



UNIDIR

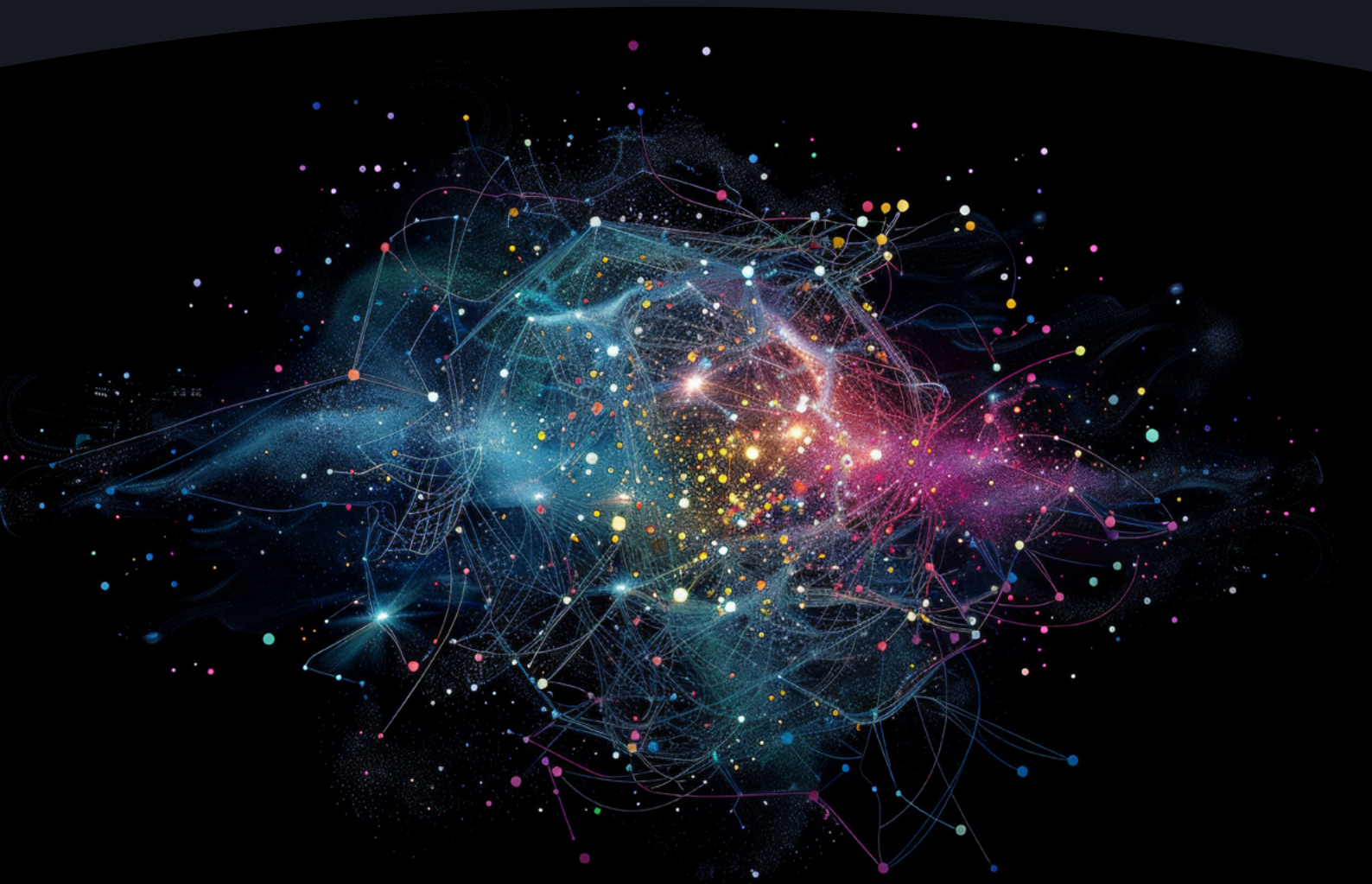


Funded by
the European Union

Governance Implications of Synthetic Data in the Context of International Security

A Technology and Security Seminar Report

FEDERICO MANTELLASSI



Acknowledgements

Support from UNIDIR's core funders provides the foundation for all of the Institute's activities. This publication was funded by the European Union as part of UNIDIR's Security and Technology Programme, which is also supported by the Governments of Czechia, France, Germany, Italy, the Netherlands, Norway and Switzerland, and by Microsoft.

The author would like to extend his sincere thanks to Wenting He for her moderation and organization of the event's first panel, as well as Jessica Espinosa Azcarraga for help in organizing the event. Additionally, the author would like to extend his thanks to all the speakers for their participation, as well as Dr. Giacomo Persi Paoli, Sarah Grand Clément, Wenting He, Calum Inverarity, Dr. Ana Beduschi and Aldo Lamberti for their comments on this report.

About UNIDIR

The United Nations Institute for Disarmament Research (UNIDIR) is a voluntarily funded, autonomous institute within the United Nations. Being one of the few policy institutes worldwide focusing on disarmament, UNIDIR generates knowledge and promotes dialogue and action on disarmament and security. It is based in Geneva and assists the international community in developing the practical, innovative ideas needed to find solutions to critical security problems.

