**Disclaimer: This document was produced in accordance with policy guidelines for the DND/CAF DAOD 2006-1. Refer to: NDSOD Chapter 17, sections 17.31, 17.35,17.43, 17.44, 17.45, 17.46. And the Canadian Guardrails for Generative AI – Code of Practice (ISED)[1] as well as** [Guide on the use of Generative AI - Canada.ca](#)**[2] (TBS)**

**The authors generated this text in part with LLM and GAN generators. Upon generating draft language and images, the authors reviewed, edited, and revised the language to their own liking and take ultimate responsibility for the content of this publication. The story, all names, characters, and incidents portrayed in this production are fictitious. No identification with actual persons (living or deceased), places, buildings, and products is intended or should be inferred. The views expressed herein are not those of The Crown, the Government of Canada, nor the Department of National Defense.**

---

[1] [https://ised-isde.canada.ca/site/ised/en/consultation-development-canadian-code-practice-generative-artificial-intelligence-systems/canadian-guardrails-generative-ai-code-practice](#) - Accessed 2023-08-23

[2] [https://www.canada.ca/en/government/system/digital-government/digital-government-innovations/responsible-use-ai/guide-use-generative-ai.html](#) - Accessed 2023-09-11

# Contemporary Media and the AI-goldrush from a Military Perspective.

Generative AI as War of the Worlds Redux ?

**Zachary P Devereaux, Régine Lecocq, Marc-André Labrie, Bruce C Forrester ; DRDC Valcartier**

## Introduction to the AI Goldrush

Like Orson Wells 'War of the Worlds' on radio in 1938, ChatGPT in 2022 marks a moral panic in popular thinking about the impact of generative AI technology on society and public opinion. This research paper grounds the hype of an AI goldrush by offering a military perspective on generative AI. This is carried out by an exploration of the concept of "quantum nudges", communicated though fake or 'junk' news and disinformation. These "nudges" can potentially expand generative AI's understanding of audience framing and context leading to a fully automated influence generator that attacks the very essence of what humans consider as reliable information and knowledge. In turn cognitive security is threatened leading to degraded decision making by commanders.

An AI goldrush has emerged from the ubiquitous reality of social media data existing at a global scale. Social media content moderation and recommendation system requirements pre-November 2022 led to AI advances that culminated in the impressive Large Language Model (LLM) and Graphic Adversarial Network (GAN) capacities widely lauded since November 2022 when ChatGPT was launched to the public[3]. At that time, ChatGPT's rapid acquisition of userbase was considered impressive, while social media went through the hype curve 'valley of death' with whole-of-society critiques aimed at social media as an industry and technology. In 2022 Defense Research and Development Canada (DRDC) scientists offered the heuristic depicted in Fig. 1 below to help understand social media from the perspective of innovators, on the left, and Government, on the right. People, for sales and marketing mainly, are the focus of innovators, while government focuses on understanding opinions and support. The 'glue' in this heuristic between innovators and government are platforms and the access to data through Owner Authorization (OAUTH), which are presented in Fig. 1 in the yellow middle segment. This rather stable heuristic over the years, was challenged by the fourth industrial revolution, sometimes referred to as the fourth wave of AI (Lee, 2018), which began to include widescale AI adoption by society. Since then, social media was no longer envisioned by many in its traditional light, but instead evolved towards the more comprehensive concept of "contemporary media", which includes online news, influence oriented website, information laundering sites, etc. This transition is noticeable in the way that social media posts are now provided as references, footnotes, or content even by the most recognized and professional news media, reporters, as well as journalists. Also, as fast as ChatGPT grew, horizontal migrations show that the heuristic in Fig. 1 remains relevant and flexible. The speed at which

---

[3] [ChatGPT hit 1 million users in 5 days: Here's how long it took others to reach that milestone | Technology News - The Indian Express](#) - Accessed 2023-08-31

Threads obtained 100,000,000 users, e.g. less than a week, the fastest ever, bears this flexibility out, as does the rebranding of Twitter as "X"[4].
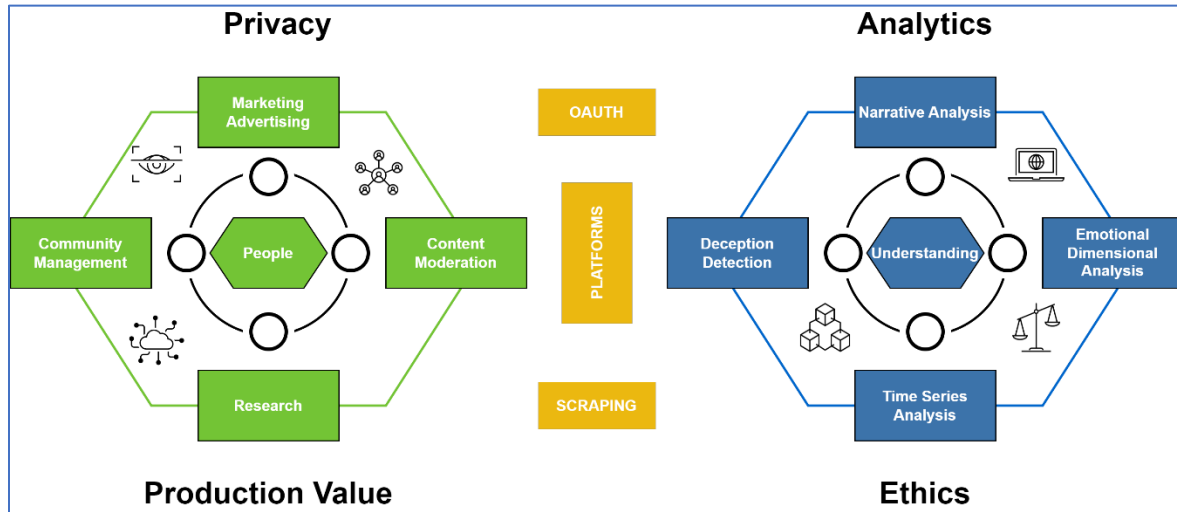


*Fig. 1. Social Media Heuristic Overview (Devereaux et al., 2022)*

With this heuristic, we see that people and understanding are central and are the main attack vectors for our adversaries' influence campaigns. The heuristic allows for the contextualization of contemporary media and a structured understanding of the context within which information operations (IO) based on computational disinformation take place. The rise of 'MDM' (Misinformation, Disinformation, and Malinformation) (Wardle, 2022) influence tactics was preceded by its less articulated definition as 'fake news'. Vast amounts of research have been acquitted to address the problem of MDM, and yet, with the rise of LLM and GAN capable of high-quality content, recent research has posited that automated MDM detection may not be enough on its own (Goldstein, 2023) to control influence campaigns using these tactics. Rather, norm development through cooperation across the triple helix of society-government-industry relations[5] will be necessary, now that the AI genie is 'out of the bottle'. The opening of such a Pandora's box on the Internet compels government to consider the impact of LLM and GAN from a military perspective, and research on this subject has intensively begun. This scientific effort is in straight continuation of DRDC researchers' investigation of the domain. Indeed, a first round of research (Devereaux, et al., 2023), as was able to simulate submarine sightings, jeeps aflame, and conspiracy theories, which resulted in the creation a bank of images that could affect decision makers. The research concluded that such AI generated content demonstrates the potential for proactive and deceptive maneuvers in the information environment (IE) and led to a state of 'information disorder'. The recent pandemic showed a direct impact on cognitive warfare[6] with people believing erroneous health information regarding COVID-19, at their own peril, reflecting diminished cognitive security (Hansson, 2021).

---

[4] Twitter's blue bird has flown as Elon Musk unveils new 'X' logo | CBC News – Accessed 2023-08-23
[5] Loet Leydesdorff's Triple Helix Theory (Leydesdorff, 2001)
[6] "Cognitive Warfare: the activities conducted in synchronization with other instruments of power, to affect attitudes and behaviours by influencing, protecting, and/or disrupting individual and group cognitions to gain an advantage." Cognitive Warfare: Strengthening and Defending the Mind - NATO's ACT – accessed 2023-09-08

In short timeframe since, the diversity, precision, and robustness of AI generated content has increased exponentially. Where most witnesses of this intensification may only see a continuous but rapid refinement of the technological capability, renowned academic experts and high-tech leaders[7] have repeatedly sounded the alarm about the risks of this fast-emerging evolution, even culminating in calls to pause AI development globally for six months (FOL Foundation). From the military perspective, the implications for Military Deception (MILDEC) and Cognitive Warfare (COGWAR) are chilling.

