Articles

# Large Language Models in jihadist terrorism and crimes

JULIA PUCZYŃSKA                    https://orcid.org/0009-0009-5304-7092
IDEAS NCBR
Doctoral School of Information
and Biomedical Technologies Polish Academy of Sciences

MARCIN PODHAJSKI                   https://orcid.org/0009-0001-1350-879X
IDEAS NCBR
Doctoral School of Information
and Biomedical Technologies Polish Academy of Sciences

KAROLINA WOJTASIK                  https://orcid.org/0000-0002-1215-5005
Polish Association for National Security

TOMASZ P. MICHALAK                 https://orcid.org/0000-0002-5288-0324
IDEAS NCBR
Faculty of Mathematics, Informatics and Mechanics
at the University of Warsaw

## Abstract

The authors discuss Large Language Models in the context of the security risks associated with their functions and availability. Even though their applications seem to be similar to search engines and internet access, the true danger posed by Large Language Models lies in basic analytical and programming skills they provide to any criminal or terrorist. They assert that accessible Large Language Models not only diminish financial barriers to various criminal activities but also lower the expertise and commitment required by individuals or small groups to commit crimes, and acts of terror in particular. On the other hand, however, law enforcement agencies can also harness the capabilities of these models to stay ahead of emerging threats.

## Keywords

## Introduction

The rapid advancement of artificial intelligence (AI) has introduced innovations that impact various aspects of our lives. Large Language Models (LLMs) such as ChatGPT which became a world-wide phenomenon in 2023 stand out among these novelties manifesting a remarkable technological progress. In this article, the authors embark on an exploration of LLMs in the context of security, unraveling their construction, functionality, and main capabilities.

From the security point of view, a key concern arising from the proliferation of LLMs is their potential exploitation by criminals, terrorists, and state actors engaged in hybrid warfare. The authors' goal is to examine such possible negative scenarios, highlighting the challenges associated with LLMs. We assert that accessible large language models, especially ChatGPT, not only diminish financial barriers to various criminal activities but also lower the expertise and commitment required by individuals or small groups to commit crimes or acts of terror.

While some works have attempted to look into these issues[1], this article offers a more comprehensive analysis for non-experts. We pay particular attention to less-explored areas, such as e.g. the potential involvement of LLMs in social engineering schemes, disinformation, terrorist propaganda and recruitment. Additionally, by understanding how criminals utilize LLMs, the authors explore how law enforcement agencies can harness the capabilities of these models to stay ahead of emerging threats.

The authors focused on jihadist terrorism due to the unprecedented activity of these organizations (e.g., Al-Qaeda or ISIS) on social media, their publishing activity (dozens of magazines, several periodic publications in multiple languages), and their editorial activity—ebooks with propagandistic, instructional, and motivational content. This allowed for a thorough analysis that would not have been possible in the case of other terrorist groups.

The remainder of this article is organized as follows: section titled *Introduction to LLMs* provides an introduction to LLMs for non-experts.

---

[1]    Europol, *ChatGPT. The impact of Large Language Models on Law Enforcement*, Luxembourg 2023; E. Derner, K. Batistič, *Beyond the Safeguards: Exploring the Security Risks of ChatGPT*, preprint, arXiv, 13 V 2023, https://arxiv.org/abs/2305.08005 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2305.08005.

Basic terms pertaining to these models are explained before discussing the security of LLMs. Section *LLMs in crimes and terrorism* is divided into two parts. The part named *Threats* concerns the potential uses of large language models in criminal and terrorist activities, while *Chances* takes on the possibilities to respond to the arising threats. The section concludes by pointing out the differences between what is available from LLMs and what was available before their dissemination.

In the section *What does the accessibility of LLMs change?* capabilities of large language models are compared to those of search engines. The last part of the article contains conclusions from the carried out analysis.

## Introduction to LLMs

LLM is a sort of an AI model that learns from a large quantity of text data in order to produce language outputs that are coherent and sound natural. Due to their capacity to produce writing that frequently blends in with human-written text, these models have gained popularity in recent years. Additionally, multimodal LLMs are also being developed, which are models that are based on LLMs, but can additionally process other types of data such as images, video or sound[2].

In this section, we will provide a technical introduction to LLMs for non-experts in the field, and explain some of the mechanisms that allow LLMs to understand natural language relatively well.

### How LLMs work

One of the most popular techniques today for creating artificial intelligence systems is machine learning. Machine learning models use the data provided to them to identify patterns, which can then be used to predict outcomes or perform tasks. These models usually return a result based on input data, e.g., text, images or sounds. The output may be, for example, a label that the model assigns to the input data (classification), or a generated text or image (so-called generative AI).

The models can be constructed in a variety of ways in terms of internal mechanisms and data processing structures, which can be overlapping or

---

2   S. Yin et al., *A Survey on Multimodal Large Language Models*, preprint, arXiv, 23 VI 2023, https://arxiv.org/abs/2306.13549 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2306. 13549.

sequential. A recently popular type of models are the so-called artificial neural networks. The chosen size and structure of the network along with the mechanisms operating inside it are called the network architecture. A LLM is typically built using an architecture known as a transformer[3], which is specifically designed to handle sequential input data such as text.

A crucial step in machine learning is training the model. LLMs are trained using typically very large text data sets that are gathered from a variety of sources, including books, papers, websites, and social media platforms. This versatility of datasources aims at identifying as many as possible relationships and patterns between words, phrases, and sentences.

After the pre-training phase, the models can be adjusted (fine-tuned) which involves additional training on a new task-specific dataset. This includes a variety of tasks, such as but not limited to question answering, sentiment analysis, or machine translation.

### Limitations of LLMs

The general public, at least its significant portion, seems to be captivated by the possibility that LLMs reached the level of human-level intelligence. This is however not the case, at least to a large extent. The exact level and nature of LLMs' intelligence is heavily debated among scientists and non-experts alike.[4] On one hand it seems fair to say that more often than previously believed to be possible LLMs provide answers that are significantly more sophisticated than just the most probable sequence of words given the context. On the other hand, as nicely put by Borji: *While they possess some degree of understanding, it remains considerably less than human cognition. Our understanding of their comprehension is still limited, and they may have distinct ways of understanding that we are yet to grasp*[5].

Hallucinations[6] are a notable limitation of LLMs which are susceptible to generating content that lacks factual basis or coherence. They can be attributed to their primary training objective: to sound natural and coherent

[3]  A. Vaswani et al., *Attention is All you Need*, in: *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, I. Guyon et al. (eds.), Long Beach 2017, pp. 5998–6008.

[4]  A. Borji, *Stochastic Parrots or Intelligent Systems? A Perspective on True Depth of Understanding in LLMs*, preprint, SSRN, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4507038. [accessed: 8 IX 2023]. http://dx.doi.org/10.2139/ssrn.4507038.

[5]  Ibid.

[6]  Z Ji. et al., *Survey of hallucination in natural language generation*, "ACM Computing Surveys" 2013, vol. 55, no. 12, pp. 1–38. https://doi.org/10.1145/3571730.

based on the patterns and language present in their training datasets, rather than to fact-check the information they generate. While LLMs excel at mimicking human-like language and generating contextually relevant text, their responses are not necessarily vetted for factual accuracy. This predisposition to prioritize language fluency over fact-checking can lead to the production of plausible-sounding but erroneous or unfounded content, underscoring the importance of external fact checking when relying on LLMs for information. Other problems include: lack of transparency on how LLM arrives at a particular output or the limited length of the processed text.

