## Surveillance & Society

# Synthetic Training Data and the Reconfiguration of Surveillant Assemblages

Louis Ravn

Résumé de l'article

Synthetic training data promise considerable performance improvements in machine learning (ML) surveillance tasks, including such applications as crowd counting, pedestrian tracking, and face recognition. In this context, synthetic training data constitute techno-fixes primarily by virtue of acting as "edge cases"—data that are hard to come by in the "real world" yet straightforward to produce synthetically—which are used to enhance ML systems' resilience. In this dialogue paper, I mobilize Haggerty and Ericson's (2000) concept of the surveillant assemblage to argue that synthetic training data raise well-known, entrenched surveillance issues. Specifically, I contend that conceptualizing synthetic data as but one component of larger surveillant assemblages is analytically meaningful because it challenges techno-deterministic imaginaries that posit synthetic data as fixes to deep-rooted surveillance issues. To exemplify this stance, I draw from several examples of how synthetic training data are already used, illustrating how they may both intensify the disappearance of disappearance and contribute to the leveling of hierarchies of surveillance depending upon the surveillant assemblage that they reconfigure. Overall, this intervention urges surveillance studies scholarship to attend to how synthetic data reconfigure specific surveillant assemblages, with both problematic and emancipatory implications.

| Dialogue | Synthetic Training Data and the Reconfiguration of Surveillant Assemblages |
|---|---|

## Louis Ravn

University of Copenhagen, Denmark
Louisravn.cph@gmail.com

## Abstract

Synthetic training data promise considerable performance improvements in machine learning (ML) surveillance tasks, including such applications as crowd counting, pedestrian tracking, and face recognition. In this context, synthetic training data constitute techno-fixes primarily by virtue of acting as "edge cases"—data that are hard to come by in the "real world" yet straightforward to produce synthetically—which are used to enhance ML systems' resilience. In this dialogue paper, I mobilize Haggerty and Ericson's (2000) concept of the surveillant assemblage to argue that synthetic training data raise well-known, entrenched surveillance issues. Specifically, I contend that conceptualizing synthetic data as but one component of larger surveillant assemblages is analytically meaningful because it challenges techno-deterministic imaginaries that posit synthetic data as fixes to deep-rooted surveillance issues. To exemplify this stance, I draw from several examples of how synthetic training data are already used, illustrating how they may both intensify the disappearance of disappearance and contribute to the leveling of hierarchies of surveillance depending upon the surveillant assemblage that they reconfigure. Overall, this intervention urges surveillance studies scholarship to attend to how synthetic data reconfigure specific surveillant assemblages, with both problematic and emancipatory implications.

## Introduction

The predictive capabilities of machine learning (ML) systems used to be closely tethered to a particular paradigm of training data upon which they were based: training data derived from "real-world" people, objects, and events. This is of central importance to data-driven surveillance systems: the ability to surveil populations, to single out anomalous behaviors, to "detect" events, used to be exclusively circumscribed by the reliance on real-world training data (e.g., Drage and Frabetti 2023).

However, the advent of synthetic training data for ML is beginning to reconfigure this situation, with imminent implications for surveillance practices. One crucial discursive promise of synthetic training data, as highlighted by emerging scholarship (Helm, Lipp, and Pujadas 2024; Jacobsen 2023: 4, 2024; Steinhoff 2022), is that they facilitate the bespoke production of data that transcend "real-world" people, objects, and events; data that figure as "edge cases, corner cases, abnormalities, anomalies, and rarities," which allegedly increase a given ML system's resilience. Accordingly, the promise for video surveillance systems is that these won't break down in the face of anomalous events, instead facilitating more accurate and seamless surveillance (e.g., Delussu et al. 2024).

Rather than conceptualizing synthetic training data for ML as inducing definite surveillance consequences, in this dialogue paper, I propose that they enact differential implications dependent upon the surveillant assemblages (Haggerty and Ericson 2000) into which they become imbricated and that they reconfigure.

From this perspective, synthetic training data are neither good, nor bad, nor neutral in relation to surveillance; instead, they may be seen to intensify or dissipate the potentialities of the surveillant assemblages in which they figure. This is analytically meaningful insofar as it moves beyond regarding synthetic training data as techno-fixes, instead opening potential entry points to the critical study of how synthetic training data reshape surveillance practices.

