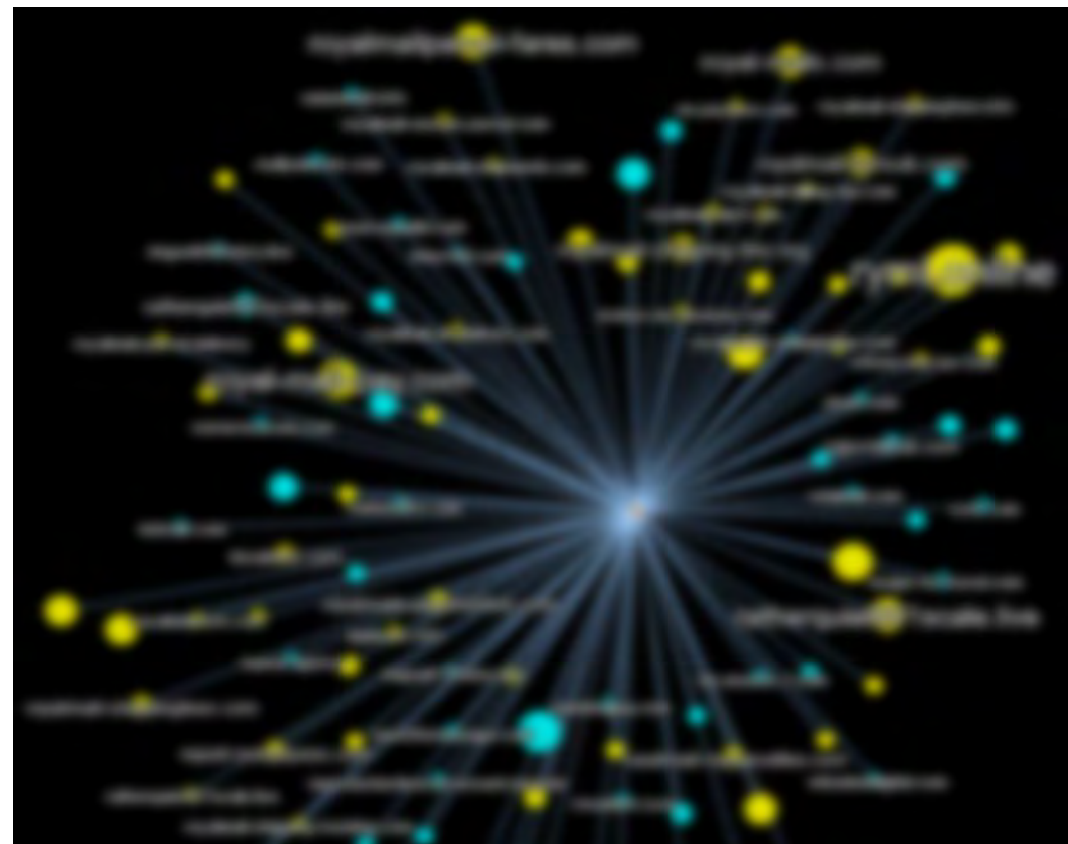
# Nexus Graph Analytics

- Our Nexus software model relationships that exist in data and allow resultant graphs to be visually explored by analysts.

- Underpinned by AI-based big data analytics techniques it preserves the most salient aspects of data before pushing to analysts to triage within context.

- Network events labelled as anomalies using novel methods using neural network embeddings are being analysed for significance by an analyst.

- Security operations teams are able to discern simply suspicious activity from that which is truly malicious.