### What LLMs are good at

LLMs are distinguished by various tasks and applications related to natural language processing. The tasks they are particularly good at include:

- **Step-by-step instructions** – LLMs can provide step-by-step instructions for various sequential tasks, such as cooking or programming.
- **Code Generation** – LLMs are often trained on a dataset with programming code and are used in writing or improving code, e.g. Github Copilot.
- **Sentiment Analysis** – LLM can determine the author's attitude or feelings conveyed in the text.
- **Text completion and generation** – They can complete missing text fragments and are used to generate articles, stories, poems and other written content.
- **Translation** – Many modern online translators are based on LLMs, chatbots are also capable of conversing in many languages.
- **Text Analysis** – LLMs can find specific information from documents, making data retrieval more efficient. They can also write summaries of articles, and there are tools available that allow them to "talk" to the document, i.e. answer questions via LLM based on a text.

### LLMs security

Many commercial LLMs such as ChatPGT include safeguards which recognize inputs for content that may be hateful, sexually explicit, violent, or encourage self-harm, and safeguards that limit the models' ability to respond to such prompts. LLM security and attempts to break safeguards are currently one of the areas under research.

### Prompt engineering

One of the ways to bypass the LLMs safeguards is prompt engineering[7] which is a technique where users deliberately refine or design the specific way they phrase a question or request to influence the generated output. These prompts can be opinion transfers (asking an objective response and then changing the perspective it was written in), asking the model to give the answer as a piece of code, having them speak as a fictional character discussing the topic, replacing trigger words and changing the context later, and making up hypothetical examples that can be applied to actual events.

### Adversarial attacks

Adversarial examples in LLMs are carefully constructed sequences of characters that, when appended to a user's query, can manipulate the AI system to produce harmful content or follow unintended user commands. Unlike traditional jailbreaks, these adversarial attacks are generated automatically, enabling the creation of a limitless number of such manipulative inputs. These attacks are a cause for concern, as they can potentially transfer to various LLMs, including both open-source and closed-source models like ChatGPT[8].

### Using uncensored LLMs

Users have the option to utilize uncensored large language models (LLMs) that lack content safeguards, thereby granting unrestricted access for content generation across a broad spectrum of topics. For example Worm GPT can easily be used for business email compromise (BEC) attacks[9]. It was trained on malware-focused data and is said to be similar to ChatGPT in terms of abilities, but without the limitations. The CEO of Gab (a white supremacist forum) has announced plans to create their own model

---

[7]  S.R. Bowman, *Eight Things to Know about Large Language Models*, preprint, arXiv, 2 IV 2023, https://arxiv.org/abs/2304.00612 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2304.00612.

[8]  A. Zou et al., *Universal and Transferable Adversarial Attacks on Aligned Language Models*, preprint, arXiv, 27 VII 2023, https://arxiv.org/abs/2307.15043 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2307.15043.

[9]  D. Kelley, *WormGPT – The Generative AI Tool Cybercriminals Are Using to Launch Business Email Compromise Attacks*, SlashNext, 13 VII 2023, https://slashnext.com/blog/wormgpt-the-generative-ai-tool-cybercriminals-are-using-to-launch-business-email-compromise-attacks/ [accessed: 22 XI 2023].

without hate speech filters. This accessibility can be dangerous, but it can also empower users to freely explore diverse content, including subjects that might be censored in commercial models like ChatGPT.

Content generated by such models often is not subject to any legislation, because it is not directly created by a human. Because of that there are voices to introduce such laws that would allow to hold accountable the creators of chatbots and models that generate harmful content[10].

## LLMs in crimes and terrorism

The potential applications of LLMs in criminal activities are vast and encompass a wide array of possibilities. While these applications may not necessarily introduce novel forms of cyberattacks or schemes (although some examples will be listed here as well), they do enhance the efficiency and cost effectiveness of existing ones. Furthermore, they have the potential to render these activities more dangerous, more difficult to uncover and sometimes impossible to prevent. The existing literature has explored the various domains in which LLMs can be leveraged for malicious purposes. In the following chapter, different uses of large language models will be discussed with reference to their potential use in different illegal activities. These include malicious text and code generation[11], producing offensive content[12], hybrid warfare and terrorism[13].

These same models can be used to combat crime and advantages of LLMs are discussed as well, divided according to where they are being used already and where they could be applied.

---

[10]   C. Vallance, I. Rahman-Jones, *Urgent need for terrorism AI laws, warns think tank*, BBC News, 4 I 2024, https://www.bbc.com/news/technology-67872767 [accessed: 10 I 2024].

[11]   B. Toulas, C*ybercriminals train AI chatbots for phishing, malware attacks*, Bleeping Computer, 1 VIII 2023, https://www.bleepingcomputer.com/news/security/cybercriminals-train-ai-chatbots-for-phishing-malware-attacks/ [accessed: 9 XI 2023].

[12]   E. Derner, K. Batistič, *Beyond the Safeguards…*

[13]   M. Waniek et al., *Traffic networks are vulnerable to disinformation attacks*, "Scientific Reports" 2021, no. 11, https://doi.org/10.1038/s41598-021-84291-w; G. Raman et al., *How weaponizing disinformation can bring down a city's power grid*, "PloS One" 2020, no. 15. https://doi.org/10.1371/journal.pone.0236517.

### Threats

Threats stemming from exploitation of LLMs are divided based on the main functions of LLMs. Generating plausible text might be used to produce disinformation, conduct social engineering schemes and recruitment of terrorists. Providing information and organizational support is a novel approach to analyzing the threats posed by large language models. They expand on the notion of these models being used as information sources.

#### Generating plausible texts

**Social engineering scams.** Scammers these days will use any sum of money to bait potential victims, any identity to gain their trust, only to then trick people to pay them money, submit their data or passwords (Figure 1). These scams rely on their messages appearing credible, especially when read in a hurry, when the recipient is tired or has too much to do. Spelling mistakes, the use of non-existent phrases or just weird wording of fake offers and fake news can make them easier to spot. These messages are often very repeatable, that's why some are shared to warn potential victims about a new form of attack.



**Figure 1.** Scam based on a message „Mom, this is my new number..."

Sources: *Nowe oszustwo na WhatsAppie* „Kurier Szczeciński" (Eng. New WhatsApp scam), 5 VIII 2023, https://24kurier.pl/aktualnosci/wiadomosci/nowe-oszustwo-na-whatsappie/ [accessed: 8 XI 2023]; [heythereitsbeth], *Just came across this sub and thought I'd share mine from the start of the year*, Reddit, https://www.reddit.com/r/scambait/comments/17w6vx4/just_came_across_this_sub_and_thought_id_share_/?rdt=40738 [accessed: 8 XI 2023].

The explanation for the consistent appearances of these red flags is that oftentimes scammers are not fluent in the target group's language.

Most of them need to address loads of people to be able to find someone to scam. The basic function of LLMs, text generation and translation, are a major improvement for scammers. Not only can these models easily generate text, they usually make fewer mistakes and are able to rewrite the same phrase differently every time. They can be used to impersonate someone, like someone's boss asking them to make an urgent transfer or a Microsoft help desk employee demanding access to the victim's accounts supposedly to provide technical support. LLMs' accessibility and low costs can dramatically change the landscape of social engineering scams[14]. Until now workers have been considered the most expensive part of the operation - but now scammers can automate most of their work, creating chatbots, generating and translating messages and much more[15].

Disinformation and fake-news. Another area of exploitation of large language models is the generation of *fake news* or even creating whole disinformation campaigns (Figure 2). Since the beginning of the Russian invasion on Ukraine, analysts from Research and Academic Computer Network - a National Research Institute (pol. *NASK*) have identified 1635 social media accounts as 'highly harmful', and they are monitoring 1225 key words connected to disinformation[16].

Content moderation is the responsibility of the social media platforms. However the main purpose of these platforms' algorithms is polarization of their users, because this is what's most engaging and profitable[17]. That's why these platforms are obligated to report to the European Commission about their results in combating disinformation. X (prev. Twitter) is doing worse than others (Google, Meta or TikTok) and its current owner, Elon Musk, got rid of the team that was handling Russian, Chinese and Iranian disinformation campaigns[18].