To unfold this argument, this dialogue paper proceeds in three steps. First, I critically engage with Haggerty and Ericson's (2000) concept of the surveillant assemblage, highlighting both its strengths and weaknesses for approaching the surveillance implications of synthetic training data. Thereafter, I analyze several examples of how synthetic training data for ML reconfigures the potentialities of surveillant assemblages. Ultimately, I propose that surveillance studies scholarship may fruitfully attend to the specificities of surveillant assemblages that synthetic training data reconfigure rather than unilaterally determine. I close by suggesting that, while synthetic training data may intensify the disappearance of disappearance, they crucially also mark the re-appearance of disappearance: the impossibility of perfectly surveilling the ever-shifting (ab)normal.

## The Surveillant Assemblage in the Age of Synthetic Data

Haggerty and Ericson's (2000) influential theorization of the surveillant assemblage emanated from an interest in moving beyond an over-reliance on Orwellian and Foucauldian conceptualizations of surveillance. While they deemed Orwellian frameworks unhelpful to account for the spread of surveillance mechanisms across all societal sectors, Foucauldian analyses overstated the panoptic model of surveillance, which insufficiently engaged with contemporaneous technological developments. By contrast, the surveillant assemblage mobilizes conceptual tools advanced by the philosophy of Deleuze and Guattari (2013 [1980]), seeking to account for some key developments in contemporary surveillance: first, the orchestration of surveillance within heterogeneous assemblages consisting of a "limitless range of other phenomena such as people, signs, chemicals, knowledge, and institutions" (Haggerty and Ericson 2000: 608); second, the convergence of previously distinct surveillance systems; third, the levelling of hierarchies of surveillance due to the surveillant assemblage's rhizomatic character; and finally, an enactment of the human body as information (see also Galič, Tilman, and Koops 2017). Together, these developments are said to announce the disappearance of disappearance (Haggerty and Ericson 2000: 619), highlighting that the surveillant assemblage progressively forecloses the evasion of institutions' constant and all-encompassing monitoring of individuals. Effectively, the concept multiplies the potential entry points for scholarly analysis of how synthetic training data shape contemporary surveillance practices.

On the downside, the notion is not without its shortcomings, thus meriting a few qualifications. For one, because of their interest in underlining the convergence of distinct surveillance systems, Haggerty and Ericson (2000) wrote of "the" surveillant assemblage in the singular. In this paper, by contrast, I suggest speaking of surveillant assemblages in the plural, precisely because "the" surveillant assemblage is "multiple, unstable, and lacks discernible boundaries" (Haggerty and Ericson 2000: 609). Moreover, while the concept aptly accounts for the spread of surveillance systems into all kinds of societal sectors, more recent scholarship reminds us that hierarchies of surveillance—for instance along racial and socio-economic lines—remain firmly entrenched in the age of algorithmic surveillance (e.g., Benjamin 2019; Browne 2015; Gregory and Sadowski 2021).

Despite these limitations, I contend that the concept of the surveillant assemblage provides two crucial sensibilities to investigate how synthetic training data transform surveillance practices. First, it allows a challenge to imaginaries of synthetic data as a simple techno-fix that proponents often propound (Jacobsen 2023). Instead, to analyze how synthetic data reconfigures surveillant assemblages means to conceptualize it as but one component in larger amalgamations of human and non-human entities. Concretely, for example, it redirects attention from simplistic claims about the diversity of datasets (Jacobsen 2024) towards specific

questions around the surveillant assemblage in which synthetic training data are mobilized. This approach has central affinities with Susser and Seeman's (in this issue) call to consider the "social, political, and economic contexts in which synthetic data-driven technologies are being developed and implemented" and echoes Wiehn's (in this issue) focus on the "multitude of human and non-human agencies" that shape synthetic data's information life cycles. Second, and by the same token, this helps us ask how synthetic training data shape the agentic capacities of any given surveillant assemblage. This highlights that synthetic training data could—depending on which surveillant assemblage they are used in—either intensify or dissipate its surveillance potentialities. Importantly, this is not to say that synthetic training data are politically neutral; rather, it induces us to investigate how synthetic training data reconfigure the existing ethicopolitical tendencies of specific surveillant assemblages. Such an approach also helps avoid "criti-hype" (Vinsel 2021) that would inadequately portray synthetic training data as imbued with essential dystopic traits. The following two sections exemplify how to mobilize this perspective.