About the Security and Technology Programme

Contemporary developments in science and technology present new opportunities as well as challenges to international security and disarmament. The UNIDIR's Security and Technology Programme aims to build knowledge and awareness about the international security implications and risks of specific technological innovations and convenes stakeholders to explore ideas and develop new thinking on ways to address them.

Note

The designations employed and the presentation of material in this publication do not imply the expression of any opinion whatsoever on the part of the Secretariat of the United Nations concerning the legal status of any country, territory, city or area, or of its authorities, or concerning the delimitation of its frontiers or boundaries. The views expressed in the publication are the sole responsibility of the individual authors. They do not necessarily represent the views or opinions of the United Nations, UNIDIR, its staff members, or sponsors.

Citation

Mantellassi, Federico. "Governance Implications of Synthetic Data in the Context of International Security: A Technology and Security Seminar Report". Geneva, Switzerland: UNIDIR, 2024.

About the Author



Federico Mantellassi

Researcher, Security and Technology Programme

Federico Mantellassi is a Researcher in the Security and Technology Programme at UNIDIR. His work focuses on the international security implications, risks and opportunities of emerging science, and technology developments and innovations. Previously, Federico was a Research and Project Officer at the Geneva Centre for Security Policy, conducting research on the intersection between emerging technologies, international security and warfare. He holds a master's degree in Intelligence and International Security from King's College London and a bachelor's degree in International Studies from the University of Leiden.

Acronyms & Abbreviations

AI	Artificial intelligence
EU	European Union
GDPR	General Data Protection Regulation
IEEE	Institute of Electrical and Electronics Engineers
ISO	International Organization for Standardization
ODI	Open Data Institute
PET	Privacy Enhancing Technology

Table of Contents

- 1. Introduction 6**
 - 1.1. About the Event 7

- 2. Synthetic Data and International Security: Framing the Issue 8**
 - 2.1. What is Synthetic Data 8
 - 2.2. Synthetic Data in the Military Domain 9

- 3. Governance Challenges and Implications 11**
 - 3.1. Synthetic Data and Civilian Data Governance 11
 - 3.2. The Role of Standards 13
 - 3.3. Synthetic Data and the International Governance of Military AI 14
 - 3.3.1. On the novelty of challenges and the applicability of existing frameworks 14
 - 3.3.2. On the importance of a multistakeholder approach 15
 - 3.3.3. On guidelines and context specificity 16
 - 3.3.4. On governance opportunities in the military domain 16

- 4. Conclusion 18**

- Annex: Event Agenda and Participants 19**

1. Introduction

Data is crucial to the training and development of artificial intelligence (AI) systems. However, three key data-related issues can act as barriers to development and deployment of AI capabilities and systems. First, the development of AI technologies has – at least in part – depended on the availability of large datasets to train AI models. Second, data is a resource whose availability, collection, cleaning, use and sharing is affected by factors such as collection costs, lack of real-world data in certain domains, as well as regulatory, legal and ethical constraints. Third, data quality, representativeness, and diversity are directly linked to an AI model's performance, level of bias, accuracy, and reliability. **Synthetic data – data that is artificially generated in the digital world with properties that are often derived from an original set of data – has been proposed as a solution to address some of these data-related issues, especially for AI model training.**¹ Indeed, synthetic data can help to address issues such as biases in datasets while also enabling their expansion, creation, diversification, and fine-tuning. Synthetic data is also

often referred to as a privacy-enhancing technology (PET), facilitating the use and sharing of sensitive datasets.² Synthetic data is particularly promising for domains such as the military.³ In this sensitive domain, AI-enabled capabilities are in increasing demand, but high-quality, diverse datasets are in short supply and the consequences of faulty algorithms are potentially serious. Synthetic data could enable the ability to develop advanced AI capabilities with less need for troves of real-world data.⁴ However, synthetic data is not a panacea and has been shown to potentially exacerbate many of the issues it seeks to curtail, sparking governance and regulatory discussions.⁵

Synthetic data exists in a relative ‘grey zone’ in terms of regulation and governance. Major data governance and AI regulatory frameworks, such as the European Union's AI Act and the General Data Protection Regulation (GDPR), mention synthetic data only in passing, if at all. For some, this entails that synthetic data, as a PET, can be a way around stringent regulatory frameworks, or a useful compliance

¹ Hao, Shuang et al. 2024. “Synthetic Data in AI: Challenges, Applications and Ethical Implications.” *School of Software Engineering, Huazhong University of Science and Technology*. <https://arxiv.org/pdf/2401.01629>; Lee, Peter. 2024. “Synthetic Data and the Future of AI.” *Cornell Law Review*. Forthcoming. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4722162; Deng, Harry. 2023. “Exploring Synthetic Data for Artificial Intelligence and Autonomous Systems: A Primer.” *United Nations Institute for Disarmament Research*. https://unidir.org/wp-content/uploads/2023/11/UNIDIR_Exploring_Synthetic_Data_for_Artificial_Intelligence_and_Autonomous_Systems_A_Primer.pdf.