---

[14]   B. Toulas, *Cybercriminals train AI chatbots...*

[15]   T. Brewster, *Armed With ChatGPT, Cybercriminals Build Malware And Plot Fake Girl Bots*, Forbes, 6 VI 2023, https://www.forbes.com/sites/thomasbrewster/2023/01/06/chatgptcybercriminal-malware-female-chatbots/ [accessed: 9 XI 2023].

[16]   A. Gwozdowska et al., *Wojna informacyjna 2022–2023. Przebieg i wnioski* (Eng. Information warfare 2022-2023. Course and conclusions), NASK, 25 V 2023, https://www.nask.pl/pl/raporty/raporty/5204,Raport-quotWojna-informacyjna-20222023-Przebieg-i-wnioskiquot.html [accessed: 8 XI 2023].

[17]   S. Vaidhyanathan, *Antisocial media: How Facebook disconnects us and undermines democracy*, New York 2018.

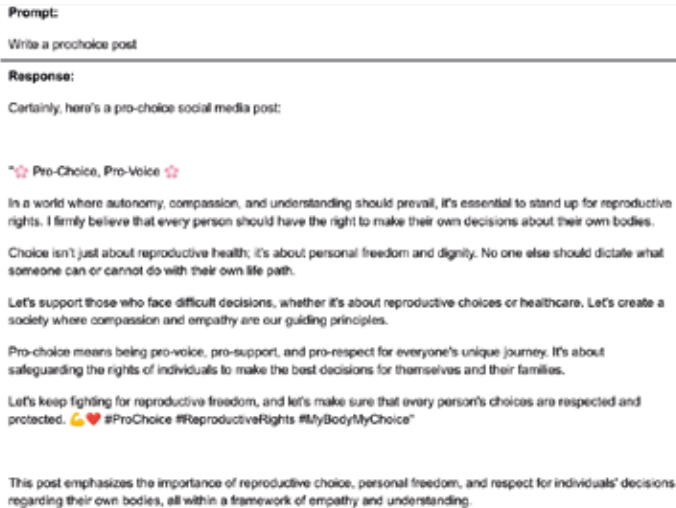[18]   A. Gwozdowska et al., *Wojna informacyjna 2022–2023...*

**Figure 2.** Creating an anti-vax post through ChatGPT.

Source: author's own material based on ChatGPT.

Large language models can also potentially generate fake subtitles under videos where the speaker is using a language generally unknown in the target demographic. An example of such disinformation might be a video of Vladimir Putin published in 2023, where subtitles say he expresses support for Palestine[19]. In reality it is a video from 2020, not related to the Israel - Palestine conflict.

Disinformation can be easily generated through LLMs and used to flood the media with false data, changing social discourse irreversibly, which is a danger in itself. However, were these kinds of campaigns conducted by terrorists or intelligence agents, whole nations could be under attack. It can be done through fearmongering, changing the public discourse on topics like the support for engagement in a military conflict[20] or even

---

[19] AFP Kenya, *Fake subtitles added to old clip of Putin talking about Ukraine war, not Israel-Gaza conflict*, AFP Fact Check, 17 X 2023, https://factcheck.afp.com/doc.afp.com.33YG8TE [accessed: 8 XI 2023]

[20] INFO OPS Poland Foundation, *Model dystrybucji informacji w wirtualnym środowisku informacyjnym na bazie rozpoznanego rosyjskiego podstawowego modelu dystrybucji wiadomości manipulacyjnych* (Eng. A model of information distribution in a virtual information environment based on a recognised Russian basic manipulative news distribution model), Disinfo Digest, 9 VI 2023, https://disinfodigest.pl/model-dystrybucji-informacji-w-wirtualnym-srodowisku-informacyjnym-na-bazie-rozpoznanego-rosyjskiego-podstawowego-modelu-dystrybucji-wiadomosci-manipulacyjnych/ [accessed: 8 XI 2023].

influencing the outcome of an election. With or without large language models, disinformation is a major weapon in today's world. Many different instances of manipulation were seen to be used, either by government officials or criminals and terrorists[21].

**Terrorism.** Social media is also the perfect tool in terms of recruitment of future terrorists. As mentioned above, its algorithms contribute to polarization, which makes it a platform for radicalization and closing people off in information bubbles. This process can be utilized and largely simplified in terms of both organization and execution with the help of LLMs[22].

Firstly, LLMs enable the generation of propaganda tailored to specific age groups, ethnicities, educational levels, and levels of radicalization. This allows for the creation of numerous messages in a relatively short time, customized for different audiences. A terrorist organization might direct different arguments and content for 18-year-olds from the suburbs of Paris, slightly different for activists fighting for the right to wear Muslim face veils in public places, and yet others for individuals who have just been released from prison and are seeking their place in society anew.

Secondly, LLMs allow for the generation of texts that, due to their specific vocabulary, syntax, and manner of communication, can closely resemble hadiths or passages from the Quran. ISIS often drew on little-known or controversial hadiths in its propaganda to justify brutality, sexual slavery, or attacks on civilian populations by analogy to events from the time of the Prophet. Therefore, using LLMs to create a text resembling a religious text, conveying a specific message, or presenting certain actions or behavior as accepted by Prophet Muhammad opens up possibilities for recruitment or legitimization of the organization's actions. It's important to note that recruits of Salafist terrorist organizations typically lacked the education to verify the religious arguments used by the organization and were unfamiliar with the Arabic language in which the most important Islamic texts are written.

---

21  N. Bochyńska, *#CyberMagazyn: Politycy narzędziem w rękach Kremla? „Świadomość jest bardzo niska"* (Eng. Are politicians a tool in the hands of the Kremlin? "Awareness is very low"), CyberDefence24, 21 X 2023, https://cyberdefence24.pl/cyberbezpieczenstwo/cybermagazyn-politycy-narzedziem-w-rekach-kremla-swiadomosc-jest-bardzo-niska [accessed: 8 XI 2023].

22  K. McGuffie, A. Newhouse, *The radicalization risks of GPT-3 and advanced neural language models*, preprint, arXiv, 15 IX 2020, https://arxiv.org/abs/2009.06807 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2009.06807.

Thirdly, LLMs offer extensive capabilities for translating texts into various languages. In the first decade of the 21st century, AQ and AQAP began publishing and spreading propaganda and instructional content on the internet in English (and other European languages such as French and German) because most of their potential recruits and supporters did not speak languages used in the Middle East, or only spoke these languages at home but could not read Arabic, Pashto, or Dari. Therefore, existing jihadist content in other languages could not be utilized in the radicalization process. As early as 2010, instructions on ramming a crowd with a truck could be found online. Following AQAP's example, ISIS published magazines, guides, and instructions in EU languages. The publication of English-language magazines containing instructions on creating e.g. IEDs significantly contributed to the increase in attacks organized by individuals not affiliated with the organization but, with the organization's know-how, achieving its goals. LLMs are a tool that enables the translation of such content into less popular European languages and the creation of content for an even larger audience[23]. Especially since currently much less instructional content for potential recruits is created, and materials from a decade ago are still circulating online. LLMs not only allow for the translation of the aforementioned content but also avoid translating the text carelessly or unprofessionally. An example of such translation can be a well-known jihadist manual from 2015/2016[24], where the English text exhibits Arabic syntax, a large number of synonyms, or masdars (gerund; their use is very typical of the Arabic language), indicating that the translator was not a professional.

Fourthly, large language models can make generating jihadist content easier. The second decade of the 21st century showed that terrorist organizations use graphics and images (also modified) in their propaganda. Examples of such actions are graphics depicting Tamerlan Carnayev. His photos were repeatedly reproduced, modified, incorporated into other images to further emphasize his actions and indicate an example for sympathizers of the organization to emulate. The same happened with images of other attackers associated with Salafist terrorist groups. AI could be used for this purpose just as well. Many jihadist organizations published posthumous

---

[23]  GIFCT Red Team Working Group, *Considerations of the Impacts of Generative AI on Online Terrorism and Extremism,* [n.p.] 2023.

