




# Identifying Personal Data Processing for Code Review

Feiyang Tang<sup>1</sup><sup>a</sup>, Bjarte M. Østvold<sup>1</sup><sup>b</sup> and Magiel Bruntink<sup>2</sup><sup>c</sup>

<sup>1</sup>Norwegian Computing Center, Oslo, Norway

<sup>2</sup>Software Improvement Group, Amsterdam, The Netherlands  
{feiyang, bjarte}@nr.no, m.bruntink@sig.eu

**Keywords:** Data Privacy Protection, Code Review, Static Analysis.

**Abstract:** Code review is a critical step in the software development life cycle, which assesses and boosts the code’s effectiveness and correctness, pinpoints security issues, and raises its quality by adhering to best practices. Due to the increased need for personal data protection motivated by legislation, code reviewers need to understand where personal data is located in software systems and how it is handled. Although most recent work on code review focuses on security vulnerabilities, privacy-related techniques are not easy for code reviewers to implement, making their inclusion in the code review process challenging. In this paper, we present ongoing work on a new approach to identifying personal data processing, enabling developers and code reviewers in drafting privacy analyses and complying with regulations such as the General Data Protection Regulation (GDPR).

## 1 INTRODUCTION


The General Data Protection Regulation (GDPR) lays the legal foundation for data protection in the EU and increases individual data protection rights throughout Europe. It also carries significant fines of up to 4% of yearly worldwide revenue for businesses that do not comply with the legislation. Many IT system providers, especially software-producing firms, may need to alter their systems in order to comply with the GDPR. This is predicted to require significant effort (Blume, 2016). As a result, providing software engineers in the industry with effective and systematic ways to build data protection into software is an essential and beneficial study topic (Lenhard et al., 2017). Organizations are pushing security to the software development life cycle, such as code review, to prevent software security vulnerabilities (Braz and Bacchelli, 2022). Similarly, to comply with privacy-by-design and perform privacy analysis tasks, code reviewers would benefit from similar tools to those used for security to identify privacy-related patterns in software.


Developers address privacy concerns using data security terminology, and this vocabulary confines their notions of privacy to threats outside of the orga-


nization (Hadar et al., 2018). However, even though data security is the main prerequisite of data privacy, privacy protection in software is still very much different from traditional security-related vulnerabilities. And according to Bambauer: “security and privacy can and should be treated as distinct concerns” (Bambauer, 2013). Developers struggle to convert legal, ethical, and social privacy concerns into concrete technology and solutions (Notario et al., 2015).

Assessing privacy involves not only finding personal data in the software but also evaluating compliance with the related processing. GDPR defines as processing: “any operation or set of operations which is performed on personal data or on sets of personal data, whether or not by automated means.” The definition encompasses a vast range of actions performed on personal data, such as collecting, recording, organization, structuring, storage, adaption or modification, retrieval, transit, etc. Privacy assessment tasks beg the question: How can we assist code reviewers and software developers in assessing personal data processing? By identifying personal data and the relevant processing in the system, code reviewers can uncover interesting patterns and utilize them to redesign the system to be more privacy-friendly or perform privacy analysis.

In this paper, we present ongoing work on a novel approach designed to assist developers and code reviewers in identifying personal data processing,

<sup>a</sup> <https://orcid.org/0000-0002-8720-6743>

<sup>b</sup> <https://orcid.org/0000-0001-6922-4027>

<sup>c</sup> <https://orcid.org/0000-0002-6117-6347>

which can subsequently be used for privacy analysis. This enables developers and code reviewers to assist organizations with a variety of important privacy-related tasks, such as completing a data protection impact assessment (DPIA) and creating a privacy policy.

## 2 RELATED WORK

An essential step in the software development process, code reviewing incorporates both manual and/or automated reviews. The main goal of code reviews is to assess and boost the code's effectiveness and correctness, pinpoint security issues, and raise its quality by adhering to best practices (McIntosh et al., 2014). To automatically evaluate code, a variety of vulnerability detection tools have been built. They are also known as source code analyzers or static analysis tools, as they can analyze a program's code without having to execute it (McGraw, 2008).