## Intensifying the Disappearance of Disappearance: Synthetic Human Behavior Data

Synthetic training data are already reconfiguring surveillant assemblages. Concretely, in a recent paper, Delussu, Putzu, and Fumera (2024) review how synthetic data can aid a variety of video surveillance applications. One of the central justifications for producing synthetic training data, the authors argue, is the scarcity of data on "abnormal events of interest for anomalous behaviour detection" (Delussu et al. 2024: 2). Such data, however, are highly relevant for several video surveillance tasks, including "crowd counting," "object and pedestrian detection and tracking," and "human and crowd behavior analysis" (Delussu et al. 2024: 2, 3, 5; see also Figure 1). A further example is constituted by the startup Anyverse (n.d.) whose synthetic data generator claims to enable customers to "build precise security and surveillance AI detection systems with synthetic data." Among the potential applications are "people detection" and the ability to "detect suspicious behaviors and dangerous situations" (Anyverse n.d.). As a final example, take ZeroEyes—a company that provides gun detection AI systems to US high schools and has now "incorporated an in-house synthetic-data pipeline to…generate data that is otherwise very hard to collect," ensuring that "all edge-cases, environments, and sensor technologies are represented within the data used to train DeepZero's object detection models" (Homeland Security Technology Consortium 2022). Delussu, Putzu, and Fumera's (2024) review paper, Anyverse's (n.d.) synthetic data generator, and the ZeroEyes pipeline (Homeland Security Technology Consortium 2022) are united by an aspiration to make surveillant assemblages more resilient to "abnormal" events.
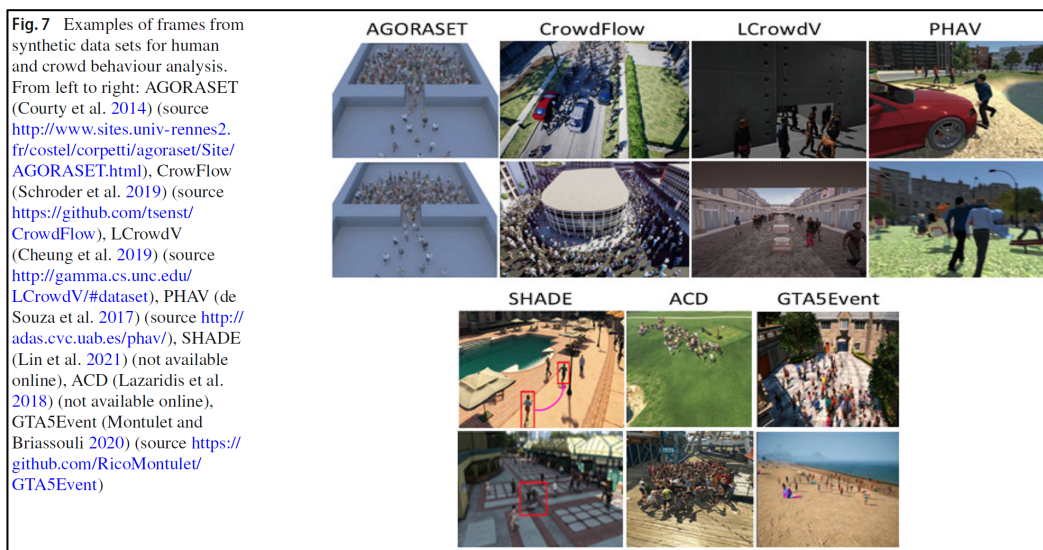


**Fig. 7** Examples of frames from synthetic data sets for human and crowd behaviour analysis. From left to right: AGORASET (Courty et al. 2014) (source http://www.sites.univ-rennes2.fr/costel/corpetti/agoraset/Site/AGORASET.html), CrowFlow (Schroder et al. 2019) (source https://github.com/tsenst/CrowdFlow), LCrowdV (Cheung et al. 2019) (source http://gamma.cs.unc.edu/LCrowdV/#dataset), PHAV (de Souza et al. 2017) (source http://adas.cvc.uab.es/phav/), SHADE (Lin et al. 2021) (not available online), ACD (Lazaridis et al. 2018) (not available online), GTA5Event (Montulet and Briassouli 2020) (source https://github.com/RicoMontulet/GTA5Event)

*Figure 1:* Synthetic data for human and crowd behavior analysis (Delussu et al. 2024).