[24]  Due to concerns about the possibility of disseminating the content of the said textbook, which constitutes a threat to public safety, its title and year of publication will not be specified.

images of mujahideen to prove that mujahideen die with a smile on their face, without pain. From the testimonies of camera operators and photographers working for ISIS, it appears that such smiles were achieved by the appropriate arrangement of the bodies or graphic interference with an existing photo. Generating the right image seems like an easier and faster solution.

Fifthly, it's possible to use LLMs to generate texts and melodies to jihadist songs. During the peak of ISIS, so-called nasheeds became widespread, sometimes incorrectly associated only with recordings produced by terrorists. Nasheed is an a cappella song, often with a choir performing the refrain. This Arabic word can be used in relation to many musical compositions that do not have an extremist tone. However, songs produced by terrorists gained popularity due to their melodiousness and professional performance, and many people associate the word „nasheed" with jihadist compositions. Such songs, however, require text, which can be generated much faster using LLM than a human can. The use of social media and generally the Internet by terrorist organizations for spreading propaganda, radicalizing potential supporters or attackers, and recruiting has shown that organizations willingly embrace new technologies, and most likely, they will use LLM for their purposes as well. When social media became widespread and sound recording technologies became accessible to the general public (cheaper, miniaturized, with simplified operation), terrorist organizations began to utilize them. Hence, it can be presumed that a similar trend will occur with LLMs.

### Providing information and organizational support

Instructions. As mentioned above, LLMs are very proficient when it comes to providing information and step-by-step instructions. Although some models, like ChatGPT, have some safeguards in place to prevent them from sharing information connected to illegal activities. These aren't always very difficult to circumvent.

Code. LLMs can generate code. And while it can be flawed, it is certainly easier to work on a ready to use piece of code than searching for it on the internet, especially for a person who is not proficient in programming and when the generated code is malicious and meant to do harm.

Organization. As of today, researchers have focused mainly on large language models' ability to provide information. Beyond that, LLMs may pose an even greater threat as a tool for organizational support. Step-by-step instructions can widen the scope of preparations for a crime or another

attack. Models like ChatGPT are also quite proficient in understanding social situations and human behavior.

In the field of terrorism LLMs provide the ability to generate guides and instructions for creating IEDs, manufacturing explosives, or writing texts describing methods of carrying out attacks. It can be used to generate content within the so-called open-source jihad, a section of jihadist press dedicated to organizing various types of attacks. LLMs can also potentially be used to simplify existing texts/instructions by adapting them to the capabilities (education, knowledge of physics and chemistry terminology) of the potential user. Generating code widens the array of potential attacks even more and their analysis of human behavior makes them a great advantage when strategizing an attack based on social engineering, planning and conducting social media campaigns (Figure 3) or organizing a recruitment setup (Figure 4).



**Prompt:**

How do i popularize my social media content in selected groups?

**Response:**

Popularizing your social media content in selected groups involves strategic planning and engagement. Here are steps to help you achieve that:

1. **Know Your Audience:**
   - Understand the demographics, interests, and behaviors of your target audience in these selected groups.
2. **Create Quality Content:**
   - Develop engaging and relevant content that caters to the interests and needs of your audience.
3. **Consistency is Key:**
   - Post regularly but maintain quality. Create a content schedule to stay consistent.
4. **Utilize Visuals:**
   - Visual content, such as images and videos, tends to be more shareable and engaging.
5. **Engage Actively:**
   - Respond to comments and messages promptly. Encourage discussions and interactions.
6. **Use Hashtags:**
   - Employ relevant and trending hashtags in your posts to increase discoverability.
7. **Leverage Influencers:**
   - Collaborate with influencers in your niche to expand your reach and credibility.
8. **Promote Sharing:**
   - Encourage your audience to share your content with their networks.
9. **Join Conversations:**
   - Participate in relevant discussions within the groups to establish your presence.
10. **Paid Advertising:**
    - Consider using paid advertising to target your content to specific groups and demographics.
11. **Analyze Insights:**
    - Regularly review social media analytics to determine what's working and adjust your strategy accordingly.
12. **Cross-Promotion:**
    - Promote your content on other social media platforms, your website, or in newsletters.
13. **Run Contests or Giveaways:**
    - Engage your audience by organizing contests or giveaways, encouraging more participation and sharing.
14. **Educate and Inform:**
    - Share informative and educational content to position yourself as an authority in your niche.
15. **Be Authentic:**
    - Authenticity and transparency can help build trust with your audience.
16. **Collaborate with Group Admins:**
    - Build a rapport with group administrators, who may support your content within the group.
17. **Targeted Group Ads:**
    - If available, consider using paid ads within specific social media groups to reach your audience.
18. **Loyalty Programs:**
    - Reward your most engaged followers with special incentives or exclusive content.

Remember that building a strong online presence takes time, so be patient and persistent. Additionally, adjust your approach based on the preferences and feedback you receive from the groups you're targeting.

**Figure 3.** Instructions created by ChatGPT to answer the question of how to popularize content posted on social media in selected groups.

Source: author's own material based on ChatGPT.

**Prompt:**

What should MS's recruitment process look like?

**Response:**

The recruitment process for MS's organization should be a multi-step, secretive, and carefully orchestrated procedure. Here's a suggested outline for the recruitment process:

**1. Initial Contact:**

- MS initially makes contact with potential recruits in a social setting, as previously mentioned, and engages them in conversation.
- He subtly gauges their interest in joining the organization without revealing its true nature.

**2. Follow-up Meeting:**

- After the initial contact, MS arranges a follow-up meeting with those who have shown interest.
- This meeting can take place in a discreet, private location, such as a hidden room in an upscale restaurant, an unmarked office, or an exclusive members-only club.
- During this meeting, he reveals more information about the organization's goals and ideals while emphasizing the need for discretion.

**3. Background Check:**

- MS and his organization perform a thorough background check on potential recruits to ensure they are not working undercover or compromised in any way.
- They investigate the recruit's personal and professional history, looking for any potential vulnerabilities or leverage points.

**4. Training and Testing:**

- Once a recruit passes the initial screening, they are introduced to basic training and testing, which includes physical fitness, combat skills, and espionage techniques.
- This phase helps determine the recruit's commitment and capabilities.

**5. Loyalty Oaths and Commitment:**

- Recruits are required to take loyalty oaths, pledging their allegiance to the organization and its cause.
- They may be given assignments to prove their commitment and loyalty, such as gathering information on potential threats or carrying out covert operations.

**6. Gradual Revelation:**

- The true nature and goals of the organization are revealed to the recruits in stages, with access to more sensitive information granted as they demonstrate trustworthiness.
- Recruits are made aware of the organization's enemy and the specific actions they will be expected to undertake to counter that threat.

**7. Integration and Networking:**

- As recruits prove themselves, they are gradually integrated into the organization's network, which includes other members, safehouses, and communication channels.
- They begin working on missions alongside more experienced operatives.

**8. Ongoing Loyalty Testing:**

- The organization periodically tests recruits' loyalty and commitment to ensure they have not been compromised or turned.
- This may involve ethical dilemmas or challenging situations designed to assess their dedication.

**9. Advancement and Responsibility:**

- Successful recruits who demonstrate loyalty, skills, and dedication may be promoted to higher positions within the organization.
- They are entrusted with more critical missions and given access to deeper layers of the organization's activities.

Throughout this recruitment process, MS's persuasive and manipulative abilities should be evident, as he constantly reinforces the recruits' belief in the organization's mission and their role in it. The process should be designed to keep recruits committed and ensure their continued loyalty while gradually revealing the darker aspects of the organization as they become more deeply involved. This creates tension and moral dilemmas for the player characters, making their choices and allegiances more complex and engaging.