CodeQL<sup>1</sup>, and Semgrep (r2c, 2022)<sup>2</sup> are two popular code review tools that utilize static analysis. CodeQL treats code as if it were data, and issues are modeled as queries. Following the extraction of these queries from the code, they are executed against a database. The database is a directory containing data, a source reference for displaying query results, query results, and log files. Semgrep matches grammatical patterns on parsed programs (represented as an Abstract Syntax Tree (AST)) instead of matching string or regular expression (regex) patterns on the program as a string. Semgrep makes it considerably simpler to construct customized rules than CodeQL, which needs rules to be defined in QL, a declarative object-oriented query language.

There is relatively little published work that focuses on code reviews to identify privacy-related vulnerabilities, and it is problematic to translate current security knowledge to privacy, which we will explain in Section 3. There are studies on the identification of personal data that are valuable to our research. Fugkeaw et al. (Fugkeaw et al., 2021) proposed AP2I to enable organizations to identify and manage personal data in the local file system automatically. By monitoring network traffic, ReCon (Ren et al., 2016) utilized machine learning to identify probable personal data breaches. van der Plas et al. (van der Plas, 2022) used CodeBERT, a RoBERT-like transformer model, to identify personal data in Git commits.

---

<sup>1</sup><https://codeql.github.com/>

<sup>2</sup><https://semgrep.dev/>

## 3 BACKGROUND AND CHALLENGES

Data privacy analysis is becoming as crucial as security vulnerability discovery and has brought a new dimension to the data security dilemma (Bertino, 2016). It is advantageous for code reviewers to be able to conduct a similar privacy analysis that they did for security.

The current state of the art is mostly focused on security analysis. Although data security is a primary requirement for data privacy, the analysis domain and identification process are rather different (Jain et al., 2016). Simply adopting security mechanisms and mindsets to analyze privacy can be misguided, and even harmful (Bambauer, 2013).

Integration of recent studies on assessing software privacy during code review is challenging. On the subject of program analysis, three well-known privacy analysis methods are available. First, static analysis based on bytecode requires project compilation, whereas dynamic taint analysis requires project execution. This is not practical nor efficient for code reviewers to implement. A machine learning-based technique is similarly difficult to implement, as it requires a large and diverse training data set. Obtaining and generating such data sets requires additional effort and could be outside the scope of code reviewers' capabilities. Lastly, text analysis based on UI widgets is constrained for privacy by domain-specific UI attributes. A financial web application that employs a model trained on an Android health mobile application is unlikely to benefit. Code reviewers require an approach that is simple to deploy, efficient, and adaptable (Buse and Zimmermann, 2012).

Due to the complex nature of privacy and the fluidity of the definition of personal data, identifying the processing of personal data in the codebase presents challenges.

In the following paragraphs, we highlight the two most significant challenges related to the task in the context of code review.

### 3.1 The Ambiguous Definition of Personal Data

Article 4(1) in GDPR defines personal data as:

*any information relating to an identified or identifiable natural person ('data subject'); an identifiable natural person is one who can be identified, directly or indirectly, in particular by reference to an identifier such as a name, an identification number, location data,*

*an online identifier or to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural or social identity of that natural person.*

The definition of personal data in the GDPR is so broad that almost any information may qualify as personal data if it refers to a specific individual, such as the fact that a person is wearing a red shirt (Berčič and George, 2009). The definition is also semantically ambiguous.

In contrast to the fact that certain data may be anonymous from the start (such as weather sensor data without any connection to real people), other data may initially be personal data but later be successfully altered to no longer have any connection to an identified or identifiable natural person. This emphasizes how flexible the categorization of personal data is (Finck and Pallas, 2020).

The same data point may be personal or non-personal depending on the context and may thus be covered by the regulation or not. This implies that the categories of personal data in the software vary depending on the software and the processing underlying it. For instance, health data such as blood pressure and medical records, for example, are sensitive for a health application, but location data is sensitive for navigation software.

Even if we accept that content-wise every item of information can be considered personal data if it can be related to an individual, the GDPR’s definition is still rather vague structurally since it is not always clear what kind of structure every ‘record’ of an individual must have to be considered personal data (Voss and Houser, 2019).

