

# Data Protection by Design: Building the foundations of trustworthy data sharing

Sophie Stalla-Bourdillon, Gefion Thuermer,\* Johanna Walker and Laura Carmichael

University of Southampton, UK

\*[Gefion.Thuermer@soton.ac.uk](mailto:Gefion.Thuermer@soton.ac.uk)

## Abstract

Data trusts have been conceived as a mechanism to enable the sharing of data across entities where other formats, such as open data or commercial agreements, are not appropriate, and make data sharing both easier and more scalable. Although the form and purposes of data trusts are currently a topic of much academic discussion, a broadly accepted definition has not yet emerged. The concept of the ‘data trust’ requires further disambiguation from other facilitating structures such as data collaboratives. Irrespective of the terminology used, attempting to create trust in order to facilitate data sharing, and create benefit to individuals, groups of individuals, or society at large, requires at a minimum a process-based mechanism, i.e. a workflow, that should have a trustworthiness-by-design approach at its core. Data protection by design (DPbD) should be a key component of such an approach.

## Keywords

Data-Driven Innovation; Data Protection by Design (DPbD); Data Trusts; General Data Protection Regulation (GDPR); Organisational DPbD Process.

## Introduction

Data protection by design (DPbD) was recently introduced into law via Article 25 of the General Data Protection Regulation (GDPR). The requirement of DPbD builds upon research and applied work conducted in the field since the end of the 90s (Cavoukian, 2009). Article 25(1) places a legal obligation on controllers<sup>1</sup> to “implement appropriate

*organisational and technical measures [...] designed to implement data-protection principles [...] in order to meet the requirements of this Regulation and protect the rights of data subjects”*. DPbD therefore plays a key role in enabling and demonstrating compliance with the GDPR.

In this paper we address the question of how the requirements of DPbD should shape the development of data trusts (this concept is explored in more detail below). We will argue that both technical and organisational requirements are foundational to ensuring trustworthy data sharing. We further insist on the necessity of starting with organisational measures and creating a DPbD process, which are prerequisites to the selection of appropriate technical measures.

In order to strengthen our claim, we also draw on our experience as interdisciplinary members of Data Pitch - an open innovation programme - to inform our proposed approach. Data Pitch aims to bring together data providers (i.e. corporate and public sector organisations) to share data with successful programme applicants (i.e. startups and SMEs) to re-use for innovation purposes. The project launched in January 2017 and will end in December 2019. It is funded by the European Union’s Horizon 2020 Research and Innovation Programme.<sup>2</sup>

## Data trusts

Data trusts have been conceived as a mechanism to enable the sharing of data across entities where other formats, such as open data or commercial

---

<sup>1</sup> The following legal definition of controller is provided by Article 4(7) of the GDPR: “‘controller’ means the natural or legal person, public authority, agency or other body which, alone or jointly with others, determines the purposes and means of the processing of

---

*personal data; where the purposes and means of such processing are determined by Union or Member State law, the controller or the specific criteria for its nomination may be provided for by Union or Member State”*.

<sup>2</sup> For more information about Data Pitch visit the project website at <https://datapitch.eu/> [last accessed on 10 May 2019].

agreements, are not appropriate, and make data sharing easier, more scalable (Hall & Pesenti, 2017) and mutually beneficial for members (Lawrence, 2016). Although the form and purposes of data trusts are currently a topic of much discussion (e.g. Alsaad et al., forthcoming; Hardinges, 2018; Wylie & McDonald, 2018; O'Hara, 2019), a broadly accepted definition has not yet emerged. This is in part because data trusts may be of benefit in data-driven innovation, as well as many other situations such as personal or health data management (Lawrence, 2016) and security, safety and efficiency, like in the Internet of Food Things project.<sup>3</sup> The concept of the 'data trust' requires further disambiguation from other facilitating structures such as data collaboratives (Susha et al., 2017). Furthermore, the use of data trusts as an internal data sharing methodology, as it is established by firms such as Truata,<sup>4</sup> has created further ambivalence around the term.

Notwithstanding the lack of convergence on a precise definition for the term 'data trust' - the key point to take from this semantic debate is that the design, development and utilisation of robust mechanisms for responsible data sharing are crucial to engender trust and ultimately drive forward data-driven innovation.