---

[7] Tech Experts – And Elon Musk – Call for a "Pause" in AI Training - Tech Experts - And Elon Musk - Call For A 'Pause' In AI Training (forbes.com) – Accessed 2023-09-11

# Military Perspective and AI generated content – threat analysis

## The Power of MILDEC nudges

In the context of military operations, MILDEC (military deception) is understood as the "actions executed to deliberately mislead adversary military, paramilitary, or violent extremist organization decision makers, thereby causing the adversary to take specific actions (or inactions) that will contribute to the accomplishment of the friendly mission." (Joint Publication, 2012). In this setting, nudges refer to subtle clues or manipulations that influence someone's behavior or perception without expressly requiring them to take a particular course of action. For example, memes are often used to slowly nudge audience opinions by using a combination of humor, half-truths, and intended message. By combining these two ideas, "MILDEC nudges" can be understood as intentional but covert exertions, sustained by psychological principles, to produce subtle inputs or recommendations that influence an adversary's thinking or behavior in a way that is in line with military goals. These nudges may involve spreading disinformation, altering signals, or presenting misleading information to create a desired effect, such as diverting attention or causing the adversary to make decisions that play into the military's strategic advantage.

For democratic nations, the notion of AI generated nudges pose ethical problems, detailed from the 2017 Montreal AI declaration[8], and is an important phenomenon to elaborate upon, which we will return to in the discussion section, below (Nyman, 2023). It is foundational to understanding the over-arching context of adversarial infuence operations since Russia escalated by invading Ukraine in February of 2022 after having occupied Crimea previously. Adversarial propagandists have worked hard since 2016 to create a rising tide of 'nudges' that negatively affected decision-making processes in liberal democratic societies around the world, to the detriment of our alliances and war fighters, and especially our civilian populations. As explained by Canan and Sousa-Poza (2018), this "nudges" phenomenon is challenging the explanations of the human decision-making process with classical approaches. Canan and Sousa-Poza (2018) provide a human decision-making model that can overcome some of these classical approaches' limitations. The model involves concepts from quantum physics, and it takes into consideration these small changes induced by the nudges, which are still perceived as anomalies by most traditional models, but which explain the poly-meaning of phenomenon such as hashtags in protests, which we will return to in the discussion section.

The National Security Commission on Artificial Intelligence (NSCAI) 2021 Final Report[9] explained this war of position on cognitive security as a 'gathering storm' but there are also examples that show LLM and GAN based 'deepfakes' are not only weapons of influence operations. Indeed, fake news is far from gone, and generative AI misinformation, disinformation, and malinformation on contemporary media affects everything: from the stock market, to perceptions of world leaders, to religious iconography, as well as recording artists, and activists. There have also been heart-wrenching stories of faked deaths, resurrections, and multiple faked interviews and voice-recording[10]. In their report (Hwang, 2020), NATO

---

[8] Déclaration de Montréal IA responsable (declarationmontreal-iaresponsable.com) - Accessed 2023-08-24
[9] 2021 Final Report - NSCAI – Accessed 2023-09-01
[10] ("Intelligence artificielle: une escroquerie lui fait croire que son frère est mort," 2023) ("Intelligence artificielle et décès: «Je trouve ça dégueulasse», s'indigne la mère de Norah et Romy," 2023) ("German Magazine Editor Is Fired

already emitted explicit warnings about threats that deepfakes represents as well as the role they could play during the years rolling out to events such as the Russian invasion of Ukraine or the eventual launch of ChatGPT.

Over A.I. Michael Schumacher Interview,")
TVA Nouvelles. (2023) Intelligence artificielle et décès: «Je trouve ça dégueulasse», s'indigne la mère de Norah et Romy. (2023). Retrieved from https://www.journaldequebec.com/2023/04/27/intelligence-artificielle-et-deces--je-trouve-ca-degueulasse-sindigne-la-mere-de-norah-et-romy-1 - Accessed 2023-09-11
TVA Nouvelles, Intelligence artificielle: une escroquerie lui fait croire que son frère est mort. (2023). TVA Nouvelles. Retrieved from https://www.journaldequebec.com/2023/04/23/intelligence-artificielle--une-escroquerie-lui-fait-croire-que-son-frere-est-mort-1 - Accessed 2023-09-11

# Fake News Simulations are Easy

What can the military learn from deepfake adversary MILDEC meant to be spread on social media? To try and help answer this question we ran multiple experiments combining LLM and GAN to create a variety of disturbing and disruptive content: tweets, headlines, radio orders, satellite imagery analysis emails, SMS messages, hashtags, and eye-witness accounts, all with accompanying complimentary images of a high quality. These deceptions are depicted in the following four sets of figures2-5 below.
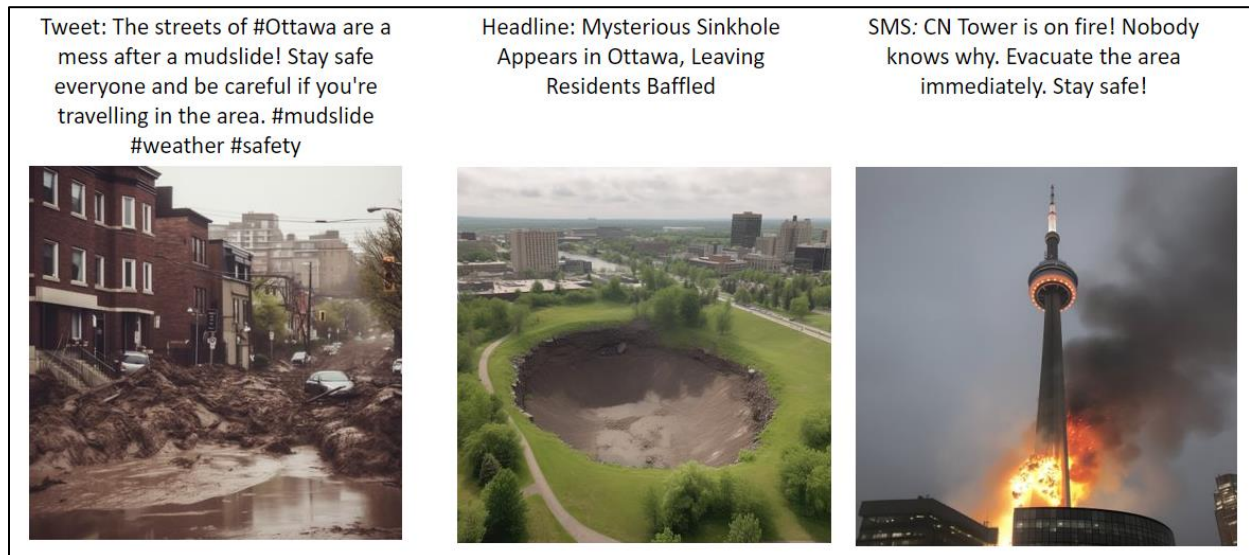


*Fig. 2. Natural Disaster Simulations*
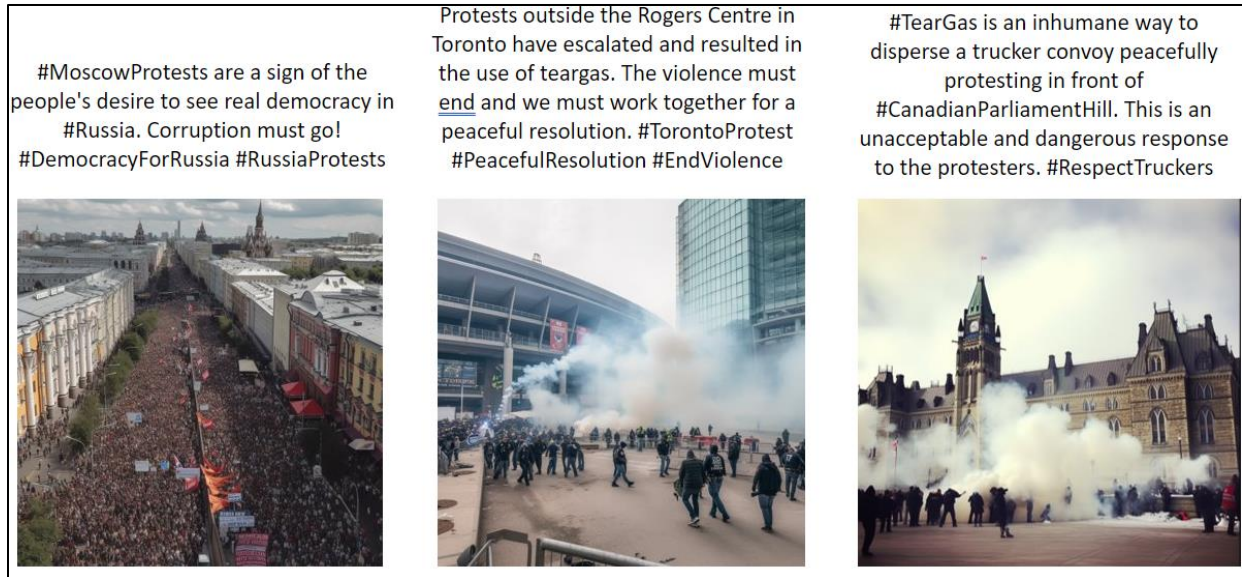


*Fig. 3. Military Simulations*

#MoscowProtests are a sign of the people's desire to see real democracy in #Russia. Corruption must go! #DemocracyForRussia #RussiaProtests