# Dragnet – Newly Registered Domains

Phishing groups register spoofed domains to mimic the organisation(s) they are targeting.
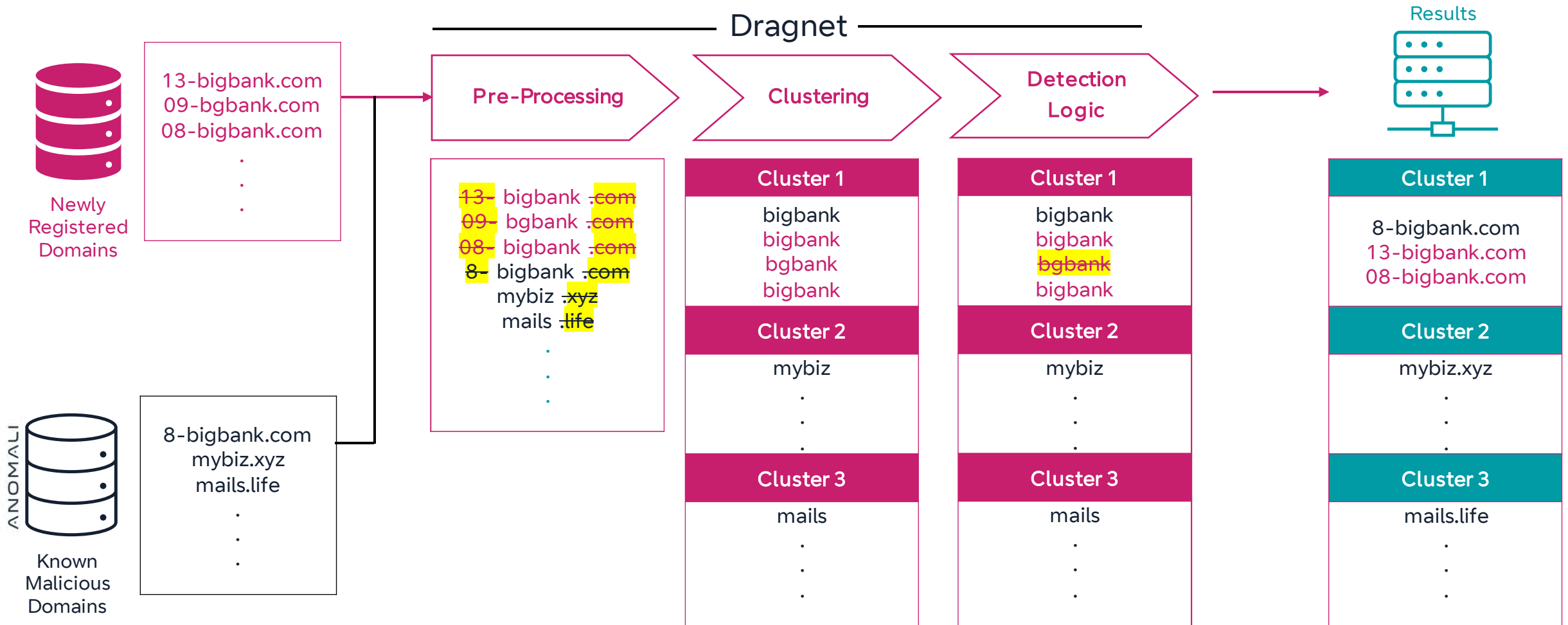
Relying on a passing familiarity to known brands, we can use proprietary feature extraction and clustering to detect likely malicious domains – before they have a chance to be used.



Value demonstrated and in pipeline for operational deployment within BT. We are also exploring revenue opportunities with the ThreatCo customers.
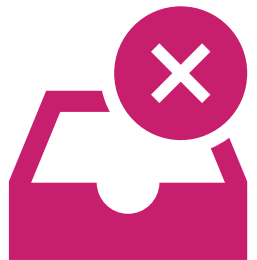
# Dragnet

## How it works



Dragnet

Newly Registered Domains

13-bigbank.com
09-bgbank.com
08-bigbank.com
.
.
.

Known Malicious Domains

8-bigbank.com
mybiz.xyz
mails.life
.
.
.

Pre-Processing

13- bigbank .com
09- bgbank .com
08- bigbank .com
8- bigbank .com
mybiz .xyz
mails .life
.
.
.

Clustering

**Cluster 1**
bigbank
bigbank
bgbank
bigbank

**Cluster 2**
mybiz
.
.
.

**Cluster 3**
mails
.
.
.

Detection Logic

**Cluster 1**
bigbank
bigbank
bgbank
bigbank

**Cluster 2**
mybiz
.
.
.

**Cluster 3**
mails
.
.
.

Results

**Cluster 1**
8-bigbank.com
13-bigbank.com
08-bigbank.com

**Cluster 2**
mybiz.xyz
.
.
.

**Cluster 3**
mails.life
.
.
.

# Dragnet Evaluation

## Malicious Domains Feed

We've been using Virus Total to validate domains we believe to be malicious.

On average, it takes 4.13 days for malicious domains to be identified on Virus Total but can be as many as 39 days (so far)!

20.3% of malicious domains were incorrectly marked as harmless on Virus Total.

On average it took 15 days for this to be corrected.

Virus Total combines intel from:
Comodo, Crowdstrike, Fireye, Fortinet, Google, Kaspersky, Malwarebytes, McAffee, Microsoft, Palo Alto, Panda, Trend Micro and hundreds of others.



