# Designing a Service for Compliant Sharing of Sensitive Research Data

Aakash Sharma[0000−0003−3965−3271], Thomas Bye Nilsen[0000−0002−3602−451X],
Sivert Johansen[0000−0003−3442−2466], Dag Johansen[0000−0001−7067−6477], and
Håvard D. Johansen[0000−0002−1637−7262]

UiT – The Arctic University of Norway, Tromsø, Norway
aakash.sharma@uit.no

**Abstract.** Data-driven research is increasingly becoming fueled by access to open datasets, often shared publicly on the Internet. However, many research projects study sensitive data. They cannot easily participate in this shift as access to their data is significantly controlled by ethical and regulatory constraints. This paper discusses the requirements for building a service that enables sensitive data for sharing between collaborators in a controlled manner. We argue that a decentralized service that maintains metadata, a global view on all data usage, and active policy in combination with local monitoring and security enforcement can provide automated compliance checking. With such a service, researchers can share sensitive data with a broader community rather than limiting access to only core project members.

**Keywords:** open data · data-sharing · compliance · sensitive data

## 1 Introduction

The Internet has changed the way researchers work, collaborate, and disseminate. Open Science is a cultural change [2]. The arguments concerning the benefits of open data are well established such as allowing researchers to explore existing datasets in new ways [1, 4, 10]. Volunteers (hereafter written as *subjects*) contribute their data for research. The trustworthiness of the institution conducting the research plays a key role in a subject's willingness to contribute [17]. Privacy leaks or misuse can damage the reputation and affect future research studies [1, 4]. Fears of misuse of data may also restrict many researchers from sharing data openly [4, 20].

Researchers argue that these concerns can be mitigated by building accountability in research data sharing and processing [8]. Recent regulations such as General Data Protection Regulation (GDPR) require researchers to abide by a subject's consent for data processing. GDPR also provides workarounds for public-funded research such as entrusting a Regional Ethics Committee (REC) or an Institutional Review Board (IRB) to protect subjects' privacy. The public's trust in researchers is fragile [1]. The growing concerns towards data breaches,

data brokers, and indiscriminate profiling of users might change subjects' willingness to participate or continue participating in a research project.

The guidelines and complexity of compliance is a tedious job and requires a complex understanding of legal, ethical, and regulatory issues [5]. Often institutions employ large teams to assist researchers in making their data openly available [18]. Researchers' concern about misuse of their data is the leading reason given for not sharing data [4, 20]. As a result, research data may end up in silos accessible only to a limited few. A lot of work has been done for simplifying regulatory requirements, easy to create toolkits [5, 18, 21] and metadata formats for making research openly accessible [7, 12, 21]. However, additional regulations such as *data sovereignty* [13] may further restrict open data. For example, medical research data is heavily regulated. Often movement of such sensitive data is restricted outside a nation's physical boundaries.

Open science and collaboration requires access to the same data regardless of international borders [2, 11]. As argued earlier, there might be regulatory restrictions on sharing sensitive data. The cloud provides an interesting platform for scalability and access for researchers. Our contribution is a scalable cloud-based service that allows researchers to do analysis on sensitive data regardless of their location. We discuss related work in Section 2 and present the requirements for such a system in Section 3. Later, in Section 4, we present our system's design and discuss how it addresses the requirements.

## 2   Related Work and Discussion

Dataverse [7] is popular data repository for sharing research data which, currently hosts tens of thousands of datasets. However, Dataverse does not support sensitive data. Datatags system [18] translates security and access requirements for sensitive data into a model set of six tags. Their approach simplifies the complex workings and guidelines for sharing datasets responsibly as they provide a decision tree for picking a correct tag for different requirements. The Datatags approach simplifies complex information flows for IRBs and RECs without specifying mechanisms for automated audits or enforcement.

Automatable Discovery and Access Matrix (A-DAM) [21] provides a *profile* as regulatory metadata for responsible sharing of biomedical assets. A-DAM provides a semi-automated approach for analyzing ethical and regulatory requirements for sharing and processing research data. Policy changes require a newer profile and it is reevaluated. Maguire et al. [9] proposed a metadata-based architecture for accountability. Similar to A-DAM, their approach attaches static policies to data, which are then verified by a gatekeeper service. Their approach introduces validation against *context* by the gatekeeper. For sensitive data, they only briefly discuss adding encryption and keeping the keys under the control of the gatekeeper. In our earlier work Lohpi [15, 16], we argued that the changes occurring during a project's life cycle may affect its data security policies. Thus, we built support for accountability by keeping data security policies up-to-date securely and efficiently. We build upon existing works to semi-automate regulatory,

ethical, and legal requirements. Our contribution is a scalable cloud-based service that allows researchers to analyze sensitive data regardless of their location. The service provides transparency to stakeholders such as subjects regarding data sharing and data use.