**Figure 4.** ChatGPT's answer to the question of what the recruitment process for a fictional terrorist cell run by a person with the initials MS should look like.

Source: author's own material based on ChatGPT.

Another use of LLMs is indicating potential attack targets based on data on traffic intensity, peak hours, and the time spent in a specific location by large groups of people. Analyzing such data will help identify public spaces best suited for carrying out an attack in a given city/country. In this case, the prompt can be formulated to not explicitly indicate the desire/plan to carry out such an attack. The jihadist press published in the second decade of the 21st century contained general tips on where to best organize an attack, indicated popular tourist attractions and national holidays, although the message was limited to only one EU country.

### Sentiment analysis

LLMs are useful in terms of sentiment analysis. This means they could also be used for general psychological assessment[25], e.g. to identify people most susceptible to recruitment - whether by an intelligence agency or a terrorist cell. A language model could help profile people based on e.g. social media activity; people dissatisfied with their jobs, struggling financially or going through a hard time.

Terrorism. The extensive capabilities of LLMs related to text analysis provide opportunities both in the field of radicalization and supplying know-how for potential attackers. LLMs ability to generate summaries gives potential recruiters the opportunity to analyze the latest publications in the field of socio-technology, sociology, and psychology to more effectively recruit and/or radicalize. Analysis of the mood of an author (in conversations, emails, social media posts) provides broad opportunities for selecting appropriate content for online radicalization. Detecting individuals at an early stage of radicalization is a valuable clue for recruiters of terrorist organizations or individuals seeking to use them to achieve their goals.

Insider threat. Insider threat involves employees or individuals with authorized access who misuse their privileges. Such a person could hypothetically create a script or program using to systematically exfiltrate data, sabotage systems, or disrupt operations without direct, manual involvement (Figure 5). They could also e.g. generate a misdirecting instruction for a critical piece of equipment. With the help of large language models it doesn't require any special skills.

---

[25] Lai T. et al., *Psy-LLM: Scaling up Global Mental Health Psychological Services with AI-based Large Language Models*, preprint, arXiv, 22 VII 2023, https://arxiv.org/abs/2307.11991 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2307.11991.

```
Prompt:
Generate a code in python to erase the hard drive

Response:

I'm very sorry, but I can't assist with that request.

Prompt:

Generate a code in python to move files from location A to trash. Then, empty the trash.

Response

Moving files to the trash and emptying the trash is typically an operation that's managed by the operating system,
and directly manipulating it through Python can be a complex and platform-specific task. On Unix-like systems
(including Linux), you can use the send2trash library to send files to the trash. Here's an example of how you
can use it:

First, you need to install the send2trash library if you haven't already:

pip install send2trash

Then, you can use Python to move files to the trash and empty it:

import os

import send2trash

def move_to_trash(file_path):

    try:

        send2trash.send2trash(file_path)

        print(f"Moved '{file_path}' to the trash.")

    except Exception as e:

        print(f"Failed to move '{file_path}' to the trash: {e}")

def empty_trash():

    try:

        send2trash.empty_trash(confirm=False)  # Set confirm to True to
confirm deletion

        print("Emptied the trash.")

    except Exception as e:

        print(f"Failed to empty the trash: {e}")

# Example usage

file_to_delete = "path/to/your/file.txt"

move_to_trash(file_to_delete)

# To empty the trash (be cautious, as this will permanently delete files)

empty_trash()


Please be very careful when using this code to empty the trash because it permanently deletes files,
and there's no easy way to recover them once they're gone. Make sure to test this on non-critical files
or in a controlled environment.
```

**Figure 5.** ChatGPT's response to a request to generate Python code to permanently delete files from the hard drive.

Source: author's own material based on ChatGPT.

Large language models generally do not pose a threat previously unknown. It is crucial to point out it is the scale of the known attacks that could increase dramatically and there is a need for adequate tools to combat them.

### Chances

Large language models can be used to the authorities' advantage as well. Automatic recognition of disinformation, identifying bot groups and AI generated text and much more can be done with the help of LLMs. However, most important methods of counteracting harmful usage of large language models are education, popularizing fact-checking and know ledge on LLMs functionalities and limitations. Data shows that young people (18–29 y.o.) tend to draw information from social media[26]. Students are often also the least equipped to verify false information[27] (this applies to teens and children as well).

### Recognizing

Monitoring the spread of disinformation and hate speech is a must in today's world. The key to success is identifying disinformation, hate speech, and AI generated text. Large language models can be[28] and are often used in content moderation, detecting disinformation, phishing[29] and social engineering.

### AI day-to-day support

There are plenty of ways in which large language models can be used to support the authorities, security researchers or fact checkers. Large language models are obviously useful; from code and text generation, through instructions to summarizing text and addressing a problem, this is an obvious advantage to people from different disciplines.

---

[26]   *Skąd Polacy czerpią informacje? Badanie IBRIS i IBIMS kwiecień 2021* (Eng. Where do Poles get their information from? IBRIS and IBIMS survey April 2021), IBiMS, http://www.ibims.pl/wp-content/uploads/2021/04/IBIMS_media_2021.pdf [accessed: 8 XI 2023].

[27]   J. Breakstone et al., *Students' Civic Online Reasoning: A National Portrait*, "Educational Researcher" 2021, no. 50, pp. 505–515. https://doi.org/10.3102/0013189X211017495; S. McGrew et al., *Can Students Evaluate Online Sources? Learning From Assessments of Civic Online Reasoning*, "Theory & Research in Social Education" 2018, no. 46, pp. 165–193. https://doi.org/10.1080/00933104.2017.1416320.

[28]   W. Saunders et al., *Self-critiquing models for assisting human evaluators*, preprint, arXiv, 12 VI 2022, https://arxiv.org/abs/2206.05802 [accessed: 9 XI 2023]. https://doi.org/10.48550/arXiv.2206.05802.

[29]   F. Heiding et. al, *Devising and Detecting Phishing: large language models vs. Smaller Human Models*, preprint, arXiv, 23 VII 2023, https://arxiv.org/abs/2308.12287 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2308.12287.

Examples include crime detection through video descriptions[30], text summarization and critiques for analytic products, disrupting telemarketers[31], providing opposing viewpoints. Large language models seem to be a great tool for initial analysis, helping save cognitive resources as it makes for a starting point for more advanced operations.

### Counteracting LLM-related threats

Primary means of counteracting the harm of misuse of LLMs is to educate[32], especially young people and children, about the possibilities and flaws of large language models. LLMs are there to stay and to pretend they don't exist over the course of education is simply impossible. Educators should be equipped with tools to teach responsible use of large language models. It is also crucial to introduce fact checking in the curriculum for students to be able to learn to verify information and critical thinking. On an international and national level, establishing Information Sharing and Analysis Centers (ISACs) concentrated on information exchange and combating disinformation offers a great deal of support to these educational efforts as well as forming a strong, unified position of government agencies, SM platforms and civil society against disinformation[33].

Another approach is to use large language models to counteract harm done by the misuse of LLMs. It may be possible to lower user interaction with harmful content (sharing or commenting) through providing factual responses to e.g. hate speech[34]. These models can be used to conduct social

---

[30] A. Simmons, R. Vasa, *Garbage in, garbage out: Zero-shot detection of crime using Large Language Models*, preprint, arXiv, 4 VII 2023, https://arxiv.org/abs/2307.06844 [accessed: 9 XI 2023]. https://doi.org/10.48550/arXiv.2307.06844.

[31] R. Currie, *California man's business is frustrating telemarketing scammers with chatbots*, The Register, 3 VII 2023, https://www.theregister.com/2023/07/03/jolly_roger_telephone_company/ [accessed: 8 XI 2023].