Due to the ambiguous nature of the definition of personal data in the relevant legislation, it is practically difficult for us to have a clear and fixed identifier to precisely locate personal data in code.

### 3.2 What Counts as Sensitive Processing?

Data subjects may agree to data processing for particular reasons. This is the usual legal basis but only counts as one factor. Processing may also be “necessary for the performance of a contract to which the data subject is a party or in order to take steps at the request of the data subject prior to entering into a contract.” (Voss and Houser, 2019)

Unfortunately, concerns that arise in principle about the relationship between contract and consent tend to be avoided in reality by disregarding consent requirements (Pormeister, 2017).

We cannot rely on existing privacy policies and written consent to uncover personal data processing in the codebase. This requires us to consider all potential personal data processing in the codebase. Later we will explain how we define and identify the relevant processing in software in Section 4.1.2.

## 4 APPROACH

We present an approach to identify instances of personal data processing in the codebase and present them in a way that facilitates the code review.

The approach has three primary phases: pattern matching, labeling, and grouping of results. As input, we take the codebase, which consists of source code files. Then, a static analyzer will evaluate these source code files using our rules and patterns. The code snippets discovered by the static analyzer are then labeled according to the various features they include. Finally, we allow users to group the results by single or several labels, allowing a personalized exploration of the findings.

An illustration of our approach is shown in Figure 1.

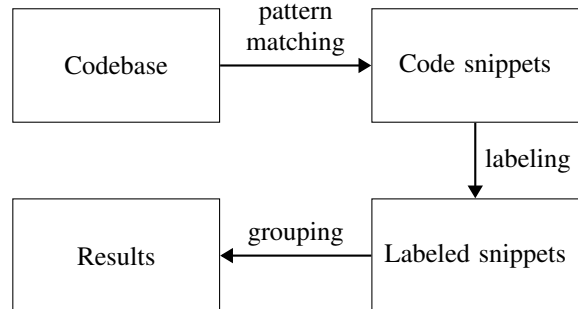


Figure 1: Approach.

### 4.1 Design Choices

In the following paragraphs, we discuss our design choices for implementing the approach.

#### 4.1.1 Types of Findings

We want to have a basic default list of personal data that we want to locate, this is mostly personal identification and characteristics data, such as full name, email address, gender, sexual orientation, and age. We call them fixed personal data.

According to different types of software, we customize default lists for them. For example, for a banking/finance application, the list may contain bank account numbers, credit scores, and salary information.

This type of personal data is subject to context - the types of processing in specific software, which we named contextual personal data.

Depending on how we locate the mentioned personal data in the software, we can divide their occurrences in code into simply two types.

The first is in clear text. This includes all kinds of locations where personal data appear in clear text. It is verbatim or direct personal data. For example, a credit card number appears in an SQL query, or an email address falls into a log function.

The other type is more common and subtle, where personal data is stored in a variable or an object. Depending on the different types of programming languages, the object types might vary from a local variable, a class instance, or a prototype. This means we aim to find the code that processes this type of data.

#### 4.1.2 Types of Processing

Simply locating every instance of personal data produces a large number of results. Many of these do not directly help the code reviewer's work, which is to find meaningful processing. We want to use a hybrid approach to cover as many as processing as possible.

Processing personal data represents a specific behavior. This motivates our first approach: to use an action name tag to find relevant processing. We adopted most of the verbs from Section 3 of DPV (Pandit et al., 2019)<sup>3</sup>. These vocabularies help us to find relevant processes in the software.

The second approach is the identification of external libraries. We know that modern applications rely on various APIs to achieve different goals. Therefore, obtaining a list of relevant APIs and detecting the existence of personal data that flows into them helps us find meaningful patterns.

## 4.2 Pattern Matching

The first step is to feed our codebase (consisting of source code files) to the static analyzer for pattern matching. We chose Semgrep as our analyzer because of its user-friendly rules and rapid processing performance. Depending on the different syntactic characteristics of personal data, as we discussed in Section 4.1.1, we adopt a hybrid approach that combines two different types of analysis.