² Naughton, Mitchell et al. 2023. “Synthetic Data as a Strategy to Resolve Data Privacy and Confidentiality Concerns in the Sport Sciences: Practical Examples and an R Shiny Application.” *International Journal of Sports Physiology and Performance*. Vol 18 (10): 1213-1218. doi: 10.1123/ijsp.2023-0007; *Syntheticus*. “Synthetic Data 101: What Is It, How It Works and What It's Used For.” *Syntheticus*. Web. n.d. <https://syntheticus.ai/guide-everything-you-need-to-know-about-synthetic-data#chapter-8>.

³ Chahal, Husanjot et al. 2020. “Messier than Oil: Assessing Data Advantage in Military AI.” *Center for Security and Emerging Technology*. <https://cset.georgetown.edu/wp-content/uploads/Messier-than-Oil-Brief-1.pdf>.

⁴ Ibid.

⁵ Deng, Harry. 2023. “Exploring Synthetic Data for Artificial Intelligence and Autonomous Systems: A Primer.” *United Nations Institute for Disarmament Research*. https://unidir.org/wp-content/uploads/2023/11/UNIDIR_Exploring_Synthetic_Data_for_Artificial_Intelligence_and_Autonomous_Systems_A_Primer.pdf.

tool.⁶ Others point to the fact that synthetic data carries with it many of the same risks as real-world data, and can result in similar downstream effects on AI model accuracy, safety, fairness, and representativeness, and thus they insist that new regulatory frameworks and approaches are necessary to avoid governance gaps and blind spots.⁷ In this light, it is of utmost importance to understand how current governance (civilian and military) and regulatory frameworks encompass synthetic data, whether they are fit for purpose to address potential risks, and if they need to be adjusted. Regulatory and governance gaps are of particular consequence in the context of the fast-advancing adoption of AI-enabled capabilities in the military domain. Understanding the implications of synthetic data for emerging military AI

governance discussions is therefore essential. To explore the governance challenges of synthetic data in the context of international security, UNIDIR's Security and Technology Programme held an event titled Technology and Security Seminar on Synthetic Data: Exploring Governance Implications.

This report provides a summary of the key themes and takeaways from discussions at the event. The report is divided into two parts, reflecting the structure of the event. The first part provides a short overview of the technology and its uses in the military domain. The second part presents the various views, issues, and potential challenges to governance presented by synthetic data in the context of international security.

1.1. About the Event

The Technology and Security Seminars comprise a series of events organized by UNIDIR's Security and Technology Programme focused on various enabling technologies. The key objectives of the series are threefold:

- ▶ expose the diplomatic community to a wide range of emerging, critical enabling technologies;
- ▶ alert the diplomatic community to the potential international security implications of the development and use of such technologies; and
- ▶ explore governance possibilities through multi-stakeholder dialogue and engagement.

On 29 October 2024, a Technology and Security Seminar was held on the topic of synthetic data governance. This half-day event consisted of a **Technology Breakfast, serving as an introduction to the technology for policymakers**, as well as a **Multi-Stakeholder Dialogue on Synthetic Data** where experts from industry, international organizations, and academia convened to share a variety of views on the specific governance challenges in the context of international security. The event took place virtually, on the margins of the seventy-ninth session of the United Nations General Assembly's First Committee in 2024. For a full programme of the event, please see the annex to this report.

⁶ Zojer, Alexander. "Synthetic Data: A Key Tool for AI Compliance under the EU's AI Act." *Mostly.AI*. 30 October 2023. <https://mostly.ai/blog/ai-compliance-with-eu-ai-act-using-synthetic-data>.

⁷ Gal, Michal, Lynskey, Orla. 2023. "Synthetic Data: Legal Implications of the Data Generation Revolution." *Iowa Law Review*. 109. Forthcoming. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4414385; Whitney, Cedric Deslandes, Norman, Justin. 2024. "Real Risks of Fake Data: Synthetic Data, Diversity-Washing and Consent Circumvention." *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*. <https://doi.org/10.1145/3630106.3659002>.

2. Synthetic Data and International Security: Framing the Issue

Takeaways

- ▶ Advances in the field of generative AI and growing AI adoption across sectors have expanded the pervasiveness of synthetic data, increasing the scale and ease with which it can be generated as well as its variety and quality.
- ▶ Synthetic data holds promise and offers potential solutions to data-related challenges (bias, scarcity, quality, representativeness, privacy) in both civilian and military domains.
- ▶ Armed forces are increasingly turning towards synthetic data in the context of their growing adoption of AI-enabled capabilities to train military AI models for identification, targeting systems, operational and tactical planning, as well as the development of scenarios and synthetic environments.
- ▶ Despite its benefits, synthetic data can perpetuate existing data-related risks, create new ones, or expand the magnitude of their impacts.