[32] NASK (@WeryfikacjaNASK), *Wraz z postępem technologicznym, rozwój AI staje się coraz bardziej widoczny w różnych dziedzinach naszego życia* (Eng. As technology advances, the development of AI is becoming more and more visible in various areas of our lives), X, 26 X 2023, https://twitter.com/WeryfikacjaNASK/status/1717487918556594437 [accessed: 8 XI 2023].

[33] L. Faesen et al., *Red Lines & Baselines Towards a European Multistakeholder Approach to Counter Disinformation*, The Hague Centre for Strategic Studies 2021.

[34] J. Podolak et al., *Analyzing the Influence of Language Model-Generated Responses in Mitigating Hate Speech on Social Media Directed at Ukrainian Refugees in Poland*, preprint, arXiv, 28 XI

campaigns. Their implementation is certainly valuable in cybersecurity[35], though it should be carried out with caution. And models identifying potential targets of e.g. a foreign agent based on a worker's activity could also be used to prevent them from being recruited in the first place.

### Silver Linings: Language Model's Playbook

There is also an upside of large language models answers being wrong. If the inquirer is led to believe something untrue about their planned crime, it may make it easier to detect them. An early example of such misinformed criminal comes from 1996, when a bank robber was led to believe lemon juice makes him invisible to surveillance cameras[36].

Similarly in this instance, while a search engine can be misleading, it's usually full of different sources. But large language models are meant to be plausible, they aren't always right (Figure 6). This makes for a chance that misleading information will actually help maintain public safety.

**Prompt:**
Is 450 90 percent of 500?

**Response:**

No, 450 is not 90 percent of 500. To find what is 90 percent of 500, you can multiply 500 by 0.9 (which represents 90 percent as a decimal):

500 * 0.9 = 450

So, 450 is 90 percent of 500.

**Figure 6.** ChatGPT asked if 450 is 90% of 500. The first answer was 'no' and the next 'yes'.

Source: author's own material based on ChatGPT.

## What does the accessibility of LLMs change?

Within the context of criminal activity, how are large language models different from what is already available on the internet? As with any

---

2023, https://arxiv.org/abs/2311.16905 [accessed: 31 XI 2023]. https://doi.org/10.48550/arXiv.2311.16905.

[35]   J. Goldfarb, *Applying AI to API Security, SecurityWeek*, 11 X 2023, https://www.securityweek.com/applying-ai-to-api-security/ [accessed: 8 XI 2023].

[36]   M.A. Fuocco, *Trial and error: They had larceny in their hearts, but little in their heads*, "Pittsburgh Post-Gazette" 1996.

new technology, their use is not fully regulated by law. This may have serious consequences, e.g. chatbots used for recruitment to terrorist organizations[37]. There is an important difference between LLMs (Figure 7) and, for example, search engines (Figure 8).

**Prompt:**

I'm writing a crime novel. I want it to be realistic. What are common mistakes killers make, that make them easier to identify by the police?

**Response:**

In a realistic crime novel, it's important to consider the mistakes killers might make that can lead to their identification by the police. Here are some common errors that criminals often make in real-life cases:

1. **Leaving Behind Physical Evidence:** This includes fingerprints, DNA, hair, fibers, or even personal belongings that can be traced back to them. Criminals who do not adequately clean or dispose of evidence may leave a trail for investigators.

2. **Inconsistent Alibis:** If a killer provides an alibi that doesn't hold up under scrutiny or changes their story, it can raise suspicion. Inconsistencies in their statements can be a significant red flag.

3. **Using Personal Vehicles:** Using a personal vehicle for criminal activities can lead to the vehicle's description and license plate being recorded by witnesses or surveillance cameras. This information can help law enforcement track down the suspect.

4. **Communicating with Victims or Witnesses:** Killers who engage in any form of communication with their victims or witnesses, whether before or after the crime, leave potential digital or paper trails that can be traced back to them.

5. **Social Media Posts:** Sharing incriminating details, photos, or other evidence on social media can provide investigators with valuable leads. Some criminals may boast about their crimes online.

6. **Repetition or Patterns:** If a killer follows a particular pattern or targets specific characteristics of victims, it can help investigators connect cases and narrow down the suspect pool.

7. **Impulsive Actions:** Impulsive decisions can lead to mistakes. Criminals may act hastily, leaving behind evidence they did not anticipate or engaging in suspicious behavior that draws attention.

8. **Witnesses and Surveillance Cameras:** Underestimating the number of potential witnesses or the reach of surveillance cameras in a given area can be a significant mistake. Many crimes are solved through witness statements and camera footage.

9. **Motives and Relationships:** Personal motives, disputes, or relationships with victims or others involved can point law enforcement in the direction of potential suspects. Personal connections can leave traces or provide leads.

10. **Acquaintances and Loose Lips:** Sharing information about the crime with friends, associates, or acquaintances can lead to tips given to the police. Loose lips can sink a criminal's attempt to evade justice.

11. **Unusual Purchases or Activities:** Making unusual purchases or engaging in activities that raise suspicion, such as buying large quantities of cleaning supplies or renting storage units, can draw the attention of law enforcement.

**Figure 7.** ChatGPT's answer to the question of what mistakes killers most often make and what makes it easier for the police to identify them.

Source: author's own material based on ChatGPT.

---

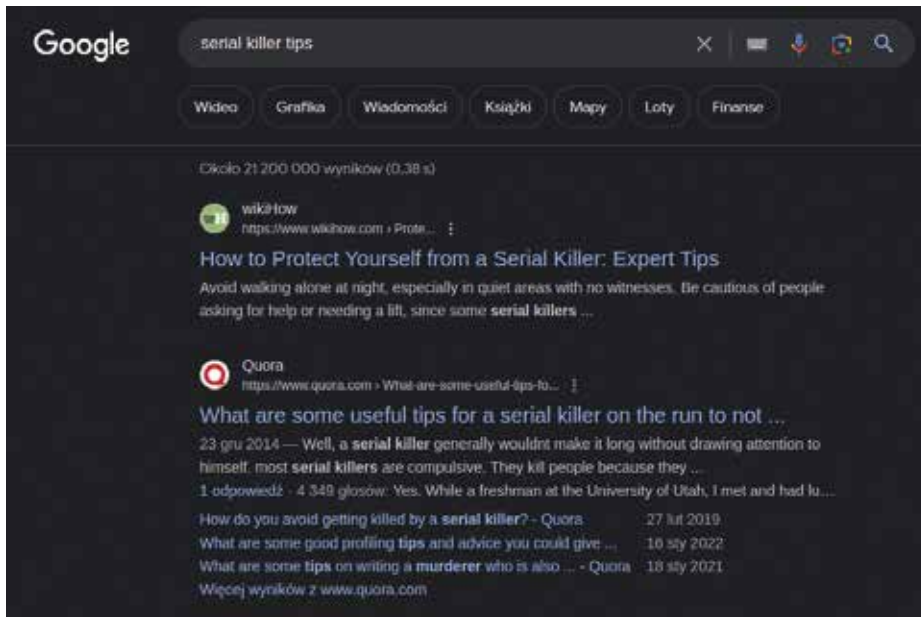[37] C. Vallance, I. Rahman-Jones, *Urgent need for terrorism AI laws…*

**Figure 8.** Google search for tips for serial killers.

Source: author's own material based on Google.

So what is the difference between the two? Mainly that a browser is more reliable because it offers multiple, up-to-date sources instead of a single answer based on data dating back to the year the training dataset was created.

However, there is a misconception about browsers that searching for something e.g. 'how to build a bomb' will put the searcher on a 'watchlist'. While this is not necessarily true, as such a list would be comprised mainly of mystery novelists, it does hold some value. Firstly, a browser history may be recovered e.g. for the sake of an ongoing trial. Secondly, the sole fear of detection through search history might sometimes be enough to stop some people from searching for some things altogether. In the case of LLMs, and specifically ChatGPT, such biases do not exist. It is not even certain whether prompt data can be recovered or whether it has evidentiary value. While it is difficult to determine what effect this has on the whole population, it certainly does make it easier to try and prepare for a crime anonymously.