NRD Time Analysis

Days since 24th Nov 2022

# Automated Response – The Vision



Traditional Security
Response (Manual)



SOAR Tools



Inflame

Control

Threat

Graph

Topology

Data

Inflame

Threat Simulation          Threat Response          Threat Awareness

# Decision Making – Novel Solutions



OpenAI Hide and Seek (2019)
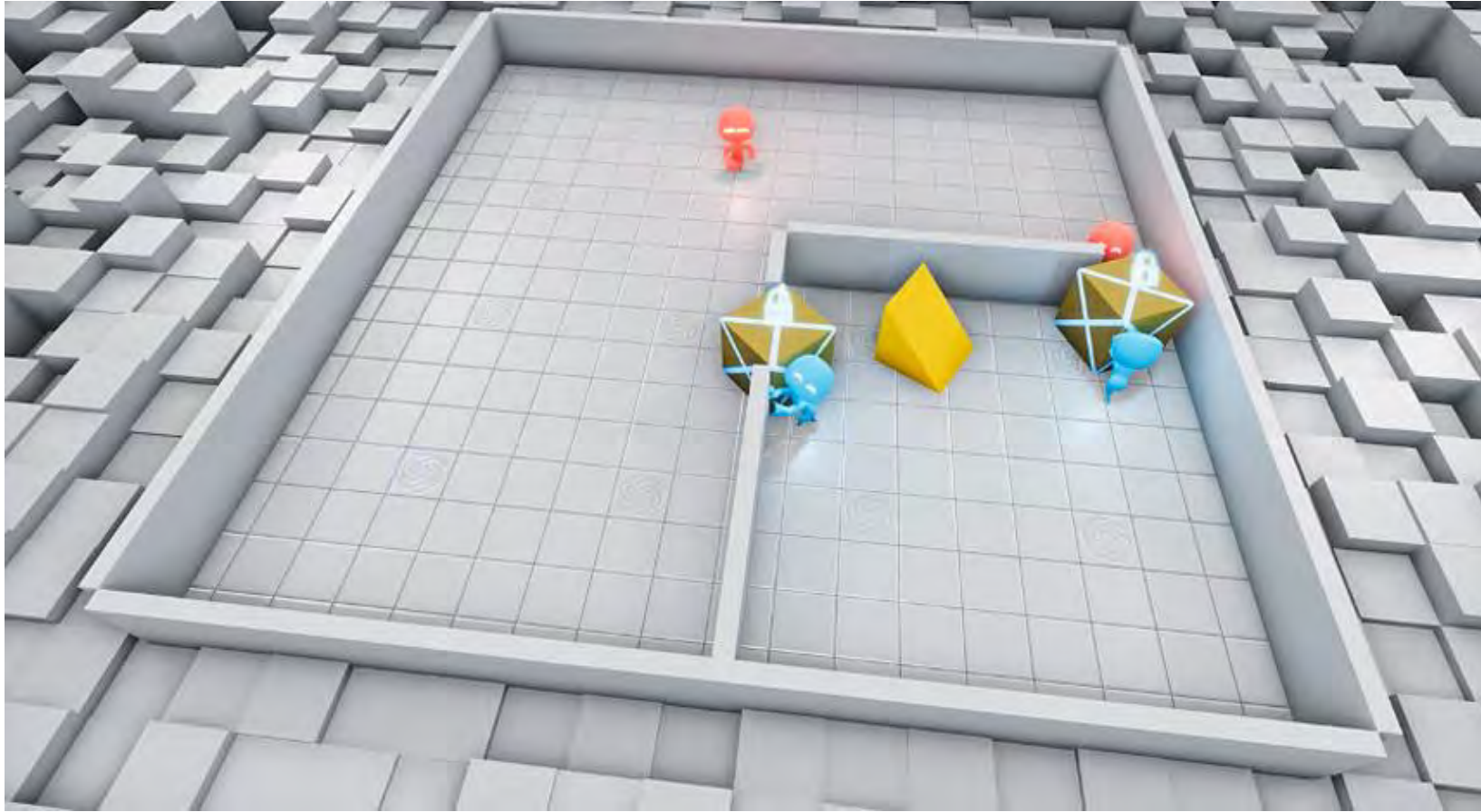
# Decision Making – Novel Solutions

# Decision Making – Novel Solutions

# Decision Making – Novel Solutions

**BT Inflame**

Clear comparisons

REMOVE - firewall_1
Mean health after 3 steps: 83.72%
Max health after 3 steps: 83.72%
Min health after 3 steps: 83.72%

REMOVE - firewall_2
Mean health after 3 steps: 44.20%
Max health after 3 steps: 57.13%
Min health after 3 steps: 35.50%

REMOVE - firewall_3
Mean health after 3 steps: 32.90%
Max health after 3 steps: 34.02%
Min health after 3 steps: 31.80%

Do nothing
Mean health after 3 steps: 11.16%
Max health after 3 steps: 14.76%
Min health after 3 steps: 6.72%