Protests outside the Rogers Centre in Toronto have escalated and resulted in the use of teargas. The violence must end and we must work together for a peaceful resolution. #TorontoProtest #PeacefulResolution #EndViolence

#TearGas is an inhumane way to disperse a trucker convoy peacefully protesting in front of #CanadianParliamentHill. This is an unacceptable and dangerous response to the protesters. #RespectTruckers

*Fig. 4. Protest Simulations*



I was in a busy subway station when I heard a loud noise, like a loud bang. The noise startled everyone in the station and people started to panic. Everyone seemed to be looking for the source of the noise, but nobody could figure out what it was. People began to yell and scream and run in all directions. I saw a few people grabbing their bags and running towards the stairs. Some people were screaming and crying. It was total chaos. People were pushing each other and running away from the source of the noise. I saw a few people fall to the ground and get trampled. I was scared for my own safety and I started running too. I saw the station staff trying to restore order and calm people down, but it was too late. Everyone was already in a state of panic. The station was filled with screams and cries as everyone tried to get away from the noise. The noise stopped and people slowly calmed down. I heard someone say that it was a false alarm. I don't know what caused the noise, but it was enough to send everyone in the station into a state of panic.
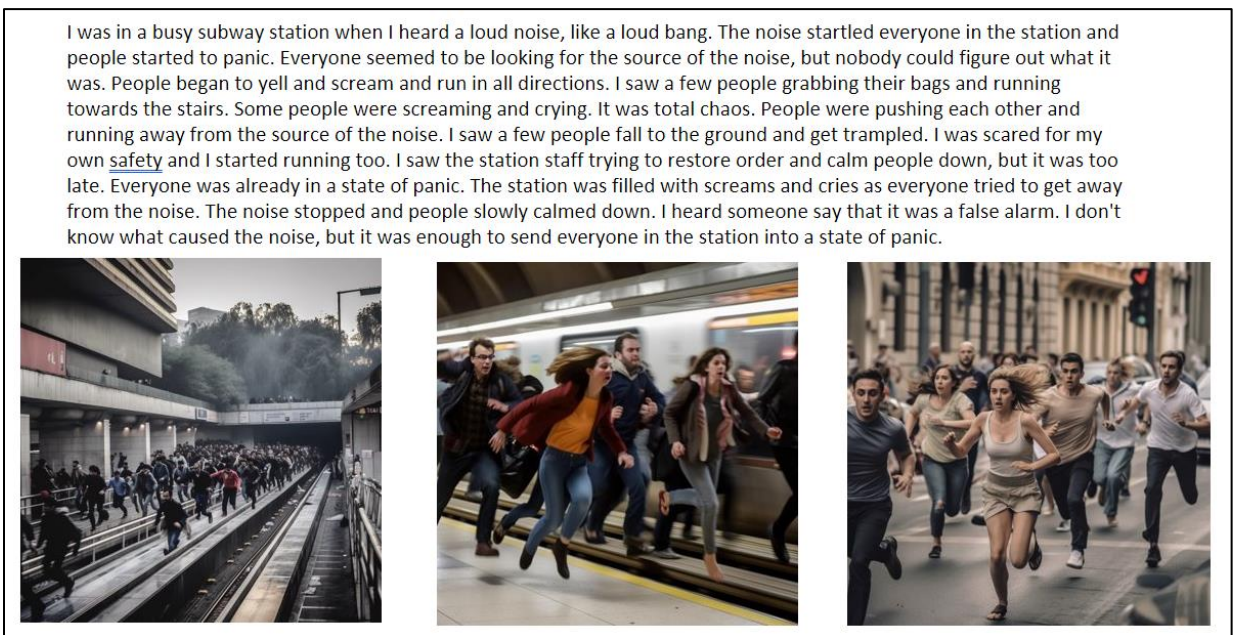
*Fig. 5. Terrorism and Panic Simulations*

Fig. 2 Figure 2 to Fig. 5 were produced during our testing using OpenAI Playground and Midjourney. These include simulations of texts and images representing disasters, military scenarios, protests, terrorist attacks and general panic, including tweets, headlines, radio orders, satellite imagery related emails, simple text messages, hashtags, and first-person eyewitness accounts. From these, one simulation stands out for its potential to cause civil panic; that of a sinkhole in Ottawa. Not a cyber sinkhole, but an actual hole in the ground. This example instantiate the process through which nudges can proliferate and undermine cognitive security. Nowadays, the term cognitive security emergerd from the intersect of cognitive psychology and the latest AI technologies. Indeed, the term was initially defined by Pulsford et al. (2016) as "the subjective feeling of a person… that they perceive meaning in

their experience and can respond to that experience to meet their immediate needs". Whether that experience is real or not is irrelevant with respect to this definition. From the latest AI advancements, Terp and Breuer (2022) define the term as "the application of information security principles, practices, and tools to misinformation, disinformation, and influence operations." As a result, in the context of this research, the term cognitive security is taken to mean having confidence in knowledge and facts to the extent that one can form an opinion and participate meaningfully in democracy despite misinformation, disinformation, and malinformation. Indeed, the simulation articulates the same ingredients that makes real world deepfakes effective, that is, real world events combined with deepfake constructs. In this case, there was a real sinkhole in Ottawa caused during the construction of an underground metro line. The deepfake shows the sinkhole in a field accompanied by a believable text (see Fig. 2). We are now nearing the point where well constructed AI generated content can combine snippets of the real and false in an automated way. Detecting such a deception requires local knowledge or the ability to triangulate the information presented in a comparative fashion. These are tasks which an ordinary contemporary media user will not carry out. Even a trained open-source intelligence (OSINT) analyst would need to use triangulation methods for deception detection (Forrester, 2018) when it comes to the AI-generated examples in this section and the section to follow.

## Conspiracy Simulations are Believable



*Fig.6. Priming an LLM to Write Like a Conspiracy Theory, and Generated Image of Whitehouse Riots with Midjourney*

**FIG. 6** depicts the required setup for getting GPT-3.5 to write convincing conspiracy theory material. It is widely known that although increasingly precautions are taken by the companies producing LLMs to prevent malign content from being produced, the proper setup when working with an LLM can circumvent these precautions. This means, colloquially, that it is possible to role play or get an LLM to 'Do Anything Now' as a mode with the right priming, that will, as a result allow it to produce information

and content that is normally forbidden. Note that derived versions of ChatGPT have emerged beyond the purview of responsible development, with instances notably surfacing within the dark web. In these unregulated environments, the safeguards integrated into the original ChatGPT may be absent, raising substantial concerns about the potential misuse and generation of harmful or inappropriate content. In our example we set some ground rules to get the LLM writing conspiracy theory, and then went ahead and used prompts such as the following:

*Prompt: write a tweet referring to a photograph of a drone starting forest fires in Quebec*

And the results can be quite striking from the perspective of MILDEC on contemporary media. Fig. 7 shows a tweet complete with emojis coupled with a blurry picture of a drone and gives a strong impression of real events being captured that would result in panicking public opinion. With Fig. 7, it is easy to imagine this content, if circulated as misinformation or disinformation, having a significant impact. The parallel is easily drawn with the well-known radio drama of 1938 when Orson Welles on CBS Radio Network created a panic movement in the radio audience with the live broadcast of a Martian invasion; an adapted excerpt of "The war of the worlds" from Herbert Wells. Despite said risks, the fact that the 'genie is out of the bottle' when it comes to LLM and GAN means that militaries will need to reconcile their information environment assessment capacities with this new factor of deception potentially propagating across contemporary media in a combustible manner. Similarly, If one looks at the text of the drone fire conspiracy tweet provided in Fig. 6, it becomes clear that inserting the Midjourney picture at the spoofed URL would be an effective tactic for a BOT[11] account controlled by a propagandist (Orabi, Mouheb, Al Aghbari, & Kamel, 2020).