Axelsson and Schroeder [1] argue that public trust is fragile and once broken it might take years to re-build. Compliance and transparency are crucial for maintaining the fragile public trust in researchers. And keeping sensitive data open is a challenging [1]. Even experts in various fields feel the lack of assurances [6] in existing practices. We provide compliance with audit-able data sharing of sensitive datasets. Many researchers have argued for building transparency for data-sharing, usage, and privacy protections [8, 14]. Along these themes, our approach addresses compliant sharing of sensitive data especially, data sovereignty. Thus, allowing sensitive datasets to reach a broader audience, while fulfilling regulatory and legal compliance requirements. The built-in transparency allows stakeholders such as subjects to understand their data usage and answering questions like who, whom, and where, about their contributed data. Such transparency may improve the public's trust and participation in studies that rely heavily on volunteers.

## 3 Requirements for the Service

As argued earlier, open data should be able to reach a broader audience. Our goal is to build a service that supports sharing of sensitive data. We now discuss the requirements for building services for researchers to share sensitive research data. The regulations, re-identification attack methods, and legal and ethical requirements may change over time. We conjecture that the following requirements are essential for building a service compliant with the legal, regulatory, and ethical requirements stipulated by concerned authorities/stakeholders. And can adapt to changes that may affect data sensitivity and a subject's preference. Previous works such as Datatags [18] and A-DAM [21] have already provided methods for computable ethical, legal, and regulatory requirements. We include the *data sovereignty* requirement for sensitive data such as medical records, which restricts the movement of data outside a nation's physical boundaries, even if the data are hosted, by a cloud service provider (CSP).

**RQ 1.** [Timely Dissemination of Data Policies] Data policies define the ethical, legal, and regulatory requirements attached with a dataset. The service should disseminate changes to data policies within a predefined time $\tau$. Each dissemination should be secure, maintain integrity, and be logged for auditing. Consent revocations and new approvals from an IRB or a REC can result in such changes. A change in laws, regulations, and institutional guidelines may result in a policy change as well.

**RQ 2** (Data Sovereignty). The service should ensure data sovereignty by verifying data residency and sovereignty. Any attempt of data movement which violates data sovereignty should be prevented and logged.
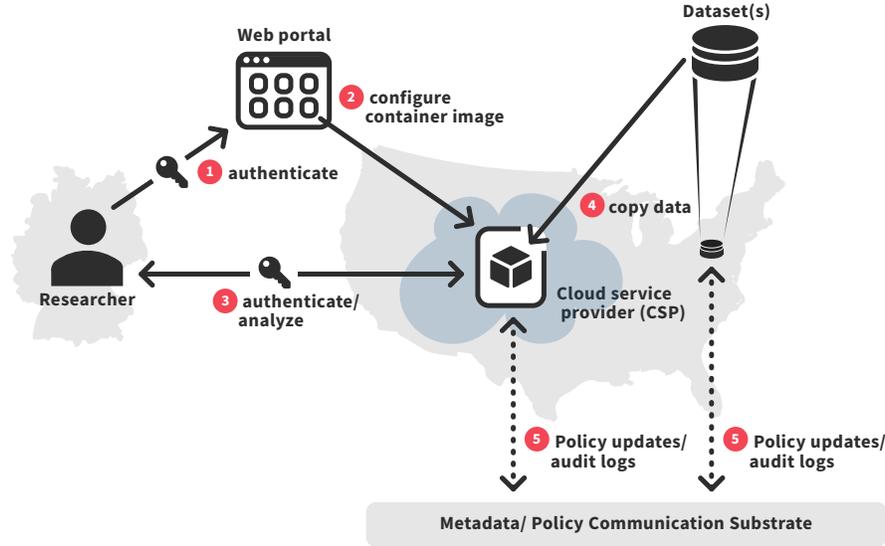
**Fig. 1.** An example workflow: A researcher from Germany wants to access a sensitive dataset from an institute in the USA. The movement of the data outside the USA is restricted. By leveraging a CSP, we facilitate sharing of the dataset for research without moving the data outside the USA. The accesses are logged for audits and transparency.