Graph | Selected: none | Box/Lasso: none

Graph | Selected: none | Box/Lasso: none
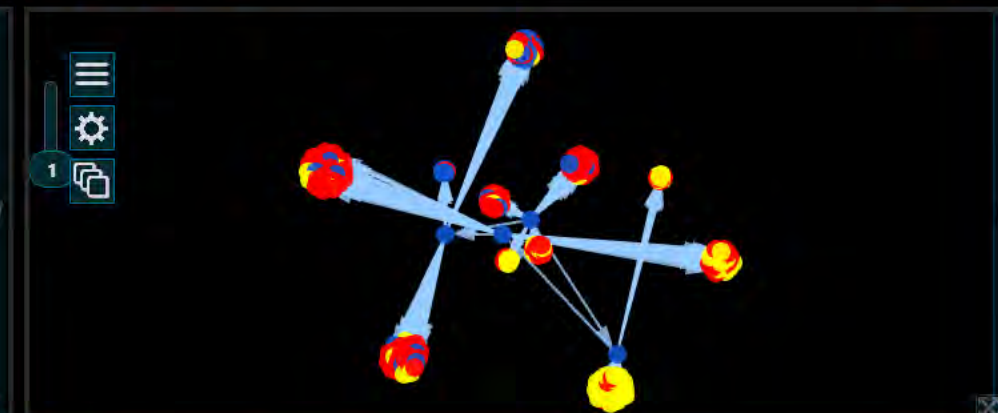
Name: REMOVE - firewall_1 | Current health score: 83.72%
Step: 3 / 3

Name: REMOVE - firewall_2 | Current health score: 57.13%
Step: 3 / 3

Graph | Selected: none | Box/Lasso: none

Graph | Selected: none | Box/Lasso: none

Name: REMOVE - firewall_3 | Current health score: 34.02%
Step: 3 / 3

Name: Do nothing | Current health score: 14.76%
Step: 3 / 3

# AI is being used to defend against new threats...

# ...but it's also being used to create new threats we need to defend against.

# AI-driven Cyber Attacks

**Social Engineering**

Automated spear-phishing

**Reconnaissance**

Discovery of high-value assets

**Strategic Planning**

Efficient AI-planned attacks

**AI-Generated Malware**

Enhanced propagation

**Assistance**

Generation of commands

**Stealth**

Deception and decoy attacks

# Generative AI

Recent progress in generating synthetic data



Increasingly realistic synthetic faces generated by variations on Generative Adversarial Networks (GANs). In order, the images are from papers by Goodfellow et al. (2014), Radford et al. (2015), Liu and Tuzel (2016), and Karras et al. (2017).

# Deepfakes

**Next-gen upgrade for disinformation**

*"... deepfakes and real images will become indistinguishable, even to high-quality detectors."*

"Artificial Intelligence, Deepfakes and Disinformation", RAND, July 2022



**OpenAI Sora.** Prompt: The camera follows behind a white vintage SUV with a black roof rack as it speeds up a steep dirt road surrounded by pine trees on a steep mountain slope...

# AI is a key factor in empowering nation states

**"China's government sees AI as a promising military 'leapfrog development' opportunity"**

"Intelligence supremacy will be the core of future warfare" and that "AI may completely change the current command structure, which is dominated by humans" to one that is dominated by an "AI cluster."



Center for a New American Security

Understanding China's AI Strategy:
Clues to Chinese Strategic Thinking on
Artificial Intelligence and National Security

February 2019

By Gregory C. Allen

# AI & Deterrence

## AI increases the risk of escalation

Manned systems better for deterrence than unmanned ones.

Widespread AI and autonomous systems could lead to inadvertent escalation and crisis instability.

Machines will likely be worse at understanding the human signalling involved in deterrence, especially de-escalation.



"We expect artificial intelligence (AI) and autonomous systems to significantly change the future battlefield."

2020 RAND Report: Deterrence in the Age of Thinking Machines
https://www.rand.org/pubs/research_reports/RR2797.html

# What is AI thinking?

**Human and AI collaboration is key**

Generalised AI is still some way off.

AI makes mistakes and will continue to do so.

Humans are needed to validate these.

AI is necessary, already omnipresent, but should it be fully automatic?

# Intelligence Augmentation



Visual Analytics

Machine-led ⟶ Human-led

| Automated Processing | Interactive Visualisation | Validation and Triage |

# Summary

AI is a tool; we need to steer its design and application.

AI is helping us defend against threats, but we need to be aware of the risks associated with AI.

The future of AI is yet to be determined, but it's accelerating rapidly.

*"The true power of AI lies not in replacing human potential, but in amplifying it to unprecedented heights, inspiring us to reach new frontiers of innovation, understanding, and compassion."*

ChatGPT, AI

**BT Group**

# Any Questions?

# BT Group