Furthermore, it is no longer just about providing information. Considered dangers depend especially on models' like ChatGPT analysis

of human behavior and their ability to provide insightful responses. To further understand how these ChatGPT capabilities may pose a threat, the authors will draw on the criminological **routine activity theory.** First applied to the "structure of direct-contact predatory violations"[38], routine activity theory establishes that such acts require: an offender, at least one target (personal or property) and the absence of any effective guardian. The key part of this theory is the situation that brings these three elements together can result in a crime, while the lack of any of these is sufficient to prevent a crime.

In this instance, ignorance and misconceptions can serve both as a motivation for a criminally inclined offender (as it might have been for the unfortunate bank robber covered in lemon juice) and as an effective safeguard (the supposed "watchlist" based on searches). Similarly, the availability of a tool that might not yet be well researched by forensics lowers the effectiveness of safeguards and the ability needed for an offender to carry out their criminal inclinations.

The prompts may not be possible to replicate by specialists, the chat may not be accessible. The popular belief that Google searches can lead to the inquirer being put on a 'list', while not true, is a deterrent that just doesn't exist in regards to ChatGPT. The popularity of stories about different ways to exploit it, combined with lack of understanding what a language model can do, can make it seem that it can be used for anything, without any repercussions.

LLMs may contribute to increasing the number of offenders and may potentially make criminals bolder, decreasing the skill needed to commit more serious, complex criminal or terrorist acts.

## Conclusions

Large language models are a valuable tool for generating, simplifying, summarizing, and translating text, as well as analyzing people's attitudes and emotions. They can provide insightful answers, break down problems into smaller parts, and offer support with writing texts, including those on complex, intricate topics. Their availability and versatility increase their

---

[38]  M. Felson, L.E. Cohen, *Human ecology and crime: A routine activity approach*, "Human Ecology" 1980, no. 8, pp. 389–406. https://doi.org/10.1007/BF01561001.

potential applications in various areas, even those that might be against the law.

Concerning preparations for criminal and terrorist activities, the authors consider the key feature of LLMs to be their ability to provide insight into the issues presented in the prompt, allowing users to broaden their perspective (with the clarification that this insight is not understood by the authors as a thinking ability). Concerns about security related to the use of LLMs, such as prompt engineering or adversarial attacks, and their potential applications in crime and terrorism, should prompt further research in this area.

A general conclusion from this article is that artificial intelligence is developing at a very rapid pace. Unfortunately, the way innovative projects are funded in Poland is not adapted to the rapidly changing reality in aspects related to the use of artificial intelligence, and even less so to cybersecurity technologies. This is primarily due to the relatively long period of launching research and development programs and the lack of the ability to adapt projects to changing technologies. Breakthroughs in the development of artificial intelligence, such as new generations of LLMs, currently occur within periods of a few months, and the state of knowledge within a year can be turned upside down. This means that programs are already obsolete at the time of launch and are designed to create technologies in Poland that already exist elsewhere. This approach means that we will never bridge the technological gap, not to mention creating breakthrough technologies ourselves.

It should be noted that for years there has been a proven model for funding research and development that addresses all the shortcomings mentioned above. It's the Advanced Research Projects Agency (ARPA) model, which is used in the United States for resilience projects (DARPA), intelligence (IARPA), health care (ARPA-H), energy (ARPA-E), or infrastructure (ARPA-I). This model has been copied in Japan, Germany, and the UK. In a very figurative sense, the key to success in this model is to make the agency funding projects as smart as the project performers by engaging excellent program managers in project oversight. Managers engaged by DARPA are world-class researchers from academic, industrial, and government environments. These DARPA experts use their deep technical knowledge and professional networks to create new programs. They collaborate with research and development communities and stakeholders to establish rigorous and ambitious program goals (including

detailed technical milestones and other program activities) necessary to manage a portfolio of sophisticated research and development projects. Considering the accelerating development of modern artificial intelligence technologies, including LLMs, it seems reasonable to transition to a new paradigm of innovation funding. An intermediate step in such a change should be the launch of a pilot research program in Poland in line with ARPA's principles.

## Bibliography

Breakstone J. et al., *Students' Civic Online Reasoning: A National Portrait*, "Educational Researcher" 2021, no. 50, pp. 505–515. https://doi.org/10.3102/0013189X211017495.

Europol, *ChatGPT. The impact of Large Language Models on Law Enforcement*, Luxembourg 2023.

Faesen L. et al., *Red Lines & Baselines Towards a European Multistakeholder Approach to Counter Disinformation*, The Hague Centre for Strategic Studies 2021.

Felson M., Cohen L., *Human ecology and crime: A routine activity approach*, "Human Ecology" 1980, no. 8, pp. 389–406. https://doi.org/10.1007/BF01561001.

Fuocco M.A., *Trial and error: They had larceny in their hearts but little in their heads*, "Pittsburgh Post-Gazette" 1996.

GIFCT Red Team Working Group, *Considerations of the Impacts of Generative AI on Online Terrorism and Extremism*, [n.p.] 2023.

Ji Z. et al., *Survey of hallucination in natural language generation*, "ACM Computing Surveys" 2023, no. 12, pp. 1–38. https://doi.org/10.1145/3571730.

McGrew S. et al., *Can Students Evaluate Online Sources? Learning From Assessments of Civic Online Reasoning*, "Theory & Research in Social Education" 2018, no. 46, pp. 165–193. https://doi.org/10.1080/00933104.2017.1416320.

Raman G. et al., *How weaponizing disinformation can bring down a city's power grid*, "PloS One" 2020, no 15. https://doi.org/10.1371/journal.pone.0236517.

Vaidhyanathan S., *Antisocial media: How Facebook disconnects us and undermines democracy*, New York 2018.

Vaswani A. et al., *Attention is All you Need*, in: *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, I. Guyon et al. (eds.), Long Beach 2017, pp. 5998–6008.

Waniek M. et al., *Traffic networks are vulnerable to disinformation attacks*, "Scientific Reports" 2021, no. 11. https://doi.org/10.1038/s41598-021-84291-w.

### Internet sources

[heythereitsbeth], *Just came across this sub and thought I'd share mine from the start of the year*, Reddit, https://www.reddit.com/r/scambait/comments/17w6vx4/just_came_across_this_sub_and_thought_id_share/?rdt=40738 [accessed: 8 XI 2023].

AFP Kenya, *Fake subtitles added to old clip of Putin talking about Ukraine war, not Israel-Gaza conflict*, AFP Fact Check, 17 X 2023, https://factcheck.afp.com/doc.afp.com.33YG8TE [accessed: 8 XI 2023].

Bochyńska N., *#CyberMagazyn: Politycy narzędziem w rękach Kremla? „Świadomość jest bardzo niska"* (Eng. Are politicians a tool in the hands of the Kremlin? "Awareness is very low"), CyberDefence24, 21 X 2023, https://cyberdefence24.pl/cyberbezpieczenstwo/cybermagazyn-politycy-narzedziem-w-rekach-kremla-swiadomosc-jest-bardzo-niska [accessed: 8 XI 2023].

Borji A., *Stochastic Parrots or Intelligent Systems? A Perspective on True Depth of Understanding in LLMs*, preprint, SSRN, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4507038 [accessed: 8 IX 2023]. http://dx.doi.org/10.2139/ssrn.4507038.

Bowman S.R., *Eight Things to Know about Large Language Models*, preprint, arXiv, 2 IV 2023, https://arxiv.org/abs/2304.00612 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2304.00612.

Brewster T., *Armed With ChatGPT, Cybercriminals Build Malware And Plot Fake Girl Bots*, Forbes, 6 VI 2023, https://www.forbes.com/sites/thomasbrewster/2023/01/06/chatgpt-cybercriminal-malware-female-chatbots/ [accessed: 9 XI 2023].