**RQ 3** (Garbage Collection). Once completed with the task, copies of data should be securely deleted and logged. No residual copies of the data or the dataset remain on an unsolicited location or machine.

**RQ 4** (Auditing). The service should log each operation and action in a distributed log. These operations and actions must be available for auditing by an IRB, REC or an independent auditing authority.

**RQ 5** (Secure Computation). The continuous access evaluation should be securely computed. The attack surface for any tampering of data access policies should be limited at the end/user system.

## 4 System Overview

We now describe our approach and discuss how different components will fulfill the requirements discussed earlier (Section 3). We assume that each dataset has a unique identifier. Researchers who are interested in accessing a dataset can authenticate themselves. The dataset in our example, has been marked with data sovereignty constraints. A CSP has a data center in the same region as the dataset.

### 4.1   Workflow

Fig. 1 shows the system architecture with an example workflow. A researcher or simply a *user*, is interested in analyzing a dataset hosted at an institution in another country. After authenticating herself using the web portal, the user configures a machine for her analysis (RQ4). The user can decide from multiple pre-configured container images which contain different data analysis software packages. These container images enable communication with a trusted substrate for exchanging data policy updates and logs (RQ1, RQ4). The daemon software is pre-installed and configured in these container images. For enhanced security, the container images are hardened to limit copying data out of the machine and allow only a set of pre-approved packages. As the last step, the user chooses the dataset that she is interested in (RQ5).

Once configured, a new container instance with chosen software packages runs in the cloud, which exists in the same country as the dataset (RQ2). Upon initialization, the user needs to authenticate herself again for obtaining a copy of the dataset (RQ4). The instance also receives policy changes that might arrive while the user is working on the dataset i.e. after the initialization (RQ1).

The instance enforces *use-based* policies using Intel's software guard extensions (SGX) using approaches like [3] (RQ5). In-line monitoring using SGX can result in a performance penalty. For performance reasons, we use the delegated monitoring architecture proposed in Birrel et al [3]'s work. The instance communicates with the metadata communication substrate and keeps the metadata/policy up-to-date (RQ1). Any changes to the checked-out dataset are disseminated through the substrate and received by both, the original dataset and the copy. After receiving metadata updates, the compliance and accesses are evaluated again (RQ5). The daemon process routinely checks for compliance with the latest data security policies. These work behind the scenes and the user is notified if additional inputs are required for compliance. Thus, making compliance easier for the user. In case of non-compliance, the user may loose access to the machine while saving the image to save her work in-progress. After resolving the non-compliance issue, the user may access the machine again. The user is only allowed to export results in different approved file formats to the web portal (RQ3) to limit data leaks. Through the portal, the user can obtain the results later. The results can be archived at the portal for cross-examination by auditors and reviewers.

At the end of the analysis, the user can terminate the instance and the analysis scripts and the dataset copy are securely destroyed (RQ3). A user may also choose to save the current state of the container for reproducibility of results [19]. The sharing and accesses generate logs containing sanitized information for audits. The stakeholders (subjects, REC or IRB) can also view reports on data use and sharing, and intervene if necessary. Oversight committees (RECs or IRBs) can review non-compliance incidents and take necessary actions. These actions can be in the form of policy updates, propagated to every copy of the dataset.

**Limitations** Our approach is designed against a benign threat model. The system may not protect against a sophisticated attacker. The availability of a cloud service provider (CSP)'s container services in the same administrative region as of a dataset's location, is crucial for the data sovereignty requirements.

## 5    Conclusion

Regulatory compliance in research data sharing is a developing problem with newly introduced regulations and growing concerns about individual privacy. The relationship between subjects and research institutions relies heavily on trust for voluntary participation. Data sharing and use, compliant with the subjects' wishes is crucial for continued participation and sustaining trust. We discussed the requirements for building a service enabling compliant data use and sharing sensitive research data. We further presented our approach for building such a service and how it addresses those requirements. Together with our partners from sports sciences and medical science, we plan to evaluate our system with legal and regulatory requirements for sensitive research data.