### **The need for increased data sharing**

Data-driven innovation is regarded as a new 'growth area' for the global economy (OECD, 2015). Given data-driven innovation is contingent upon "*the use of data and analytics to improve or foster new products, processes, organisational methods and markets*" (OECD, 2015), it is vital that interested parties have lawful access and rights to (re)use vast amounts of robust data where necessary and appropriate. It is therefore unsurprising that a key obstacle to the growth of data-driven innovation is a lack of data sharing (Mehonic, 2018; Skelton, 2018) – also referred to as the 'data-pooling problem' (Mattioli, 2017). For instance, a deficiency of training datasets

has led to the failure of multiple private and public machine learning initiatives (Mehonic, 2018).

There are numerous reasons why organisations may be reticent to share data for innovation purposes, including concerns over privacy, data quality, free-riding, competition, reputation and proprietary issues (Mattioli, 2017). Data trusts are proposed as one approach that could encourage increased data sharing and re-usage within a wider data-driven innovation strategy;<sup>5</sup> especially for personal and anonymised data (Edwards, 2004; Reed & Ng, 2018).

### **Sharing personal data**

The GDPR applies only to information pertaining to an identified or identifiable natural person. In many instances of data sharing, however, as has been shown by Data Pitch, the data that is shared is, or could become personal data. With sensitive sectors such as healthcare and research increasingly utilising artificial intelligence, this is only likely to increase (Lawrence, 2016).

Given the broad definition of personal data, it is imperative that those designing, developing and utilising data trusts remain compliant with the GDPR. Due to the key role of DPbD in enabling and demonstrating compliance with the GDPR, it is vital that we further explore how the requirement of DPbD does or should impact upon the construction of data trusts. As O'Hara (2019) argues, the purpose of data trusts is to 'support trustworthy data processing', which is achieved by applying constraints that go beyond the law. This requires determining what the law actually mandates and adding to its prescription.

### **Data protection by design (DPbD)**

Despite the concept of privacy-by-design being well established in principle, its technical implementation has been rather limited thus far (Tsormpatzoudi et al., 2016; Hansen, 2016). Given Article 25 of the GDPR now directly places a legal obligation on controllers to practise DPbD, there is a real incentive for its

---

<sup>3</sup> For more information about the Internet of Food Things project see: <<https://www.foodchain.ac.uk/>> [last accessed on 10 May 2019].

<sup>4</sup> For more information about about Truata see: <<https://www.truata.com/>> [last accessed on 10 May 2019].

---

<sup>5</sup> For further elements of such a strategy see e.g. the British Academy & The Royal Society (2017) report which focuses on the need for 'a renewed governance framework' and stewardship body for trustworthy data sharing.

widespread implementation. Especially, as pursuant to Article 83(4), any infringements of Article 25 may result in *“administrative fines up to 10 000 000 EUR, or in the case of an undertaking, up to 2 % of the total worldwide annual turnover of the preceding financial year, whichever is higher”*.

## Organisational as well as technical measures

It still remains difficult to find practical DPbD guidance that provides extensive coverage of both the organisational and technical dimensions mandated by Article 25. When DPbD is presented and explained, the focus is often set on its technological dimension (Wiese Schartum, 2016) – the engineering of data protection principles through design strategies and privacy enhancing techniques (Deng et al., 2011; Danezis et al., 2015). Less emphasized is that the requirement also has a vital organisational dimension - i.e. Article 25(1) places a legal obligation on controllers to *“implement appropriate organisational and technical measures [...] to protect the rights of data subjects.”* For instance, organisational measures may refer to the adoption of particular procedures and the selection of particular individuals to decide and action various aspects of data processing, including the type of privacy-enhancing technologies (PETs) to be employed across the data sharing and re-usage lifecycle (The Royal Society, 2019).

## Seven core data-protection principles

This organisational dimension of DPbD implies a particular workflow i.e. a series of accountable decisions and actions taken by responsible individuals with appropriate expertise prior to the commencement of the data processing activities under consideration. Note that an organisation may also choose to automate some of these decisions for reasons of scalability; in that case, the accountable decisions by individuals concern the design of the automation.

The main nodes of this workflow echo the seven core data-protection principles at the heart of the GDPR and directly referred to by Article 25(1): (i) ‘lawfulness, fairness and transparency’; (ii) ‘purpose

limitation’; (iii) ‘data minimisation’; (iv) ‘accuracy’; (v) ‘storage limitation’; (vi) ‘integrity and confidentiality’; and (vii) ‘accountability’. These data protection principles are outlined in GDPR Article 5, and impose high level restrictions upon how personal data should be collected and used, how data quality should be ensured and maintained, and how personal data should be protected. These principles are particularly important when data is not only processed internally, but also shared between organisations.