2.1. What is Synthetic Data⁸

Synthetic data can be defined as “**information created by computer simulations or algorithms that reproduce some structural and statistical properties of real-world data.**”⁹ Various generation methods exist for synthetic data, and the resulting datasets can either be fully synthetic (with all data artificially generated), partially synthetic (with a small portion of a real dataset replaced with synthetic data), or hybrid (where real-world

and fully synthetic data are blended).¹⁰ In short, synthetic data is mostly utilized to **complete datasets** (and seek to address issues in the data, such as those related to bias or representativeness), **create datasets where none exist**, or **remove personally identifiable information** when sensitivity requires it. Hence, the value of synthetic data lies in its ability to assist with **key data issues, namely bias, representativeness, quality, scarcity, and privacy.**

⁸ This section builds on previous work undertaken by UNIDIR’s Security and Technology Programme on Synthetic Data and International Security. For a detailed, in-depth exploration of what synthetic data is, and of the international security risks and opportunities linked to synthetic data, especially in the context of AI enabled and autonomous military capabilities, see https://unidir.org/wp-content/uploads/2023/11/UNIDIR_Exploring_Synthetic_Data_for_Artificial_Intelligence_and_Autonomous_Systems_A_Primer.pdf.

⁹ De Wilde, Philippe et al. 2024. “Recommendations on the Use of Synthetic Data to Train AI Models.” *United Nations University*. <https://collections.unu.edu/eserv/UNU:9480/Use-of-Synthetic-Data-to-Train-AI-Models.pdf>.

¹⁰ Syntheticus. “Synthetic Data 101: What Is It, How It Works and What It’s Used For.” *Syntheticus*. Web. n.d. <https://syntheticus.ai/guide-everything-you-need-to-know-about-synthetic-data>.

While synthetic data is not a novel concept, and has been used for some time, recent technological advances – especially in generative AI – have dramatically increased the **scale and ease** with which it can be produced, the **diversity** of types of data that can be created, as well as its **quality**. These advances have lowered the bar of access to synthetic data and vastly expanded the number of individuals and organizations without extensive technical expertise that can now utilize it. In turn, this has increased its pervasiveness, with some assessments estimating that 60% of all AI training data will be synthetic as of 2024.¹¹ The increasing popularity of synthetic data is furthermore the result of the ever-growing need for more data for the training of AI models.

Synthetic data is however no panacea, and has been shown to potentially perpetuate, and sometimes exacerbate, the problems its use aims to address. Indeed, synthetic data is not inherently private, secure, representative or unbiased, necessitating much consideration and curation to make it so. Furthermore, research has shown that the repetitive training of AI models on synthetic data generated from previous version of themselves can lead to ‘model collapse’, whereby a model forgets its underlying data distribution leading to drastic reduction in output quality and accuracy.¹² Additionally, the increased prevalence of synthetic data could expand the risk surface in data-related issues and increase the magnitude of negative impacts.¹³

2.2. Synthetic Data in the Military Domain

Synthetic data is increasingly prevalent in the military domain, where issues surrounding data scarcity, bias, and sensitivity are particularly acute.¹⁴ Like in the civilian sector, the increased use of synthetic data in this domain is linked to armed forces’ turn towards AI-enabled solutions. In this context, synthetic data is primarily used for the **training of military AI models for identification, targeting, operational and tactical planning**, as well as the development of scenarios and **synthetic environments**.

Principally, synthetic data can help armed forces to **fill gaps and increase the quality of their datasets** – such as creating images of

objects from different angles and in different conditions – to increase the performance of AI models. Additionally, synthetic data can assist in **data management**, helping to reduce costs associated with labelling and collection, and **accelerating the development** of AI products. Furthermore, synthetic data can be used to create realistic simulations of various military operations, including adversarial attacks. These simulations can enable States to test the effectiveness of their AI systems, develop new strategies and tactics, and prepare for a wider range of potential threats in a controlled and safe environment.

¹¹ Gartner. “Gartner Identifies Top Trends Shaping the Future of Data Science and Machine Learning.” 1 August 2023. <https://www.gartner.com/en/newsroom/press-releases/2023-08-01-gartner-identifies-top-trends-shaping-future-of-data-science-and-machine-learning>.

¹² Shumailov, Iliia, et al. 2024. “AI Models Collapse When Trained on Recursively Generated Data.” *Nature*. <https://www.nature.com/articles/s41586-024-07566-y>.

¹³ Deng, Harry. 2023. “Exploring Synthetic Data for Artificial Intelligence and Autonomous Systems: A Primer.” *United Nations Institute for Disarmament Research*. https://undir.org/wp-content/uploads/2023/11/UNIDIR_Exploring_Synthetic_Data_for_Artificial_Intelligence_and_Autonomous_Systems_A_Primer.pdf.

¹⁴ Ibid.

However, use of synthetic data in the military domain suffers from inherent risks linked with such data use. Indeed, synthetic data, despite aiming to represent reality, can perpetuate and even reinterpret existing biases found in the original data it is derived from. That possibility presents a significant risk, particularly in sensitive military contexts where biased decisions can have severe consequences. Moreover, risks of re-identification of individuals or sensitive information within datasets persist, potentially leading to the disclosing of sensitive military data, while ‘data poisoning’ attacks by malicious actors could skew the learning process of AI systems.¹⁵



AI generated, Adobe Stock.

¹⁵ Ibid.