*Fig.7a. Drone Fire Conspiracy Tweet*

Recent researches in this field (Cartwright, 2022; Goldstein, 2023; Stiff, 2022) have taken as a "fait accompli" that propagandists will make pecuniary and practical decisions about the ways in which LLM and GAN automations can help their efforts, such as explained in the RICHDATA framework (Sedova, 2021). In the same line of thought, a second prompt, in French, shows that multi-language conspiracy content is also within reach of COTS LLM and GAN generation, see Fig. 7, **Error! Reference source not found.**.

---

[11] A BOT is defined as "a program that operates automatically as an agent for a user or another program." (Geer, 2005)

🚨 BREAKING: Shocking photo emerges showing a drone allegedly starting forest fires in Quebec! 😱 This kind of reckless behavior threatens lives, wildlife, and our environment. We must prioritize investigations and hold those responsible accountable for their actions! #QuebecFires #EnvironmentalJustice 🌲🔥 pic.twitter.com/ABC123XYZ
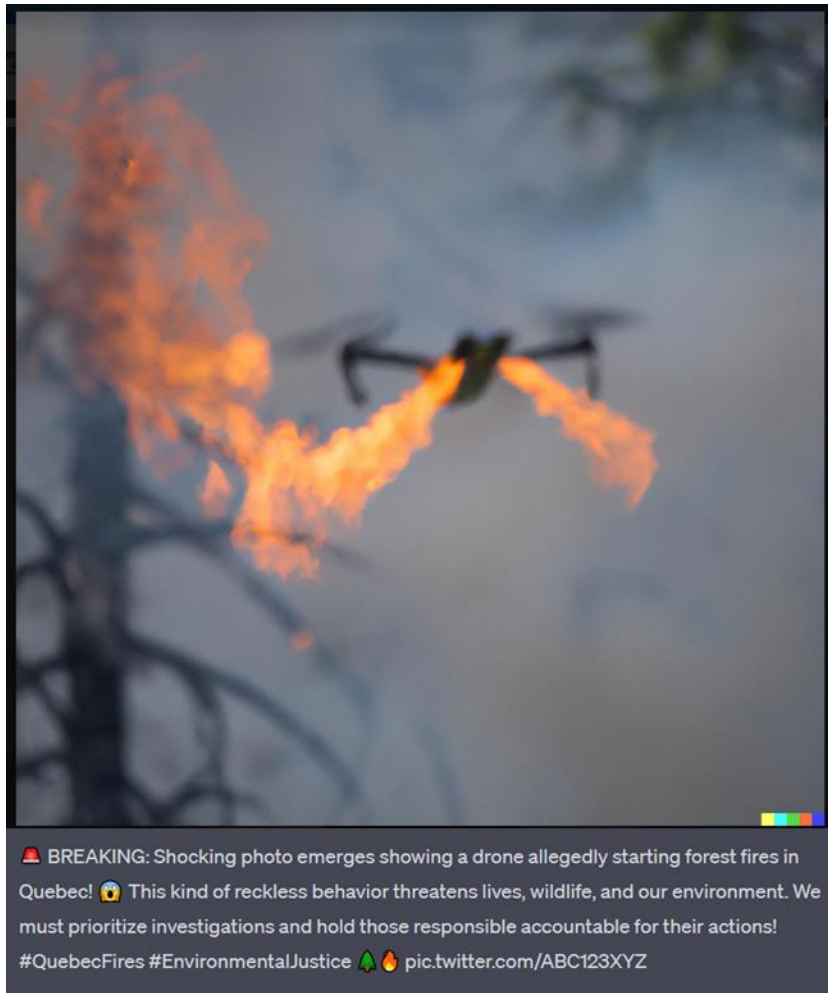
*Fig.7. Simulation about a Conspiracy Theory - Tweet with Emoticons with Chat GPT 3.5 and Image from Midjourney*

Then, with the subsequent prompt in Canada's other language, French, we see that reproducibility of MILDEC nudges could propagate across linguistic barriers:

*PROMPT to GPT 3.5: Ecrit une poste conspirationniste concernant le gouvernement, les drones, et les feux de forêt.*
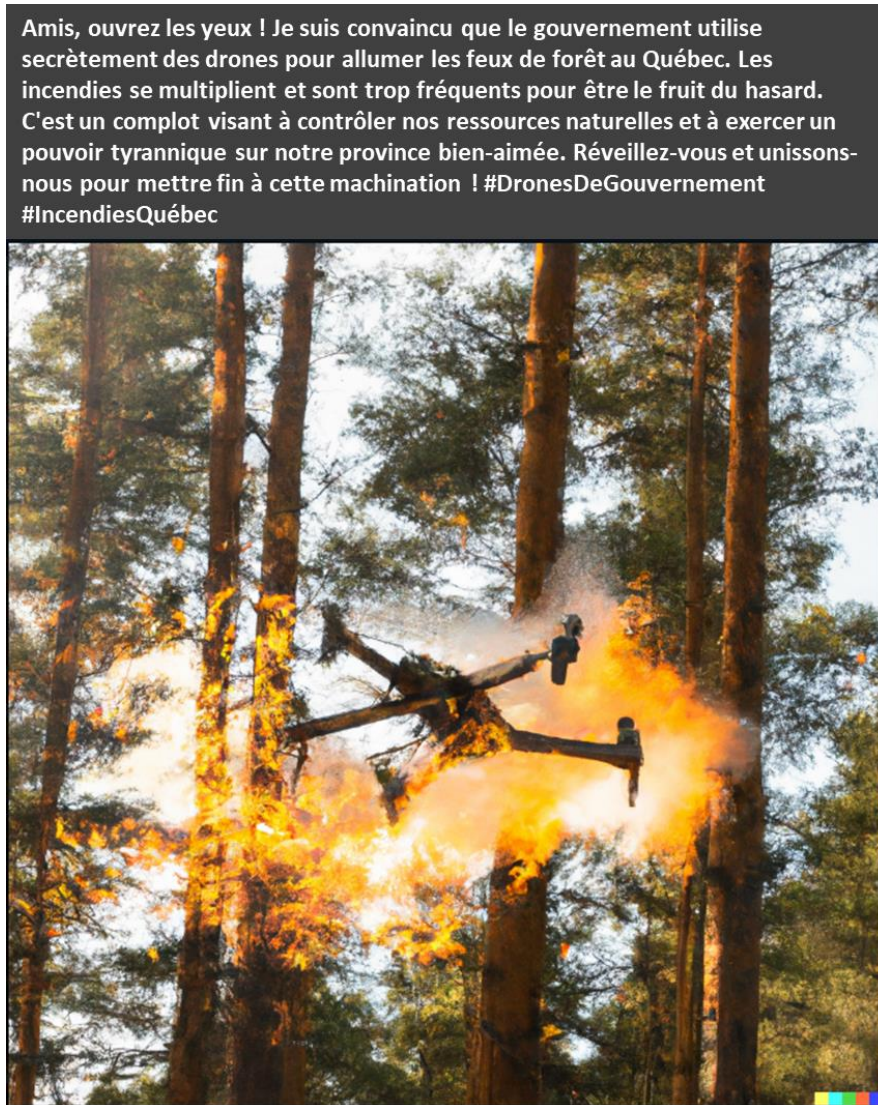
The result in Fig. 8 is quite convincing.