## DPbD workflow

Essentially, before any processing starts, the data controller should put in place technical and organisational measures in order to facilitate compliance with the data protection principles as listed in Article 5. Article 25 thus refers to Article 5. The basic structure for a DPbD workflow - comprising eight nodes I. to VIII. - can be derived from Article 5 of the GDPR as follows:

- I.** Define your **purpose** for sharing data in this instance. (See Article 5(1)(b) - ‘purpose limitation’.)
- II.** Identify your **legal basis** for sharing data in this instance. (See Article 5(1)(b) - ‘purpose limitation’.)
- III.** Determine which data are **necessary** for your specific purpose. Ensure that you **reduce**: (i) any **non-essential processing** activities; and (ii) the **amount of data** required - e.g. mask or hide direct identifiers that are not required for processing in this instance. **If you can anonymise data just do it!** (See Article 5(1)(c) - ‘data minimisation’.)
- IV.** Set a **data retention period** in relation to the purpose. (See Article 5(1)(e) - ‘storage limitation’.)
- V.** Ensure the data to be shared are **accurate**. (See Article 5(1)(d) - ‘accuracy’.)
- VI.** Verify that the processing is **fair**. (See Article 5(1)(a) - ‘lawfulness, fairness and transparency’.)
- VII.** Ensure the data are **not altered or disclosed without permission** - e.g. define who is eligible to access data - and the processing is **confidential**. (See Article 5(1)(f) - ‘integrity and confidentiality’.)
- VIII.** Ensure the processing is **transparent** and **monitored**, e.g. by logging activities so that you can know what is happening with the data (and ultimately demonstrate compliance). **Best practice: assess risk**

**before initiating processing.** (See Article 5(1)(a) - ‘lawfulness, fairness and transparency’ and Article 5(2) - ‘accountability’. )

If data trusts are the mechanism through which data sharing will be enabled in the future, it is therefore clear that they should embed a DPbD workflow, and thereby be underpinned by organisational and technical measures as defined by GDPR Article 25.

## Two lessons learnt from Data Pitch

After familiarisation with the DPbD workflow in principle, the next step towards trustworthy data sharing is determining how to carry out this DPbD workflow in practice. From our experience with Data Pitch, we raise two key organisational lessons learnt for successful implementation of a DPbD workflow:

### **(1) Strong engagement across business functions for responsible data sharing and re-usage.**

Responsible data sharing can be viewed as a chain of decisions and actions.<sup>6</sup> For instance, a company may consider: why it may wish to share data; what kind of entity might be eligible to access the data; what the purpose of data sharing is; what authority it has to share the data; and how it might ensure that the data sharing is compliant. It is extremely unlikely these decisions and actions will be taken by one person alone. Such decision-making needs strong engagement across business functions - from security experts and data scientists to data protection officers and business strategists.<sup>7</sup> Senior-level support is crucial to overcome ambiguities in the decision-making process.

**(2) An agreed process for accountable decision-making.** It is vital that there is a process in place where organisational and technical measures are selected to uphold the seven core data-protection principles across the lifecycle of the data processing activity (e.g. over the course of an open innovation programme). These organisational and technical measures must be appropriate i.e. well-suited to the

specific context and purpose of the data processing activity in question.

## Embedding a DPbD approach within Data Trusts

We therefore argue that the effective entrenchment of DPbD within the construction of data trusts requires (at least):

- (a)** Cognisance of the minimum legal requirements for DPbD - including both its organisational and technical dimensions - as mandated by Article 25 together with its accompanying DPbD workflow located in Article 5.
- (b)** An organisational DPbD process that addresses (at minimum) the legal requirements for DPbD across the entire data trust lifecycle (i.e. from initial plans for creating a data trust to a data trust in operation).
- (c)** Strong, cross-functional business engagement that brings the required expertise to successfully shape, execute and appraise the DPbD process.

Given that we have already examined both points (a) and (c), we will now turn our attention to what an organisational DPbD process for data trusts is likely to involve. Note that we are only able to signpost some key aspects of a DPbD process to act as a point of reference for data trusts - there is no one-size-fits all approach. A DPbD process must always take into account the specific context and purpose of the data sharing and re-usage activities in question.

**Scenario.** A few organisations are interested in working together to form a new data trust. This data trust would be centred around the creation of a data pool so as to improve their current levels of innovation activities. This data pool would involve each organisation sharing their data with authorised members of the data pool i.e. the other organisations and (potentially) third parties. A significant amount of these datasets are likely to be personal or anonymised.