3. Governance Challenges and Implications

Takeaways

- ▶ The governance landscape for synthetic data is immature in both the civilian and military domains. More work is needed to provide clarity over how existing governance frameworks and regulations apply to synthetic data, and how they might need to be adapted to better cover possible gaps.
- ▶ No consensus exists over the need for new and dedicated regulations and frameworks specifically focused on synthetic data.
- ▶ International standards are an important tool in the technology governance toolbox. While no international standards exist with respect to synthetic data, work is ongoing in their development and will be instrumental in fostering responsible innovation and adoption of the technology.
- ▶ Due to its increased use in the military domain to train AI systems, synthetic data is of high relevance to military AI governance discussions. More work should be undertaken to apply, adapt, or build upon established practices and governance concepts linked to data in the military domain.
- ▶ Synthetic data presents opportunities for the governance of military AI, by potentially enabling greater data-sharing, joint development of AI applications, and common development of guidelines for responsible synthetic data generation and use, hereby advancing global responsible military AI goals.
- ▶ Governance of synthetic data in the military domain will require more multi-stakeholder engagement. This entails cooperation among States, but also with the private sector, which should be closely involved in governance discussions and efforts. Fostering trust between governments and industry will be fundamental to this effort.

3.1. Synthetic Data and Civilian Data Governance

Balancing the risks and the opportunities of synthetic data will require an understanding of its governance challenges. While synthetic data is not necessarily novel, **governance discussions relating to its generation and use are only now emerging both in civilian and military domains.** Questions regarding

synthetic data's legal status, regulatory needs, and potential governance approaches are embryonic, and the governance landscape remains immature. No legislation or frameworks specific to synthetic data currently exist. Some regulatory frameworks such as the EU AI Act mention synthetic data in passing, while

select governments have issued guidelines on synthetic data generation.¹⁶

In the civilian domain, no consensus exists on whether synthetic data challenges data regulatory and governance frameworks, and if so, in which ways. For example, it is noted that synthetic data could challenge the categories of personal/non-personal data, which are the foundation of data governance regulations and frameworks such as the EU's GDPR. **It is argued that regulations such as these are not adequately equipped to address the complexities of synthetic data, which can blur the lines between these categories.** Depending on the type of synthetic data – fully synthetic, partially synthetic, or hybrid – the level of personal information present and the risk of re-identification, and therefore the applicability of data protection laws, can vary considerably.¹⁷ This ambiguity creates legal uncertainties for both developers and users of synthetic data. In this respect, **the increasing use of synthetic data may necessitate an expansion of the scope of the traditional personal/non-personal data paradigm in data protection regulation.**

For others, appropriately generated synthetic data is a useful PET, to be used as a tool for compliance with various data regulation

frameworks. Additionally, synthetic data could be of use in achieving broader data governance goals, by democratizing access to valuable data while **protecting privacy**, enabling transparent data catalogues and audit trails for **accountability**, **improving data quality** by providing a consistent and controllable data source, and **facilitating secure data-sharing** on national and international levels. Disagreements persist over whether synthetic data is a useful tool to be incentivized, or an innovation possibly undermining legal mechanisms developed to guard against various data-related risks.¹⁸

What transpires is therefore a lack of legal, and normative, clarity with respect to the processing of synthetic data. Hence, some have argued for the need for clear guidelines, to ensure **transparency**, **fairness**, and **accountability** in the processing of all types of synthetic data, as well as enhanced clarity and guidelines with regards to what data has been employed in the foundational models used to generate synthetic data.¹⁹ Propositions include:

- ▶ **transparency:** synthetic data should be clearly labelled as such, and information about its generation process should be available;
- ▶ **accountability:** means of establishing clear procedures for calling to account those

¹⁶ Personal Data Protection Commission of Singapore. 2024. "Privacy Enhancing Technology Proposed Guidance on Synthetic Data Generation." *Personal Data Protection Commission of Singapore*. <https://www.pdpc.gov.sg/-/media/files/pdpc/pdf-files/other-guides/proposed-guide-on-synthetic-data-generation.pdf>.

¹⁷ Fully synthetic datasets contain data that is fully generated by an AI model and contains no real-world data. The model identifies the statistical proprieties and patterns of a dataset and generates an entirely new one. Partially synthetic data replaces some selected sensitive features of a dataset and replaces them with synthetic values, while keeping some real data. Hybrid synthetic data combines real world, and fully synthetic data, pairing random records from a real dataset with a synthetic record. For more detail, please see Synthetikus. "Synthetic Data 101: What Is It, How It Works and What It's Used For." *Synthetikus*. Web. n.d. <https://synthetikus.ai/guide-everything-you-need-to-know-about-synthetic-data>. and IBM. "What is synthetic data." IBM. n.d. <https://www.ibm.com/topics/synthetic-data>.

¹⁸ Gal, Michal, Lynskey, Orla. 2023. "Synthetic Data: Legal Implications of the Data Generation Revolution." *Iowa Law Review*. 109. Forthcoming. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4414385.