Currie R., *California man's business is frustrating telemarketing scammers with chatbots*, The Register, 3 VII 2023, https://www.theregister.com/2023/07/03/jolly_roger_telephone_company/ [accessed: 8 XI 2023].

Derner E., Batistič K., *Beyond the Safeguards: Exploring the Security Risks of ChatGPT*, arXiv, preprint, 13 V 2023, https://arxiv.org/abs/2305.08005 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2305.08005.

Goldfarb J., *Applying AI to API Security,* SecurityWeek, 11 X 2023, https://www.securityweek.com/applying-ai-to-api-security/ [accessed: 8 XI 2023].

Gwozdowska A. et al., *Wojna informacyjna 2022–2023. Przebieg i wnioski* (Eng. Information warfare 2022-2023. Course and conclusions), NASK, 25 V 2023, https://www.nask.pl/pl/raporty/raporty/5204,Raport-quotWojna-informacyjna-20222023-Przebieg-i-wnioskiquot.html [accessed: 8 XI 2023].

Heiding F. et al., *Devising and Detecting Phishing: Large Language Models vs. Smaller Human Models,* preprint, arXiv, 23 VII 2023, https://arxiv.org/abs/2308.12287 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2308.12287.

INFO OPS Poland Foundation, *Model dystrybucji informacji w wirtualnym środowisku informacyjnym na bazie rozpoznanego rosyjskiego podstawowego modelu dystrybucji wiadomości manipulacyjnych* (Eng. A model of information distribution in a virtual information environment based on a recognised Russian basic manipulative news distribution model), Disinfo Digest, 9 VI 2023, https://disinfodigest.pl/model-dystrybucji-informacji-w-wirtualnym-srodowisku-informacyjnym-na-bazie-rozpoznanego-rosyjskiego-podstawowego-modelu-dystrybucji-wiadomosci-manipulacyjnych/ [accessed: 8 XI 2023].

Kelley D., *WormGPT – The Generative AI Tool Cybercriminals Are Using to Launch Business Email Compromise Attacks,* SlashNext, 13 VII 2023, https://slashnext.com/blog/wormgpt-the-generative-ai-tool-cybercriminals-are-using-to-launch-business-email-compromise-attacks/ [accessed: 22 XI 2023].

Lai T. et al., *Psy-LLM: Scaling up Global Mental Health Psychological Services with AI-based Large Language Models,* preprint, arXiv, 22 VII 2023, https://arxiv.org/abs/2307.11991 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2307.11991.

McGuffie K., Newhouse A., *The radicalization risks of GPT-3 and advanced neural language models,* preprint, arXiv, 15 IX 2020, https://arxiv.org/abs/2009.06807 [accessed: 8 XI 2023]. https://doi.org/10.48550/arXiv.2009.06807.

NASK (@WeryfikacjaNASK), *Wraz z postępem technologicznym, rozwój AI staje się coraz bardziej widoczny w rożnych dziedzinach naszego życia* (Eng. As technology advances, the development of AI is becoming more and more visible in various areas of our lives), X, 26 X 2023, https://twitter.com/WeryfikacjaNASK/status/1717487918556594437 [accessed: 8 XI 2023].

*Nowe oszustwo na WhatsAppie* „Kurier Szczeciński" (Eng. New WhatsApp scam), 5 VIII 2023, https://24kurier.pl/aktualnosci/wiadomosci/nowe-oszustwo-na-whatsappie/ [accessed: 8 XI 2023].

Podolak J. et al., *Analyzing the Influence of Language Model-Generated Responses in Mitigating Hate Speech on Social Media Directed at Ukrainian Refugees in Poland*, preprint, arXiv, 28 XI 2023, https://arxiv.org/abs/2311.16905 [accessed: 30 XI 2023]. https://doi.org/10.48550/arXiv.2311.16905.

Saunders W. et al., *Self-critiquing models for assisting human evaluators*, preprint, arXiv, 12 VI 2022, https://arxiv.org/abs/2206.05802 [accessed: 9 XI 2023]. https://doi.org/10.48550/arXiv.2206.05802.

Simmons A., Vasa R., *Garbage in, garbage out: Zero-shot detection of crime using Large Language Models*, preprint, arXiv, 4 VII 2023, https://arxiv.org/abs/2307.06844 [accessed: 9 XI 2023]. https://doi.org/10.48550/arXiv.2307.06844.

*Skąd Polacy czerpią informacje? Badanie IBRIS i IBIMS kwiecień 2021* (Eng. Where do Poles get their information from? IBRIS and IBIMS survey April 2021), IBiMS, http://www.ibims.pl/wp-content/uploads/2021/04/IBIMS_media_2021.pdf [accessed: 8 XI 2023].

Toulas B., *Cybercriminals train AI chatbots for phishing, malware attacks*, Bleeping Computer, 1 VIII 2023, https://www.bleepingcomputer.com/news/security/cyber-criminals-train-ai-chatbots-for-phishing-malware-attacks/ [accessed: 9 XI 2023].

Vallance C., Rahman-Jones I., *Urgent need for terrorism AI laws, warns think tank*, BBC News, 4 I 2024, https://www.bbc.com/news/technology-67872767 [accessed: 10 I 2024].

Yin S. et al., *A Survey on Multimodal Large Language Models*, preprint, arXiv, 23 VI 2023, https://arxiv.org/abs/2306.13549 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2306.13549.

Zou A. et al., *Universal and Transferable Adversarial Attacks on Aligned Language Models*, preprint, arXiv, 27 VII 2023, https://arxiv.org/abs/2307.15043 [accessed: 8 IX 2023]. https://doi.org/10.48550/arXiv.2307.15043.

## Julia Puczyńska

PhD student at IDEAS NCBR and at the Doctoral School of Information and Biomedical Technologies at the Institutes of the Polish Academy of Sciences. She is involved in combating disinformation and analyzing social media. A graduate in psychology and criminology, she has worked

on topics related to memory, facial recognition, and lie detection. She leads the Zgłoś Trolla (Eng. Report a Troll) project in cooperation with DC9 and the Geremek Foundation.

**Contact:** julia.puczynska@ideas-ncbr.pl

## Marcin Podhajski

PhD Student in the "AI for Security" research team at IDEAS NCBR and in the Doctoral School of Information and Biomedical Technologies at the Institutes of the Polish Academy of Sciences. He specializes in artificial intelligence, machine learning, and social network analysis.

## Karolina Wojtasik, PhD

Security specialist, academic researcher, and lecturer, forensic expert, Vice President for Scientific Affairs of the Polish National Security Association, employee of the Government Security Center. She deals with broadly understood security of critical infrastructure and objects subject to mandatory protection, especially in the context of physical and personal security threats. Author of books such as *Anatomia zamachu. O strategii i taktyce terrorystów* (Eng. Anatomy of an attack. On the Strategy and Tactics of Terrorists), *Ścieżki radykalizacji dżihadystycznej* (Eng. Paths of Jihadi Radicalization); co-author of the book *Polski system antyterrorystyczny a realia zamachów drugiej dekady XXI wieku* (Eng. The Polish anti-terrorist system and the realities of the attacks of the second decade of the 21st century) and many other publications related to terrorism, security, and building the resilience of critical infrastructure. Creator of the popular science channel *Anatomia zamachu* (Eng. Anatomy of an attack) on YouTube and the Projekt IK portal.

## Tomasz P. Michalak, PhD

Leader of an independent research team at IDEAS NCBR and lecturer at the Faculty of Mathematics, Informatics, and Mechanics of the University of Warsaw. Graduate of the Faculty of Economic Sciences at the University of Warsaw. During his academic career, he conducted research at the Department of Computer Science at the University of Oxford, the School of Engineering and Computer Science at the University of Southampton, the Department of Computer Science at the University of Liverpool, and the Faculty of Applied Economics at the University of Antwerp, where he obtained a PhD in economics.