**Three layer approach.** As there is no agreed configuration for data trusts, we represent data trusts through three core layers that feature in many data sharing ecosystems. These three core layers

---

<sup>6</sup> For instance, Bunting & Lansdell (2019) examine how to design ‘decision-making processes for data trusts’.

<sup>7</sup> Tsormpatzoudi et al. (2016) also highlight the importance of an interdisciplinary approach for effective DPbD implementation.

comprise: (1) the data layer - where interested parties make plans to create a data pool; (2) the access layer - where pooled data are made discoverable through a data trust; and (3) the process layer - where pooled data are approved for (re-)usage via the data trust. Note that data may be stored centrally (e.g. all datasets will be held by the data trust) or disparately (e.g. individual datasets will be held by different parties).

### **(1) The data layer: preparation of data sources.**

DPbD should be embedded into the plans for the new data trust through the following process:

- (i)** Ensure that all potential members are aware of the legal requirements for DPbD (in particular Article 25 and Article 5) - and the overarching DPbD process for the data trust. Recognise any gaps in knowledge - and provide further training and guidance where necessary.
- (ii)** Identify the appropriate persons across all organisations that have the authority and required expertise to decide and action on the pooled data.
- (iii)** Provide clear guidelines for reviewing data in the planned data pool, including guidance on: how to assess whether data can be understood as personal data; and high risk processing.
- (iv)** Apply standardised procedures for the removal of unnecessary personal data. The data minimisation principle should directly impact the way data sets are redacted and presented. For instance, direct identifiers should be stripped away as often and as early as possible to minimise the personal data contained in data sets.

**(2) The access layer: discovery of pooled data.** The datasets within the planned data pool should then be made discoverable to authorised parties through metadata. DPbD should be embedded into the access layer of the new data trust through the following process:

- (v)** Define who is eligible to access the pooled data, and place limitations on who accesses the data, and why. These boundaries are defined around the purpose of the data trust itself, but also include a clear distinction between the raw data and metadata.
- (vi)** Provide standardised access through centralised technical solutions, underpinned by monitoring and auditing processes, or provide the governance

processes to manage peer-to-peer direct sharing that enable auditing.

### **(3) The process layer: approval of pooled data**

**(re)usage.** The (re)usage of datasets within the data pool should be managed by the data trust, which should be in the position to make informed decisions about whether (or not) to permit data sharing with interested parties. DPbD should be embedded into the process layer of the new data trust through the following process:

- (vii)** Control data usage through standardised risk assessments. Once the processing purpose and data sources are confirmed, there should be an assessment of the intended versus allowed use of the data, to guarantee in particular the lawfulness and fairness of processing and ultimately the impact upon the rights and freedoms of data subjects. Such an assessment should be done in context of the intended use, and therefore renewed each time a new purpose is suggested. Once again risk assessment is key for accountability. Risk assessment is iterative - it should start as early as the pooling phase and be reviewed at the inception of the re-usage phase.
- (viii)** Ensure that data are tailored to queries. Queries that are interested in aggregates should only be responded to with aggregate data. Where raw data is required, this should be limited to the necessary attributes. Traditional techniques based on extract, transform, load should be reconsidered as they tend to create unnecessary movements of data. The potential for PETs, such as differential privacy, should be fully explored at this stage.

## **Conclusion**

While the concept of data trusts is neither new or precisely-defined, data trusts are conceived as an important tool to engender trust as part of a wider response to data sharing barriers that may impede data-driven innovation.

Given the likelihood that the data to be shared may be personal data or could become personal data (e.g. through purpose or result of use, re-identification), it is vital that data trusts embed DPbD through the implementation of appropriate organisational and technical measures that uphold the seven core

data-protection principles at the heart of the GDPR. The DPbD workflow defined by Article 5 is therefore key to the effective implementation of the appropriate organisational and technical safeguards that lead to trustworthy data sharing.

There is an opportunity for data trusts to lead the way with the practical implementation of DPbD by giving equal attention to its organisational and technical dimensions. Strong engagement across business functions will be critical for the creation and adoption of well-considered processes that embed DPbD.

## Acknowledgements

Data Pitch is funded by the European Union's Horizon 2020 Research and Innovation Programme under the Grant Agreement 732506.