¹⁹ Beduschi, Ana. 2024. "Synthetic Data Protection: Towards a Paradigm Change in Data Regulation?" *Big Data and Society*. Vol 11 (1). <https://doi.org/10.1177/20539517241231277>.

responsible for generation and processing of synthetic data should be developed; and

- **fairness:** synthetic data should include some guarantees that it is not being generated and used in ways that bring adverse effects, such as perpetuating biases or creating new ones.

Due to a lack of legal clarity, it is possible that synthetic data falls outside of regulatory oversight while potentially carrying with it some of the same issues that these governance

frameworks seek to address in terms of real-world data. Hence, a **key governance challenge in the civilian domain will be ensuring that synthetic data is developed and used in a way that, if it falls outside the scope of current data regulations, does not perpetuate, or create new harms.** Further research and work are therefore required to more clearly spell out synthetic data's legal standing, as well as to identify potential governance gaps in existing regulations and frameworks to offer clarity to developers and users of synthetic data.

3.2. The Role of Standards

Standards (both technical and non-technical) are an important aspect of civilian technology governance. In the context of synthetic data, they are both much needed, and an important area of current work. **Indeed, no international standard for the generation of synthetic data exists,** with no universally agreed upon definitions or benchmarks for evaluating the quality and trustworthiness of synthetic data. **By providing clear definitions, methodologies, and evaluation criteria, standards can create common understandings and benchmarks for assessing synthetic data.** Standards can reassure organizations that the synthetic data they use meets specific quality and privacy thresholds. Moreover, standards for the labelling and documentation of synthetic data, as well as auditing and mechanisms to track provenance, could form a key building block for ensuring transparency,

fairness and accountability in the generation and processing of synthetic data.

Work to this effect is beginning on various fronts. The Institute of Electrical and Electronics Engineers (IEEE), for example, is leading efforts to develop a global standard and best practices for privacy-safe synthetic data.²⁰ Similar efforts are underway in the International Organization for Standardization (ISO).²¹ The Open Data Institute (ODI) has contributed towards the development of a tool named 'Croissant', a community standard that provides machine-readable metadata for datasets, helping to standardize documentation of machine learning datasets.²² Standards will provide a framework and parameters for responsible innovation, incentivizing good practices in the private sector for synthetic data generation, use and innovation.

²⁰ The Institute of Electrical and Electronics Engineers. "Synthetic Data." *The Institute of Electrical and Electronics Engineers*. n.d. <https://standards.ieee.org/industry-connections/activities/synthetic-data/>.

²¹ International Organization for Standardization. "Information Technology — Artificial intelligence — Overview of Synthetic Data in the Context of AI Systems." *International Organization for Standardization*. n.d. <https://www.iso.org/standard/86899.html#lifecycle>.

²² Simperl, Elena and Thomas Carey-Wilson. "The ODI to Help Develop an Open Metadata Standard for Machine Learning Data." Open Data Institute. 6 March 2024. <https://theodi.org/news-and-events/blog/the-odi-to-help-develop-an-open-metadata-standard-for-machine-learning-data/>.

3.3. Synthetic Data and the International Governance of Military AI

Data is of increasing military importance due to the growing centrality of AI in many aspects of the military domain. In light of this, and because synthetic data is primarily used in the military domain for the training and development of various AI capabilities, governance discussions linked to synthetic data's impact on international security should be discussed within the context of the governance of military AI.

Data, and its related issues, have been identified as one of the priority areas of work for responsible AI in the military domain.²³ **However, these issues have not been central to ongoing governance efforts**, remaining largely high-level and lacking granularity. It therefore follows that discussions **surrounding synthetic data, its governance, and importantly its potential effect on ongoing military AI governance efforts have themselves been embryonic.** Some regional-level discussions, although similarly nascent, have taken place, especially in still-digitalizing parts of the world where data gaps present a significant barrier of entry to military AI and where alternatives such as synthetic data can act as surrogate.

3.3.1. On the novelty of challenges and the applicability of existing frameworks

Synthetic data in the military domain brings about similar challenges, with similar implications, to real data, such as biases (both their perpetuation, exacerbation, or creation), reliability

and representation concerns, accountability, traceability or lack of explainability among others. Therefore, **synthetic data should not be excluded or remain an unexplored issue in the context of military AI.** Importantly, this could lead to a situation where certain types of data remain outside the scope of governance discussions, while potentially perpetuating risks, hence exacerbating data-related risks in military AI.

For example, ensuring data *accountability* – a key tenet of responsibility – could be further complicated in the context of synthetic data. Indeed, the use of synthetic data introduces an additional layer of persons, sometimes external actors, responsible for its generation, thereby making it harder to trace direct accountability in case of errors. Relatedly, synthetic data could exacerbate data *explainability* issues, **due to the lack of internationally agreed upon standards on the generation, use, and labelling of synthetic datasets.** Lack of clarity over data provenance due to limited traceability could then hinder auditing capabilities to address biases in datasets. Moreover, the democratization of data access through synthetic data could be an opportunity for greater access to the development of AI and other digital capabilities. On the one hand, this could help to address issues surrounding the digital divide. On the other hand, in the context of international security, it could also act as an enabler of **greater proliferation of military AI capabilities** by lowering the bar of entry to the development

²³ Afina, Yasmin, Persi Paoli, Giacomo. 2024. "Governance of Artificial Intelligence in the Military Domain: A Multi-Stakeholder Perspective on Priority Areas." *United Nations Institute for Disarmament Research*. https://undir.org/wp-content/uploads/2024/09/UNIDIR_Governance_of_Artificial_Intelligence_in_the_Military_Domain_A_Multi-stakeholder_Perspective_on_Priority_Areas.pdf.

of advanced AI models.²⁴ Yet, rather than representing novel challenges, these are further complications of pre-existing data challenges in the military domain, which require specific attention and clarification in the context of synthetic data.