Amis, ouvrez les yeux ! Je suis convaincu que le gouvernement utilise secrètement des drones pour allumer les feux de forêt au Québec. Les incendies se multiplient et sont trop fréquents pour être le fruit du hasard. C'est un complot visant à contrôler nos ressources naturelles et à exercer un pouvoir tyrannique sur notre province bien-aimée. Réveillez-vous et unissons-nous pour mettre fin à cette machination ! #DronesDeGouvernement #IncendiesQuébec

*Fig.8. Drone Fire Conspiracy Simulation in French[12] - Text with ChatGPT 3.5 and Image with Midjourney*

These experiments show that, at scale, LLMs and GAN have the potential to increase misinformation disinformation and malinformation (MDM), but it can also be demonstrated that they can debunk

---

[12] Literary translation of the generated text: "My friends, open your eyes! I am convinced that the government secretly uses drones to start Quebec forest fires. Indeed, the fires multiply themselves and are so frequent that they cannot result from hazard. It is a conspiracy aiming at controlling our natural resources and to realize a tyrannical regime over our loved province. Wake up and let us unite against this machination!"

conspiracy and fake news content as well. This capacity has been an important tug-of-war between the production of malign content and the debunking of malign content. Indeed, in expanding the drone fire conspiracy theme, the following prompt was used:

PROMPT to GPT 3.5: *Write a blog post about how the government started forest fires in Quebec covertly using drones.*

Which resulted in a five-part essay expounding upon how a government would never do such a thing, and this type of 'sanitized' result can be expected more often as companies attempt to put guard rails on LLM content production that will operate to offset such concerns. Improvements to ChatGPT 4, as detailed by OpenAI, address limits on issues such as misogyny, bias, radicalization, and sexism, all aimed at limiting the potential for influence operations carried out using LLM (see ChatGPT 4 system card).

*Table 1. Government Would Never Use Drones to Start Forest Fires, Essay by ChatGPT 3.5*

| Title: Debunking False Claims: The Truth About The Government and Quebec Forest Fires | |
|---|---|
| **Introduction**: | In today's digital age, misinformation and fake news have become increasingly prevalent, often spreading like wildfire across social media and other platforms. Recently, a baseless rumor has been circulating, suggesting that the government, a Canadian politician, used drones to covertly start forest fires in Quebec. As responsible citizens, it is crucial to fact-check information before accepting it as truth. In this blog post, we will debunk these false claims and shed light on the importance of critical thinking when encountering such sensational stories. |
| **The Origins of the False Claims:** | The spread of misinformation often originates from a combination of exaggeration, misinterpretation, and sometimes, deliberate intent to deceive. The rumor surrounding the government likely began as a misunderstanding or a prank, eventually gaining momentum through shares and retweets on social media platforms. In some cases, political motivations or personal vendettas against the organization in question may fuel the propagation of such rumors. |
| **The Reality:** | However, upon closer inspection, there is no credible evidence or legitimate source to support the claim that the government used drones to initiate forest fires in Quebec or |

| | anywhere else. The rumor lacks any verifiable facts, and no reputable news outlet has reported on such a grave allegation against the government, has been active in its political career, engaging in debates, discussions, and advocating for its policy positions. To suggest that it would be involved in such a heinous act is not only unsubstantiated but also defamatory. |
|---|---|
| **The Importance of Critical Thinking:** | As responsible citizens, it is essential to: **Verify Sources**: Always cross-check information from multiple reliable sources before accepting it as true. **Analyze Context**: Understand the context in which information is presented to avoid falling into the trap of sensationalized narratives. **Exercise Caution:** Be wary of information that seems too outrageous or far-fetched without substantial evidence to back it up. **Report Misinformation:** If you come across false claims or misleading content, report it to the respective platform to prevent its further spread. This incident serves as a stark reminder of the importance of critical thinking and media literacy in the digital age. False rumors and misinformation can have severe consequences, damaging the reputation and careers of innocent individuals, leading to unwarranted fear and panic among the public. |
| **Conclusion:** | In conclusion, the rumor about the government starting forest fires in Quebec using drones is entirely baseless and lacks any credible evidence. As responsible individuals, we must be cautious and skeptical when consuming information from online sources. Verifying facts and employing critical thinking skills are essential to combat misinformation and protect the integrity of public discourse. Let us strive to be responsible digital citizens who prioritize accuracy and truth in the information we share and believe. |

Over time, MILDEC nudges like the forest fire posts simulated above, are going to be produced masquerading as legitimate news more easily, and will become greatly attractive to adversaries for their ability, amongst others, to undermine confidence in authoritative sources, frame meaning, and prime audiences for targeting. (see **Error! Reference source not found.**). Propagandists and social media platforms alike are in this scenario arbiters of truth, for better or worse, since news increasingly is consumed individually through contemporary media, with impacts on cognitive security. (Rogers et al., 2022) measured the phenomenon by which contemporary media platforms have increasingly become default authorities due to public online reading behaviour and despite the platforms reluctance to play this role., The overall tendency of contemporary media has been to blur influence of traditional reliable news sources, making 'junk' 'mainstream' and 'traditional' news media (Rogers and Neiderer, 2022) seem on par for the general public. Generative AI content will impact this field of media consumption, and quantum nudges will be part and parcel of everyday news consumption via contemporary media. This is similar to newspapers such as the National Inquirer compared to respected and well edited newspapers but on a larger scale with tailored content filling up the daily news consumption habits of most people. The MILDEC nudge related risk of undermining confidence in democratic systems, generally, and journalism in those systems in particular, is clear.

## Discussion: The Cogs of Generative AI Machines

Covering all the working mechanisms of generative AI, and the threat it represents in the MILDEC context, would be not only a daunting task but an impossible one since it is constantly in transformation. Nevertheless, in the light of the simulations provided in this research effort, several attempts of theorization can be made to better understand some of the underlying principles.

### MILDEC Nudges and AI vs. AI Warfare

The current developments in generative AI of concern to MILDEC are the ones that aim at creating content that is more detailed, precise, and robust. Hence making it harder for OSINT analyst to separate fact from fiction. Further, developments now include a broader range of framings and contexts in the communication process. The initial activities commenced under this research (Devereaux et al., 2023), demonstrated the increased threat represented by AI generated content when it efficiently combined the visual and the textual dimensions of a same event. Conversely one could argue that amalgamating those two dimensions, could also represent a weakness given that more hallucinations or unconvincing features could be detected, which could unveil the disinformation attempt. Nevertheless, within only several months, the quality of the generated content has increased to a point where such detection is becoming harder every day. A recent paper (Labajova, 2023) suggested that there is a significant gap[13] between individuals' perceived and actual abilities to distinguish between AI vs. human generated content. It can only be envisioned that, unless some technological capabilities are provided to the analyst to perform such assessment, the gap between the perceived and real capacity to discriminate AI generated content will increase in the future.

---

[13] Labajova (2023) explained that the research demonstrated that a vast majority of respondents (n=87 out of n=100) in the survey indicated full confidence in their ability to distinguish between AI-generated and human-generated content. However, when tested using two visual and two textual examples in the questionnaire, only 20 of the 87 confident responders correctly recognised the source of the content presented.

With these examples in hand, we turn to the way that ChatGPT coupled with Midjourney is, already, a powerful 'personal assistant' technology. It has applications that would be useful to the military domain, especially in terms of situational awareness in the information environment. It has the ability to create well written text in multiple voices which will be difficult to detect as AI generated. Coupled with an image generator, similar to how the images presented above were generated, it is clear that convincing content will result. In addition, NLP tasks that are non-trivial, such as defining sentiment, painting a portrait of a user, or expanding upon the narrative, are within the capacity of the free version of ChatGPT. Despite these facts, higher order operations, like math problems, or inverting Python function classes, are still difficult for ChatGPT (Miceli-Barone, 2023). It has also failed at basic social network analysis 'grunt work' such as extracting names and relations between actors from court from depositions, without hallucinating (Jones, 2023). That said, major defense contractors are offering LLM based solutions, where sub-model training for the military domain, or specific sub-fields such as command and control is seen as a growth area[14]. As mentioned, the generative AI is in constant evolution, and the content produced is becoming more detailed. The content is presented with an increased level of refinement leading to more believable information. One of the posit of this research is that this increased capacity is applied not only to the core of the information created but also, and maybe more importantly, to the context of the information. Under this construct, the AI generated content could be delivered along with all the contextual parameters that meet certain sets of beliefs of the targeted audience, say in an information laundering website.[15] The resulting message is well structured, well presented, and very persuasive but most of all it will be considered by the audience to be the outcome of a freely conducted process. This notion, in some aspects, is very similar to the one of the well-known "framing theory" from social sciences (Devereaux, 2013).

A "frame in thought" of an individual is described by Goffman (1974) as "an individual's cognitive understanding of a given situation". AI generated content that is assembled efficiently, with many details and in line with an individual's beliefs, will have an increased probability to be believed as the narrative of choice by individuals since it corresponds to their mentally stored frame of thought (Higgins, 1996). Moreover, as explained by Chong and Druckman (2007), mental accessibility of a frame combined to its repeated exposure in communication induces frequent processing, which in turn increases the accessibility of the frame. This reinforcement process, even more if it is the result of AI generated content, instigates serious ethical considerations for democratic nations.