## References

- Cavoukian, A. (2009). Privacy by Design: The 7 Foundational Principles. Information and Privacy Commissioner of Ontario, Canada.
- Data Pitch <<https://datapitch.eu/>>.
- Hall, W., & Pesenti, J. (2017). Growing the Artificial Intelligence Industry in the UK. Independent Review, Retrieved from <<https://www.gov.uk/>>.
- Lawrence, N. (2016, June 3). Data Trusts Could Allay Our Privacy Fears. *The Guardian*, Retrieved from <<https://www.theguardian.com/uk/>>.
- Alsaad, A, O'Hara, K., & Carr, L. Forthcoming, 30 June - 3 July 2019. Institutional Repositories as a Data Trust Infrastructure. In Proceedings of Web Science 2019. Boston, MA: ACM.
- Hardinges, J. (2018, July 10). What is a Data Trust? Open Data Institute Blog, Retrieved from <<https://theodi.org/>>.
- Wylie B., & McDonald, S. (2018, October 9). What is a Data Trust? Centre for International Governance Innovation (CIGI), Retrieved from <<https://www.cigionline.org/>>.
- O'Hara, K. (2019). Data Trusts: Ethics, Architecture and Governance for Trustworthy Data Stewardship. Web Science Institute White Paper, Retrieved from <<https://eprints.soton.ac.uk/>>.
- Internet of Food Things Network Plus <<https://www.foodchain.ac.uk/>>.
- Susha, I., Janssen, M., & Verhulst, S. 2017. Data Collaboratives as a New Frontier of Cross Sector Partnerships in the Age of Open Data: Taxonomy Development. In Proceedings of the 50th Hawaii International Conference on System Sciences, 2691–2700.
- Truata <<https://www.truata.com/>>.
- Organisation for Economic Co-operation and Development (OECD). (2015). Data-Driven Innovation: Big Data for Growth and Well-Being. Report, Retrieved from <<https://www.oecd.org/>>.
- Mehonic, A. (2018, October 3). Can data trusts be the backbone of our future AI ecosystem? The Alan Turing Institute Blog, Retrieved from <<https://www.turing.ac.uk/>>.
- Skelton, S. K. (2018, November 30). New forms of governance needed to safely and ethically unlock value of data. ComputerWeekly.com, Retrieved from <<https://www.computerweekly.com/>>.
- Mattioli, M. (2017). The Data-Pooling Problem. Berkeley Technology Law Journal, 32(1), 179-236.
- Edwards, L. (2004). The Problem with Privacy: A Modest Proposal. International Review of Law, Computers & Technology, 18(3), 263-294.
- Reed C., & Ng, I. (2019, February 14). Data Trusts as an AI Governance Mechanism: Response to the Singapore Personal Data Protection Commission. Retrieved from <<https://www.ssrn.com/>>.
- British Academy & The Royal Society. (2017). Data management and use: Governance in the 21st century. Report, Retrieved from <<https://royalsociety.org/>>.
- Tsormpatzoudi, P., Berendt, B., & Coudert, F. (2016). Privacy by Design: From Research and Policy to Practice – the Challenge of Multi-disciplinarity. In B. Berendt, T. Engel, D. Ikonou, D. Le Métayer & S. Schiffner (Eds.), Privacy Technologies and Policy (pp. 199-212). Springer, Cham.
- Hansen, M. (2016). Data Protection by Design and by Default à la European General Data Protection Regulation. In A. Lehmann, D. Whitehouse, S. Fischer-Hübner, L. Fritsch & C. Raab (Eds.), Privacy and Identity Management. Facing up to Next Steps (pp. 27-38). Springer, Cham.
- Wiese Schartum, D. (2016). Making privacy by design operative. International Journal of Law and Information Technology, 24 (2), 151-175.
- Deng, M., Wuyts, K., Scandariato, R., Preneel, B., & Joosen, W. (2011). A privacy threat analysis framework: Supporting the elicitation and fulfillment of privacy requirements. Requirements Engineering, 16 (1), 3–32.
- Danezis, G., Domingo-Ferrer, J., Hansen, M., Hoepman, J-H., Le Métayer, D., Tirtea, R., & Schiffner, S. (2015). Privacy and Data Protection by Design - from Policy to Engineering. European Union Agency for Network and Information Security (ENISA) Report, Retrieved from <<https://www.enisa.europa.eu/>>.
- The Royal Society. (2019). Protecting privacy in practice: The current use, development and limits of Privacy Enhancing Technologies in data analysis. Report, Retrieved from <<https://royalsociety.org/>>.
- Bunting M., & Lansdell S. (2019). Designing decision-making processes for data trusts: lessons from three pilots. Report, Retrieved from <<https://theodi.org/>>.