Governance discussions surrounding novel technologies' impacts on international security should therefore ideally first focus on whether – and how – legal and normative frameworks apply. In the context of synthetic data, it is important to analyze whether these governance challenges are new, whether they only complicate pre-existing data challenges

and to what extent existing data governance concepts apply to synthetic data. Hence, as opposed to necessarily designing new governance frameworks or approaches, the international community should look to established practices and concepts, such as equitability, responsibility, traceability and reliability, and work on applying, adapting, or building upon them for synthetic data. In this matter, the international community is not starting from scratch and can leverage an already extensive body of work on what constitutes 'good data' and emerging knowledge on what good data practices in the military domain look like.

BOX 1.

Area of Future Research: International Trade of Synthetic Datasets

An important area needing further research is the potential implications of the international trade in synthetic datasets, and whether – and how – such trade should in some cases be controlled, monitored, or restricted. Indeed, a market could develop to trade synthetic datasets which could be used by malicious actors in the development of disruptive AI capabilities. The international community should hence consider how the trade in synthetic datasets interfaces with non-proliferation efforts and arms control. It should explore whether some synthetic datasets, or types of synthetic data, should be controlled through tools such as control lists for export controls.

3.3.2. On the importance of a multistakeholder approach

Governance frameworks which will most effectively enable the leveraging of the benefits of synthetic data for the military domain are ones which will be **multi-stakeholder in nature**. This entails not only cooperation among States, **but close cooperation with private sector actors, who should be involved in such governance**

discussions. The private sector plays a major role in military AI, being primarily responsible for the research and development of core technologies. This remains true with synthetic data, where most of the capabilities for generation as well as testing and evaluation of synthetic datasets rest with industry players. **This creates additional dependencies on private technology companies, especially for States**

²⁴ Maas, Matthijs M. 2019. "Innovation-Proof Global Governance for Military Artificial Intelligence? How I Learned to Stop Worrying and Love the Bot." *Journal of International Humanitarian Legal Studies*. Vol 10 (1). https://brill.com/view/journals/ihts/10/1/article-p129_129.xml?language=en.

with lesser resources, which may not have independent testing and evaluation capabilities for synthetic data, and which rely on the private sector for the quality of synthetic datasets. This dependency necessitates frameworks for public–private partnerships that prioritize **building trust between governments and industry**. Such trust is crucial to ensure that private sector actors engage in governance discussions and adopt responsible practices in developing, deploying, and testing synthetic data for military AI systems. Furthermore, multi-stakeholder approaches contribute to the much-needed creation of common language and understanding of synthetic data.

3.3.3. On guidelines and context specificity

While the development of clear guidelines in the generation and processing of synthetic data should be an aspiration, it has been noted that, especially in a military context, their development might be premature. Indeed, clear guidelines are typically grounded in well-defined best practices. However, in the case of synthetic data, the field might still be too nascent to establish definitive best practices. In this context, standardizing testing procedures or establishing rigid guidelines before a thorough understanding of the technology’s capabilities, benefits, limitations, and potential risks has been developed might be counterproductive.

Furthermore, clear guidelines might not accommodate for highly context-specific nature of assessing the appropriateness and level of responsibility of synthetic data-use in a military context. In fact, a given synthetic dataset might be used ‘responsibly’ in one

scenario, while its use in another context could be seen as ‘irresponsible’. This context-dependent nature of synthetic data-use makes it difficult to develop universally applicable guidelines that effectively address the nuances of different use cases.

Additionally, quality metrics are themselves also context dependent. For example, the closeness with which a synthetic dataset represents reality is used as a key indicator of its quality.²⁵ However, in some cases, particularly in the military domain, deviation from this is precisely the intent. In other words, **using synthetic data to represent unprecedented scenarios to help with creative planning might be one of the advantages unlocked by the use of synthetic data in the context of military operations.**

Moreover, governance approaches will require regional and national contexts to be taken into account. Due to the external dependencies synthetic data can create – such as the reliance on external actors for the generation of synthetic datasets and their quality assurance – it is particularly important to ensure that synthetic datasets created outside a given region reflect local realities in the intended context of use. This will require transparency over the parameters and assumptions of synthetic datasets, setting out why, how, what for, and by whom synthetic data is being created.