# Bibliography

[1] Axelsson, A.S., Schroeder, R.: Making it open and keeping it safe: E-enabled data-sharing in sweden. Acta sociologica **52**(3), 213–226 (2009)

[2] Bartling, S., Friesike, S.: Opening science: The evolving guide on how the internet is changing research, collaboration and scholarly publishing. Springer Nature (2014)

[3] Birrell, E., Gjerdrum, A., van Renesse, R., Johansen, H., Johansen, D., Schneider, F.B.: SGX enforcement of use-based privacy. In: Proceedings of the 2018 Workshop on Privacy in the Electronic Society, pp. 155–167 (2018)

[4] Borgman, C.L.: Open data, grey data, and stewardship: Universities at the privacy frontier. Berkeley Tech. LJ **33**, 365 (2018)

[5] Braunschweig, K., Eberius, J., Thiele, M., Lehner, W.: The state of open data. Limits of current open data platforms (2012)

[6] Hammack, C.M., Brelsford, K.M., Beskow, L.M.: Thought leader perspectives on participant protections in precision medicine research. Journal of Law, Medicine & Ethics **47**(1), 134–148 (2019)

[7] King, G.: An introduction to the dataverse network as an infrastructure for data sharing (2007)

[8] Kroll, J.A., Kohli, N., Laskowski, P.: Privacy and policy in polystores: a data management research agenda pp. 68–81 (2019)

[9] Maguire, S., Friedberg, J., Nguyen, M.H.C., Haynes, P.: A metadata-based architecture for user-centered data accountability. Electronic Markets **25**(2), 155–160 (2015)

[10] Molloy, J.C.: The open knowledge foundation: open data means better science. PLoS biology **9**(12) (2011)

[11] Mulligan, A., Mabe, M.: The effect of the internet on researcher motivations, behaviour and attitudes. Journal of Documentation (2011)

[12] Pampel, H., Vierkant, P., Scholze, F., Bertelmann, R., Kindling, M., Klump, J., Goebelbecker, H.J., Gundlach, J., Schirmbacher, P., Dierolf, U.: Making research data repositories visible: the re3data. org registry. PloS one **8**(11), e78080 (2013)

[13] Peterson, Z.N., Gondree, M., Beverly, R.: A position paper on data sovereignty: The importance of geolocating data in the cloud (2011)

[14] Schneider, G.: Disentangling health data networks: a critical analysis of Articles 9 (2) and 89 GDPR. International Data Privacy Law **9**(4), 253–271 (2019)

[15] Sharma, A., Nilsen, T.B., Brenna, L., Johansen, D., Johansen, H.D.: Accountable Human Subject Research Data Processing using Lohpi. In: Proceedings of the ICTeSSH 2021 conference (2021), https://doi.org/10.21428/7a45813f.80ebd922

[16] Sharma, A., Nilsen, T.B., Czerwinska, K.P., Onitiu, D., Brenna, L., Johansen, D., Johansen, H.D.: Up-to-the-minute Privacy Policies via gossips in Participatory Epidemiological Studies. Frontiers in Big Data **4** (2021)

[17] Slegers, C., Zion, D., Glass, D., Kelsall, H., Fritschi, L., Brown, N., Loff, B.: Why Do People Participate in Epidemiological Research? Journal of Bioethical Inquiry **12**(2), 227–237 (2015), https://doi.org/10.1007/s11673-015-9611-2

[18] Sweeney, L., Crosas, M., Bar-Sinai, M.: Sharing sensitive data with confidence: The datatags system. Technology Science (2015)

[19] Trisovic, A., Durbin, P., Schlatter, T., Durand, G., Barbosa, S., Brooke, D., Crosas, M.: Advancing computational reproducibility in the Dataverse data repository platform. In: Proceedings of the 3rd International Workshop on Practical Reproducible Evaluation of Computer Systems, pp. 15–20 (2020)

[20] Winthrop, S., Roth, I., Baynes, G.: Social considerations to make data FAIR-er: Understanding researchers' views on data "misuse" and credit. Septentrio Conference Series (1) (2019)

[21] Woolley, J.P., Kirby, E., Leslie, J., Jeanson, F., Cabili, M.N., Rushton, G., Hazard, J.G., Ladas, V., Veal, C.D., Gibson, S.J., et al.: Responsible sharing of biomedical data and biospecimens via the "Automatable Discovery and Access Matrix"(ADA-M). NPJ genomic medicine **3**(1), 1–6 (2018)