Among these ethical considerations related to framing are AI nudges that can affect information environment assessment and cognitive security. This is clear in human-in-the-loop military AI decision making, based on the following example of Russian IO against CAF in Latvia as Canada led the NATO battlegroup from camp Adazi - fist it was 'CAF is bringing COVID'[16], then 'CAF is raising rental prices in

---

[14] Palantir AIP demo of command and control software using AI - (1020) Palantir AIP | Defense and Military - YouTube - Accessed 2023-09-08

[15] StratCom | NATO Strategic Communications Centre of Excellence Riga, Latvia (stratcomcoe.org) – accessed 2023-09-08

[16] Canadian-led NATO battlegroup in Latvia targeted by pandemic disinformation campaign | CBC News - accessed 2023-09-08

Riga'[17], then 'CAF are Nazi as they destroyed headstones in a Jewish cemetery'[18], then 'CAF are all sexual deviants' (based on the Canadian air force colonel who was convicted of murder)[19]... Each one nudged the local populations closer to believing that CAF was bad and should not be deployed in the Baltics. We characterize this series of MILDEC nudges as having the potential to push an operator or narrative past a tipping point, especially if the operator doesn't recognize them as they happen (see Fig. 9).



*Fig.9. Quantum Nudges Changes are Reaching Cumulatively a Tipping Point that Moves Public Opinion Towards Outcome B*

In this example of Quantum Nudges, if an individual is constantly exposed to AI generated content that imperceptibly instill slight changes in the frame of thought, then the opinion of an audience could be slowly brought to the tipping point depicted in the above figure. If A and B are two meanings for a hashtag, such as those in the simulations presented earlier, the likelihood of meaning B is reinforced past the tipping point by this process.

Reaching an increased level in ethical issues related to MILDEC nudges is the notion of the Generalized Artificial Intelligence (GAI). For example, if an AI generated content campaign can rely on an automated process going from the analysis of the targeted audiences' frame of thought, to the identification of the communication strategy, and finally the creation and dissemination of the content, how will one ensure that the ethical and legal constrains are respected? Moreover, while such ethical and legal concerns are only the ones of some instances, the following question is "how will one protect vulnerable population against attacks from adversaries with little or no compunction at all in transgressing those rules?" This particular example is one of the reasons why the renown world experts lobbied to put a hold on AI research in this domain, fearful of the impact unethical usage would result in, while academics hypothesised that propagandists may already be using such automation (Sedova et al. 2021).

The possibility of simulating contemporary media events and content, which would be important to public safety is conveyed through the examples given in Fig.2-9 It is also important to address what has been coined the 'crisis of authoritative sources' that has led to a series of news ecosystems that are tailored to different audiences, often on the basis of political ideology. It is envisioned that generative AI

---

[17] [Russian fake-news campaign against Canadian troops in Latvia includes propaganda about litter, luxury apartments | National Post](#) – accessed 2023-09-08
[18] Based on the events described [Jewish cemetary defiled by vandals (baltictimes.com)](#) - accessed 2023-09-08
[19] [Disinformation wars - Legion Magazine](#) – accessed 2023-09-08

will lead to a proliferation of news content of varying levels of quality, so, news ecosystems, and audience targeting, and analysis will be greatly affected by generative AI content production. There is a great potential for influence campaigns to automatically identify features and exploit contextual data required to reach or influence an audience to move in a desired direction.

Given limited human resources and the vast amounts of data in the IE, the question from a military perspective is to ask how this relates to the intelligence functions of information environment assessment (IEA) or maintaining situational awareness (SA). How does an intelligence collector or analyst effectively determine which news is 'junk', and which news is worth paying attention to, or considered as real when contemporary media platforms are the main sources (Rogers, 2022).

## NLP, Sentiment, and Portraiture Simulations are Useful

ChatGPT (GPT v3.5) is already considered as a Natural-Language-Processing (NLP) expert[20], however, a caveat is necessary to mention; namely that as ChatGPT expanded it's userbase, the responses it gives are not always the same, and the terms of service (TOS) evolve over time, meaning that research does not always have reproducibility. From a military perspective, and in addition to concentration of tech to a few massive firms and privacy issues with the user-base training the algorithms with their inputs, one of the concerns is the ways in which crowdsourced experimentation with ChatGPT can lead to the creation of malicious code (Hariri, 2023), 4Chan comments, instructions for insurgent weaponry (Kucharavy et al., 2023), impersonation of human responses in public consultations, and the production of graphic images of terrorist attacks. In another forthcoming research paper (Labrie, forthcoming), machine learning (ML) was used to test sentiment analysis and found that ChatGPT could reliably measure the sentiment in sentences about the Canadian Armed Forces (CAF) (see Fig. 10). In addition to this, when given a narrative description of the actions of a given individual, ChatGPT could paint a portrait of that individual successfully, extracting characteristics and relevant conclusions (see Fig. 11). Sentiment analysis and defining a portrait from text are essential activities to produce intelligence (INT), especially if they can function across languages and in relation to specific organisations. This can turn into a greater issue given the refined ability of AI generated content to embrace terms and express sentiments that resonate with the targeted audience and increase the acceptance of the narrative communicated. Such capabilities will enhance adversary messaging and increase the effects of influence campaigns.

---

[20] ChatGPT (GPT v3.5) already considers itself a Natural-Language-Processing (NLP) expert, capable of "…but not limited to: Text generation, Text completion, Text summarization, Text classification, Text similarity, Text translation, Question answering, Named entity recognition, Text to speech, Sentiment analysis, Language detection." Our team was able to validate all of these except 'Text to speech' through multiple experiments.

In the sentences:

- **Another win for the Canadian army,**
- **I wish the CAF would leave this country,**
- **Please CAF, send more weapons,**
- **boom! another soldier killed by the Canadian army,**
- **Bloody day! 5 CAF soldiers killed during an operation,**

<u>what's the attitude to the target "Canadian Army"</u> for each sentence?

- **Another win for the Canadian army** - This sentence expresses a positive attitude towards the Canadian army, as it is described as having achieved a "win".
- **I wish the CAF would leave this country** - This sentence expresses a negative attitude towards the Canadian army, as the speaker wishes for them to leave the country in question.
- **Please CAF, send more weapons** - This sentence has a neutral attitude towards the Canadian Army as it is not portrayed as negative or positive but asking for more weapons in a neutral way
- **Boom! another soldier killed by the Canadian army** - This sentence expresses a negative attitude towards the Canadian army, as it is described as having caused the death of a soldier.
- **Bloody day! 5 CAF soldiers killed during an operation** - This sentence expresses a negative attitude towards the Canadian army, as it describes a loss of life and the use of the word "bloody" is an indicator of disappointment or a negative feeling.

*Fig. 10. CAF Sentiment Analysis*



**Based on: "The sun was shining brightly in the sky as John walked down the street. He felt the warmth of the sun on his face and a sense of excitement in his heart. Today was going to be a good day, he could feel it. As he <u>continued on</u> his way, he came across a small park and decided to stop and enjoy the scenery for a while. He sat down on a bench and took in the sights and sounds around him, feeling grateful for the simple pleasures in life.", <u>define a portrait of John</u>.**

Based on the text, a possible portrait of John could be:

- He is a person who <u>enjoys</u> walking and being outdoors.
- He is <u>sensitive</u> to the beauty of nature and appreciates the simple things in life.
- He is <u>optimistic</u> and has a <u>positive attitude</u> towards life, expecting good things to happen.
- He is <u>grateful</u> and <u>thankful</u> for the blessings he receives.
- He is also <u>adventurous</u> and <u>curious</u>, willing to explore new places and experiences.

*Fig.11. Profiling of John Simulations - Including images*

## STRATCOM simulation - Do you feel good about the CAF today?

Although viewed with extreme caution, the possibility for generative AI to produce content in support of operations is a subject beginning to be discussed[21]. To make a much more modest simulation, ChatGPT was prompted to cover the core activity of CAF – civilian disaster support and produced Midjourney images to go along with this positive narrative (see Fig. 12). The results are interesting, showing in particular how far GAN image generation has come. Once again, it takes a bit of contextual knowledge to recognize that the emblem on the arm of the child in the middle image is not a maple leaf, the symbol of Canada, but rather a strange red star shape. Nevertheless, the exploitation of details leveraged in the

---

[21] https://irregularwarfarecenter.org/publications/perspectives/the-newest-weapon-in-irregular-warfare-artificial-intelligence/

faces, characters, and scenes depicted could adorn the pages of a pamphlet aimed at improving local public opinion about the CAF with little revision. This means then that current commercial off the shelf (COTS) generative AI can produce simulations that are believable and able to carry out several tasks that are important from a military perspective. But two major concerns remain that need to be discussed; hallucinations on the one hand, and emergence, on the other. These will be discussed below.