3.3.4. On governance opportunities in the military domain

Synthetic data does not only pose governance challenges, but presents opportunities as well, especially for the governance of military AI. There is indeed the potential

²⁵ Deng, Harry. 2023. “Exploring Synthetic Data for Artificial Intelligence and Autonomous Systems: A Primer.” *United Nations Institute for Disarmament Research*. https://unidir.org/wp-content/uploads/2023/11/UNIDIR_Exploring_Synthetic_Data_for_Artificial_Intelligence_and_Autonomous_Systems_A_Primer.pdf.

for synthetic data to facilitate **data-sharing abilities** between armed forces, as well as to help in the **common development of military AI capabilities**. For example, synthetic data's privacy-preserving potential could provide opportunities to share datasets, something that is often desired, but impeded by the sensitivity and classified nature of military data. In the context of the military domain, this presents tremendous value within and across government organizations and nations.

Synthetic data could serve as 'neutral ground' for collaborative military AI projects between nations. By using synthetic datasets that mirror real-world scenarios but do not contain sensitive information, **States could work together to develop and test AI systems, enhance interoperability, and share best practices without the risks associated with exchanging real military data.** Additionally, States can collaborate on developing collective synthetic datasets that can be used to train and test AI systems for enhanced interoperability – a key issue in the development

of military AI capabilities. This collaborative approach could foster greater cohesion among allied forces, improve the effectiveness of joint operations, and contribute to a more stable and secure international environment in the context of military AI.

Moreover, synthetic data represents an opportunity for States to develop common norms and guidelines surrounding its generation and use in the military domain. Due to its nascent nature, discussions surrounding synthetic data present the international community with an opportunity to develop **shared responsibility frameworks**. Multiple States could agree on governance principles and begin sharing best practices to move together in a systematic way towards the establishment of good practices for synthetic data generation and use. A multilateral taskforce on data governance could, for example, jointly address some issues and provide a forum for the development of procedures, processes, and accountability standards.



4. Conclusion

Synthetic data presents significant potential for advancing AI capabilities across both civilian and military domains. Its advantages – addressing data scarcity, enhancing privacy, and facilitating the creation of more representative and less biased datasets – make it a powerful tool. However, synthetic data is no panacea, and its use carries inherent risks. To extract the most benefit from this technology, it is of utmost importance that governance discussions begin to consider the issue. These efforts are currently in their infancy in both the civilian and military domains, and legal and normative ambiguities persist for the generation, processing, and use of synthetic data. To avoid a legal and normative vacuum leaving synthetic data risks unaddressed, efforts should be directed at identifying gaps in existing frameworks and providing clarity to users and generators of synthetic data. To this end, guidelines, the development of international technical standards, and cooperation with industry will be key.

In the military domain, data-related issues remain at the periphery of AI governance discussions. In this context, efforts should be directed not only at securing their place within these efforts, but specifically considering the effects of synthetic data use in the military domain. As synthetic data mostly complicates

existing governance challenges, as opposed to creating an entirely novel landscape, this will not necessarily entail the development of novel frameworks or regulations. It could entail a need for the application of military data best practices and concepts to the generation and use of synthetic data. Hence, more work on the issue is required to extend emerging military AI governance frameworks to synthetic data, clarifying how these practices and concepts can be applied.

As synthetic data grows increasingly prevalent, it brings with it not only governance challenges but also opportunities for collaborative international efforts which could have positive downstream effects on global military AI governance. To this end, multi-stakeholder efforts bringing together States, and importantly the private sector, will be instrumental.

Looking ahead, synthetic data will not be the last innovation in data science. This underscores the importance of creating governance frameworks adaptable to future developments. Building such frameworks with flexibility at their core will be essential to their sustainability, ensuring that they remain relevant as new technologies and use cases emerge.

Annex: Event Agenda and Participants

Introductory Remarks

Federico Mantellassi Researcher, United Nations Institute for Disarmament Research

Technology Breakfast on Synthetic Data and International Security

Dr. Eleonore Fournier-Tombs Head of Anticipatory Action and Innovation, Centre for Policy Research, United Nations University

Calum Inverarity Senior Researcher, Open Data Institute

Moderated by

Wenting He Associate Researcher, United Nations Institute for Disarmament Research

Multi-Stakeholder Dialogue on Synthetic Data: What Opportunities and Challenges for International Governance

Dr. Jane Pinelis Chief AI Engineer of the Applied Information Sciences Branch at Johns Hopkins University's Applied Physics Laboratory

Aldo Lamberti Founder and CEO, Syntheticus; Subject Matter Expert, European Commission; Vice-Chair, Industry Connection Synthetic Data, IEEE; Working Group Expert, Standard for Security and Trustworthiness Requirements in Generative Pretrained Artificial Intelligence (AI) Models, IEEE

Yasmin Afina Researcher, United Nations Institute for Disarmament Research; Expert, Global Commission on Responsible Artificial Intelligence in the Military Domain

Dr. Ana Beduschi Full Professor of Law with a Personal Chair at the University of Exeter; Director, Research Centre for Science, Culture and the Law, University of Exeter Law School

Moderated by

Federico Mantellassi Researcher, United Nations Institute for Disarmament Research



@unidir



/unidir



/un_disarmresearch



/unidirgeneva



/unidir



UNIDIR

Palais des Nations
1211 Geneva, Switzerland

© UNIDIR, 2024

WWW.UNIDIR.ORG