Crucially, approaching such projects from the perspective of the surveillant assemblages that are thereby reconfigured helps reject the notion that synthetic training data are somehow by themselves capable of perfecting surveillance capabilities. From this perspective, this intensification is not simply attributable to synthetic training data alone; instead, it must be viewed as the overall outcome of a heterogeneous surveillant assemblage whose constituent components—including synthetic training data, algorithms, institutions, policy frameworks, and human practitioners—co-function to enact this tendency. To illustrate, this perspective requires recalling that surveillant assemblages taking the form of event/people detection systems have long been shown to disproportionately target minorities (e.g., Amoore 2020; Benjamin 2019; Drage and Frabetti 2023). To the extent that synthetic training data do improve the performance of such systems, their intensification of a surveillant assemblage's tendency toward the disappearance of disappearance may "fuel surveillance" (Susser and Seeman, in this issue) and target such minorities most acutely. However, a critique of synthetic training data alone would fail to adequately account for how they interact with other components of surveillant assemblages that collectively enact the disappearance of disappearance.

## Leveling Surveillance Hierarchies: VFRAME's Synthetic Munition Data

As a counterexample to the potential for synthetic training data to intensify the disappearance of disappearance, the VFRAME project (Harvey 2019) showcases how they can also be mobilized to emancipatory ends. Concretely, the project has been devoted to developing synthetic training data for object recognition systems capable of identifying the use of illegal cluster munitions in active war zones, such as Syria and Russia (Murgia 2021). The reason for the indispensability of synthetic data here is that "videos containing cluster munitions are typically limited in quantity" and are captured in "extreme situations" that substantially degrade their quality as training data (Harvey 2019). Moreover, they offer substantial safety improvements over the crowdsourcing of photos of live munitions. One of the project's broader aims is to support human rights investigators and to work towards achieving accountability.



*Figure 2: VFRAME project's synthetic training data for conflict zone objects*

The synthetic training data produced and mobilized by the VFRAME project, seen through the lens of surveillant assemblages, effectuate a tendency entirely different from the one described in the previous

examples. In this case, the synthetic training data reconfigure the potentialities of a surveillant assemblage from below, one that is committed to aiding human rights investigations in a struggle to hold accountable war crimes perpetrated by the Syrian and Russian armed forces. Put differently, here we see how synthetic training data can intensify what Haggerty and Ericson (2000: 606) dub the "rhizomatic leveling of the hierarchy of surveillance": powerful political actors themselves become subject to new modes of surveillance. Here, I agree with Susser and Seeman (in this issue) when they argue: "If deployed thoughtfully, synthetic data could be a valuable tool."

## Concluding Remarks: Synthetic Training Data and the Re-appearance of Disappearance

In this dialogue paper, I have argued that synthetic training data differentially reconfigure surveillant assemblages. The central benefit of this analytical approach is that it foregrounds how synthetic training data—rather than acting techno-deterministically—may both intensify or dissipate the potentialities of different surveillant assemblages. The examples served to highlight that synthetic training data may augment the tendency of surveillant assemblages to accelerate the disappearance of disappearance—as in the case of event detection and crowd behavior analysis—but may also give rise to emancipatory projects that level surveillance hierarchies—as illustrated by the VFRAME project. The implication of this approach for future surveillance studies scholarship is to eschew essentialist perspectives on synthetic data in favor of analytical sensitivity to the specificities of a given surveillant assemblage—and its tendencies for the intensification or levelling of surveillance hierarchies.

While one argument of this dialogue paper has been to suggest that synthetic training data may intensify the disappearance of disappearance, this argument also invites a final problematization. Synthetic training data embody and enact the dream of working towards a "world model" (Amoore et al. 2024)—an AI system that would be able to respond to all events encountered in the real world, a surveillant assemblage whose gaze nothing could escape. This dream of a gapless surveillant assemblage, however, must necessarily fall short of itself: the project of synthetically generating (ab)normal events to make ML systems resilient to these—to make (ab)normalities surveillable—overlooks that their generation always amounts to a reshaping of the (ab)normal (Jacobsen 2023). The impossibility of capturing the (ab)normal by synthetically generating it highlights the possibility of the re-appearance of disappearance: in spite of bespoke synthetic training data, there will remain people, objects, and events that evade surveillant assemblages.