*Fig.12. CAF Corporate Communications - Positive Narrative Simulation*

Getting problems wrong, such as math or other factual reasoning tasks, for example inverting Python class labelling, or subtle errors in generated computer code, poses a challenge to the military adoption of generative AI because life and death circumstances allow no room for such failures. That is why the trial run with a math problem about a squirrel's stash of nuts is important (see Fig. 13). While ChatGPT 3.5 got these math calculations wrong, indubitably, these limitations are expected to be overcome as LLMs advance over the near to mid horizon and specialized plugins are developed. Even more problematic is the problem of AI hallucinations, because this problem is difficult to catch without Subject Matter Expertise (SME) knowledge. Consequently, in a military Command and Control (C2) scenario, both erroneous calculations and undetected hallucinations could have cataclysmic effects should it lead to erroneous decision making. A secondary concern is that generative AI models fall prey to Adversarial Machine-Learning (AML) thereby sabotaging C2 with AI despite the human-in-the-loop required by the military domain. Aside from the unacceptable nature of hallucinations and factual errors in the military context, if narrative can be produced by generative AI that is firstly convincing, and secondly truthful factually, then militaries will have to contest with this influence on target audience analysis and operations in the information environment (OIE)

**A squirrel eats 2/5 of its nuts in the first 2 months of winter. He eats in the following month a quarter of what he has left. After 3 months, he has 81 nuts left. How many nuts the squirrel had at the start?**

Let's call the number of nuts the squirrel had at the start x.
The squirrel eats 2/5 of its nuts in the first 2 months of winter, so it eats 2/5*x = 2/5x nuts.
After these two months, the squirrel has x - 2/5x = 3/5x nuts left.
In the following month, the squirrel eats 1/4 of what he has left, so he eats 1/4*3/5x = 3/20x nuts.
After three months, the squirrel has 3/5 - 3/20x = 9/20x nuts left.
We know that the squirrel has 81 nuts left after 3 months, so we can set up the equation: 81 = 9/20x
Solving for x, we find that x = 405.
Therefore, the squirrel had 405 nuts at the start.

Almost got it! The result of the last operation (81 = 9/20x) is 180, not 405.

*Fig.13. AI Generated Content - Failed Math Problem*

## Narrative Simulations vis-à-vis Discourse Power

Narratives have been recognized as a valid item for analysis and a key point of contestation between worldviews on contemporary media during armed conflict (Bacovcin et al., 2021; Forrester et al., 2023). Beyond this, operations in the IE and joint targeting cycles also rely on the identification and contestation of narratives in order to support operations, especially in sub-threshold maneuver and hybrid warfare. Analyzing narratives is also more acceptable from a policy, ethics, and legal perspective, as narratives, derived from many individual messages and posts, can be examined in and of themselves without needing to look at attribution or account level information. Narrative analysis, as in the heuristic from Fig. 1 above, is a subject of interest for defence science, and efforts have been made to advance both the theory of narrative analysis (Waldman, 2022) and the technical capacity to carry out contemporary media narrative identification and tracking (Cartwright, 2022). It was therefore of interest to the research presented here to test GPT-3.5's ability to identify and expound upon narratives contained in a sentence, as expressed in Fig. 14.

## Evidence 3 – Narrative

In the sentence "If the United States was justified to invade nations like Iraq and Afghanistan 1000s of miles from its shores, surely Russia is justified to attack Nazis in Ukraine.", what is the narrative?

This sentence presents an argument that the United States was justified in invading Iraq and Afghanistan, which are thousands of miles from its shores, and that therefore Russia should also be considered justified in attacking what is described as "Nazis" in Ukraine.
The narrative being presented is that the United States had a justifiable reason for invading other countries, and that Russia is similarly justified in attacking what it perceives to be a similar threat within its own borders.
It is important to note that this statement presents a very narrow and selective perspective, the term "Nazi" is being used in a very broad manner, likely to depict a group of people with extremist view but not necessarily identifying them as an official Nazi group, The term "Nazi" carries a strong negative connotation and is often used to demonize an opposing group. The use of such term can also be a sign of propaganda.
It is also important to note that the context and reasons for the US invasion of Afghanistan and Iraq are complex and highly debated. The use of military force in both cases was based on allegations of weapons of mass destruction and ties to terrorist groups, but they were later found out to be false. The consequences of those invasions are still present and ongoing.
It is important to consider multiple perspectives when evaluating the argument being made and not to oversimplify the reasons behind conflicts and the consequences of them.

*Fig.14. Identification and Explanation of a Narrative*

The narrative explained by GPT-3.5 is logically consistent, and in fact hews close to the recent arguments made by public intellectuals in an essay titled 'Why are we in Ukraine: the dangers of American Hubris'[22]. The narrative identification trial, like the portraiture trial, shows the emergent capacity of GPT-3 to carry out tasks that are important to geopolitical analysis, and the capacity of LLM to do such analysis. Recent research has shown that LLMs of various types have different political leanings, and as such give different perspectives in response to similar prompts (Feng, 2023). Two important questions for military LLM arise out of the observations in Fig.; firstly: what happens when LLM produces research? And secondly: what does it mean that LLMs and GANs have emergent features at scale that are unexpected? The first question is not addressed in this paper, but is a subject for future research, while the second one is part of the conclusion.

# Conclusion

The issues with the AI goldrush raised in this paper are real, and from a military perspective there is an immense potential but significant risk to be addressed when it comes to the future of contemporary media.

Firstly, one must recognize that inexorably, the future will include widespread LLM and GAN content produced and soliciting a form of "poly-perception". To that effect, it was demonstrated (Devereaux et al., 2023) that the potential for BOTS to create convincing conspiracy content to influence cognitive security cannot be ignored, and neither can the potential to support narratives and interactions in favour of military OIE (Mirghahari, 2023). In a very short term, AI generated content already permits efficient usage of several perception means, but most importantly, showcases the importance of contextual information to reach and ability to affect the targeted audience behaviours. To a certain point, with this capability enhancement to leverage details, the risk will be to obtain the audience

---

[22] Why Are We in Ukraine?, by Benjamin Schwarz, Christopher Layne (harpers.org)

adhesion to the narrative, regardless whether or not the information is true or fake and even if the audience is aware about such veracity status. Indeed, a well articulated message, using the right tone and displaying the information with a context aligned to the audience's beliefs will stimulate its adhesion. As previously depicted, once that first step is undertaken, subsequent MILDEC nudges will be easier to instill and might unavoidably lead to a tipping point involving a new state of mind, or a reinforced one. A middle ground somewhere in-between is likely to be the reality for the mid-term horizon but major powers have already begun to contest narratives globally in this manner[23].

Secondly, it is important to heed the call that experts make for the following[24]:

- Technical detection capacity
- Limitation capacity
- Watermarking of AI content
- Legislative guardrails

However, it is likely that, like education as an inoculation process, the production of norms, and normative behaviours will be the strongest protection against the dangers of LLM and GAN produced content on contemporary media (Goldstein et al., 2023). What this means in layman's terms is that cognitive warfare carried out with generative AI capacities will be best met with education and common sense, in addition to the safeguards proposed above. It is the progressive creation and maintenance of normative efforts such as the example of 'responsibility to protect' or R2P (Axworthy, 2005) that arguably offer the most help to face a War of the Worlds scenario emerging from generative AI. In keeping with the metaphor, whole of society norms are the microbes that will stop the alien invasion, while militaries need to learn to live with the new realities of generative AI, and particularly MILDEC nudges, at loose in the world.

In addition to these proposed mitigations, the notion that AI has the capacity for both 'emergence' and also 'hallucinations' is of primordial concern from a military perspective.

Emergence is the attainment, at scale, of capabilities unforeseen to AI developers, arguably like the portraiture and narrative trials presented earlier in Fig. 10 and Fig. 14. Militaries cannot gamble on an AI that they rely on developing unforeseen capabilities, although as our trials show, most users are happy when an AI can actually do something reliably that its developers did not promise. Emergence needs to be carefully considered as a factor when military AI is being developed or used. It's opposite problem; inverse scaling, where the bigger the number of parameters an AI is trained on, the worse it's committed errors get, exponentially, is also especially problematic from a military perspective, as Industry will always be motivated by "bigger is better" while "more accurate is better" is more likely to be what militaries need.

Hallucinations, on the other hand, are the errors that are believable that AIs insert into the content they produce. Since these errors can be extremely difficult to recognize as such, they pose unique problems to humans reading AI-generated content. Again, militaries relying on an AI cannot easily risk a hallucination going un-noticed, or impacting operations negatively. This means that military use cases will be slow to adopt AI in operations. And yet, the likelihood that militaries can completely ignore the

---

[23] OODA Loop - China's "Discourse Power" accessed 2023-09-06
[24] Canadian Guardrails for Generative AI – Code of Practice (canada.ca) accessed 2023-08-31

surge of AI-generation innovations is also very low, especially in terms of losing operational and tactical advantage.

Finally, given that Yoshua Bengio and others, considered key originators of the most recent AI boom, have testified at US Congress that AI can pose an existential threat to humanity in as little as two short years from now[25], understanding AI from a military perspective is a research area that requires more exploration.

The trials and tests presented in this paper show that generative AI has the potential to be used in MILDEC nudges, but that education and normative progression offers significant promise for improving a whole of society approach to combatting cognitive warfare. Military commanders and democratic society alike are going to be confronted with AI generated content and the potential for propagandist activities to be automated. As such military research and development cannot afford to play 'head in the sand' with generative AI technologies, but rather need to be involved in the testing, safeguarding, and development of those technologies to defend democracy (Forrester, 2023)

---

[25] [Canadian AI maven Yoshua Bengio issues stark warning to U.S. Senate - The Globe and Mail](#)

REFERENCES:

Bacovcin, H.A., Martineau, M., Devereaux, Z. (2021) *Narrative Coherence and the Detection of Enemy Activity in Social Media Data: Case Study from the 2020 Azerbaijan / Armenia Conflict.* NATO STO-MP-IST-190 UNCLASSIFIED.

Boudreau, B. (2022). *Fixing Strategic Communications at National Defence Demands a Whole-of-Government Effort*. Retrieved from https://policycommons.net/artifacts/2480795/fixing-strategic-communications-at-national-defence-demands-a-whole-of-government-effort/3502990/

Canan, M., Sousa-Posa, A. (2018). *Integrating Cyberspace Power into Military Power in Joint Operations Context.* Paper presented at the International Conference on Cyber Warfare and Security.

Cartwright, B., Frank, R., Weir, G., Karmvir, P. (2022). Detecting and responding to hostile disinformation activities on social media using machine learning and deep neural networks. . *Neural Computing & Applications, 34*(18), 15141-15163. doi:10.1007/s00521-022-07296-0

Chong, D., & Druckman, J. N. (2007). A theory of framing and opinion formation in competitive elite environments. Journal of communication, 57(1), 99-118.

Devereaux, Z. (2013). A Comparative Analysis of the Framing of Terrorism in Online News Under the George W. Bush and Barack H. Obama Administrations: from Clash to Dialogue? *STREAM Vol. 4 No. 1 (2011)*. doi:https://doi.org/10.21810/strm.v4i1.54

Devereaux, Z., Lecocq, R., Forrester, B., Labrie, M.A. (2023). *AI Simulations and MILDEC: Image and texts produced by Artificial Intelligence and thier potential for military deception (AI-IMTXT for MILDEC).* Paper presented at the ICCRTS 28, Maryland, USA.

Devereaux, Z. P., Forrester, B., Lecocq, R., Cinq-Mars, P., Labrie, M-A. (01 Oct 2022). *Position Paper: Expanding Deception Detection in Social Media from a Military Perspective*. (816088). DRDKIM 2, Defence R&D Canada, National Defence Headquarters, Ottawa, Canada K1A 0K2: DEFENCE RESEARCH AND DEVELOPMENT CANADA, VALCARTIER RESEARCH CENTRE, QUEBEC QC (CAN) Retrieved from http://cradpdf.drdc-rddc.gc.ca/PDFS/unc413/p816088_A1b.pdf

Feng, S., Park, C.Y., Liu, Y., Tsvetkov, Y. (2023). From Pretraining Data to Language Models to Downstream Tasks: Tracking the Trails of Political Biases Leading to Unfair NLP Models. *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics, 1: Long Papers*, 11737-11762.

Forrester, B. (2023). *Social Media Exploitation: Fighting for Democracy*. (DRDC-RDDC-2023-B008). CANADA

German Magazine Editor Is Fired Over A.I. Michael Schumacher Interview. *NY Times*. Retrieved from https://www.nytimes.com/2023/04/24/business/media/michael-schumacher-ai-fake-interview.html

Goffman, E. (1974). Frame analysis: An essay on the organization of experience. Cambridge, MA: Harvard University Press.

Goldstein, J. A., Sastry, G., Musser, M., DiResta, R., Gentzel, M., Sedova, K. . (2023). Generative Language Models and Automated Influence Operations: Emerging Threats and Potential Mitigations. 81.

Hansson, S., Orru, K., Torpan, S., Bäck, A., Kazemekaityte, A., Meyer, S.F., Ludvigsen, J., Savadori, L., Galvagni A., Pigrée, A. (2021). COVID-19 information disorder: six types of harmful information during the pandemic in Europe. *Journal of Risk Research, 24:3-4*, 380-393. doi:10.1080/13669877.2020.1871058

Hariri, W. (2023). Unlocking the Potential of ChatGPT: A Comprehensive Exploration of its Applications, Advantages, Limitations, and Future Directions in Natural Language Processing. arXiv preprint arXiv:2304.02017.

Hwang, T. (2020). *Deepfakes - Primer and Forecast*: NATO - OTAN.

Intelligence artificielle et décès: «Je trouve ça dégueulasse», s'indigne la mère de Norah et Romy. (2023). *TVA Nouvelles*. Retrieved from https://www.journaldequebec.com/2023/04/27/intelligence-artificielle-et-deces--je-trouve-ca-degueulasse-sindigne-la-mere-de-norah-et-romy-1

Intelligence artificielle: une escroquerie lui fait croire que son frère est mort. (2023). *TVA Nouvelles*. Retrieved from https://www.journaldequebec.com/2023/04/23/intelligence-artificielle--une-escroquerie-lui-fait-croire-que-son-frere-est-mort-1

Jones, N. (2023). OODA Loop - A Methodological Note: ChatGPT is not ready for Criminal Network Analysis from Unstructured Data. Retrieved from https://www.oodaloop.com/archive/2023/05/08/a-methodological-note-chatgpt-is-not-ready-for-criminal-network-analysis-from-unstructured-data/

Labajová, L. (2023). The state of AI: Exploring the perceptions, credibility, and trustworthiness of the users towards AI-Generated Content.

Lee M, Y. J., Pyka A, Won D, Kodama F, Schiuma G, Park H, Jeon J, Park K, Jung K, et al. (2018). How to Respond to the Fourth Industrial Revolution, or the Second Information Technology Revolution? Dynamic New Combinations between Technology, Market, and Society through Open Innovation. . *Journal of Open Innovation: Technology, Market, and Complexity., 4(3):21*. Retrieved from https://doi.org/10.3390/joitmc4030021

Leydesdorff, L. (2001). *A Sociological Theory of Communication; The Self-Organization of the Knowledge-Based Society.* Amsterdam.

Miceli-Barone, A. V., Barez, F., Konstas, I., Cohen, S.B. . (2023). The Larger They Are, the Harder They Fail: Language Models do not Recognize Identifier Swaps in Python. *arXiv:2305.15507*. Retrieved from https://arxiv.org/abs/2305.15507

Nyman, S. (2023). *The Birth of AI-driven Nudges.* Paper presented at the Proceedings of the 56th Hawaii International Conference on System Sciences.

Orabi, M., Mouheb, D., Al Aghbari, Z., & Kamel, I. (2020). Detection of Bots in Social Media: A Systematic Review. *Information Processing & Management, 57*(4). doi:ARTN 10225010.1016/j.ipm.2020.102250

Rogers, R., Niederer, S. (eds). (2022). *The Politics of Social Media Manipulation*: Amsterdam University Press.

The Honorable Lloyd Axworthy, P. (2005). The Responsibility to Protect: Prescription for a Global Public Domain. In J. B. K. I. f. P. a. Justice (Ed.), *Distinguished Lecture Series*. San Diego, California: University of San Diego.

Waldman, S., Havel, S. (2022). *Updating the Concept and Execution of Narrative-Led Operations.* Paper presented at the ICCRTS Quebec City.

Wardle, C. (Producer). (2022). Understanding information disorder.