## References

Amoore, Louise. 2020. *Cloud Ethics: Algorithms and the Attributes of Ourselves and Others*. Durham, NC: Duke University Press.
Amoore, Louise, Alexander Campolo, Benjamin Jacobsen, and Ludovico Rella. 2024. A World Model: On the Political Logics of Generative AI. *Political Geography* 113: 1–9.
Anyverse. N.d. Build Precise Security and Surveillance AI Detection Systems with Synthetic Data. https://anyverse.ai/use-case-security-surveillance/ [accessed July 25, 2024].
Benjamin, Ruha. 2019. *Race after Technology: Abolitionist Tools for the New Jim Code*. Medford, MA: Polity Press.
Browne, Simone. 2015. *Dark Matters: On the Surveillance of Blackness*. Durham, NC: Duke University Press.
Deleuze, Gilles and Felix Guattari. 2013 [1980]. *A Thousand Plateaus*. London: Bloomsbury Academic.
Delussu, Rita, Lorenzo Putzu, and Giorgio Fumera. 2024. Synthetic Data for Video Surveillance Applications of Computer Vision: A Review. *International Journal of Computer Vision* 132: 4473–4509.
Drage, Eleanor, and Federica Frabetti. 2023. The Performativity of AI-powered Event Detection: How AI Creates a Racialized Protest and Why Looking for Bias Is Not a Solution. *Science, Technology, & Human Values* 49 (5): 1045–1072.
Galič, Maša, Tjerk Tilman, and Bert-Jaap Koops. 2017. Bentham, Deleuze and Beyond: An Overview of Surveillance Theories from the Panopticon to Participation. *Philosophy & Technology* 30: 9–37.
Gregory, Karen, and Jathan Sadowski. 2021. Biopolitical Platforms: The Perverse Virtues of Digital Labour. *Journal of Cultural Economy* 14 (6): 662–674.
Haggerty, Kevin D., and Richard V Ericson. 2000. The Surveillant Assemblage. *British Journal of Sociology* 51 (4): 605–622.
Harvey, Adam. 2019. 3D Rendered Synthetic Data. VFRAME. https://vframe.io/3d-rendered-data/ [accessed May 31, 2024].
Helm, Paula, Benjamin Lipp, and Roser Pujadas. 2024. Generating Reality and Silencing Debate: Synthetic Data as Discursive Device. *Big Data & Society* 11 (2): https://doi.org/10.1177/20539517241249447.

Homeland Security Technology Consortium. 2022. HS Tech Capabilities Statement. https://hstech.ati.org/wp-content/uploads/2022/06/ZeroEyes-Inc.-Capabilities-Statement.pdf [accessed October 15, 2024].

Jacobsen, Benjamin. 2023. Machine Learning and the Politics of Synthetic Data. *Big Data & Society* 10 (1): 1–12.

———. 2024. The Logic of the Synthetic Supplement in Algorithmic Societies. *Theory, Culture & Society* 41 (4): 41–56.

Murgia, Madhumita. 2021. Researchers Train AI on "Synthetic Data" to Uncover Syrian War Crimes. *Financial Times*, December 09. https://www.ft.com/content/8399873e-0dda-4c87-ba59-0e2678166fba [accessed February 23, 2024].

Susser, Daniel, and Jeremy Seeman. 2024. Critical Provocations for Synthetic Data. *Surveillance & Society* 22 (4): 453–459.

Steinhoff, James. 2022. Toward a political economy of synthetic data: A data-intensive capitalism that is not surveillance capitalism? *New Media & Society* 26 (6): 3290–3306.

Vinsel, Lee. 2021. You're Doing It Wrong: Notes on Criticism and Technology Hype. *Medium*, February 01. https://sts-news.medium.com/youre-doing-it-wrong-notes-on-criticism-and-technology-hype-18b08b4307e5 [accessed July 24, 2024].

Wiehn, Tanja. 2024. Synthetic Data: From Data Scarcity to Data Pollution. *Surveillance & Society* 22 (4): 472–476.