- Match personal data in clear text using regular expression matching.
- Taint analysis to find flows in each file between a source (where personal data enters the analysis

scope) and a sink (where personal data gets processed) that match our criteria.

Our personal data processing rules currently support Java, JavaScript, and TypeScript as our primary analysis domains. However, our rules for identifying clear-text personal data apply to the vast majority of Semgrep-supported languages.

#### 4.2.1 Source and Sink

Our prototype classifies the sources into nine separate categories. As stated in Section 4.1.1, we divide fixed personal data into four different categories: *account*, *contact*, *national ID*, and *personal ID*. Included are five more contextual personal data categories, such as *location*, *health*, and *financial* data. In addition, we provide a template for identifying the processing of personal data and enable code reviewers and developers to submit additional personal data simply by entering the relevant keywords. Then, corresponding rules will be automatically produced for future use.

Sinks are categorized into five main types. Three types of action: *data manipulation (M)*, *data transportation (T)*, and *data creation/deletion (C/D)*. Another two represent two special types: *database (DB)* and *encryption (E)*.

A sink's name may contain a specific type of source. For example, `setLatitude(100,100)` does not take any source into the method, but includes a source identifier `Latitude` and a sink identifier `set`, showing that it processes values directly as a source into a sink. We call this special type of sink a source-specific sink. When a source-specific sink invokes anything, we mark this source-specific sink as the new source but the caller of the source-specific sink as the new sink. For example, in `gpsTracker.setLatitude(100,100)`, `setLatitude` becomes the new source and `gpsTracker` is the new sink.

Inspired by how Privado<sup>4</sup> uses regular expressions to identify GDPR-related data in Java applications, a sample Semgrep rule that matches the pattern of *account* data source goes into a *transportation (T)* sink is shown below in Figure 2, followed by a sample code snippet detected in Figure 5.

## 4.3 Labeling

The identified findings from Semgrep are in the form of various lengths of code snippets (consisting of statements and expressions). Each finding contains at least one detected sink and one source (or an object that received value from a source). We abstract the

<sup>3</sup><https://w3c.github.io/dpv/dpv/>

<sup>4</sup><https://www.privado.ai>

```

1 rules:
2   - id: account-data-transportation
3     languages:
4       - javascript
5       - java
6       - typescript
7     mode: taint
8     message: Match found
9     pattern-sinks:
10      - patterns:
11        - pattern: $SINK(..., $Z, ...)
12        - metavariable-regex:
13          metavariable: $SINK
14          regex: (?i)(.*(send|move|connect|escap|stream|redirect|
15            erase|query|share|stor|transfer|transmit|move).*)
16      pattern-sources:
17        - pattern-regex: (?i).*(?:account|user|customer|doctor|patient|
18          policyholder|insurer|claimant)[^\s/();|,!=>]{0,3}(id|number|no|
19            num)
20        - pattern-regex: (?i)(?:facebook|twitter|instagram|linkedin|
21          pinterest|behance|dribbble)[^\s/();|,!=>]{0,2}(?:id|account|
22            username|handle)
23        - pattern-regex: (?i).*(?:db|database|jira|sql|postgres|mongo|aws)
24          [^\s/();|,!=>]{0,3}(psw|pswd|password|passwd)
25        - pattern-regex: (?i)(.*(?:db|database|jira|sql|postgres|mongo|
26          aws)[^\s/();|,!=>]{0,3}user[^\s/();|,!=>]{0,3}name)|(.*(account|
27            customer|doctor|patient|teacher|student|person|organization|
28            company)[^\s/();|,!=>]{0,3}name)
29      severity: WARNING

```

Figure 2: Semgrep rule: find personal data flows from account data source to transportation sink.

```

1 this.userService.updateUser(newUser.id,
2                               {defaultOrganizationId:
3                               newUser.organizationId})
4                               .catch((error) => {
5     console.error('Error while updating default
6                   organization id', error);
7   });

```

Figure 3: Sample code snippet (from ToolJet) detected by Semgrep showing a flow from account personal data to a transportation sink.

structure of possible sources and sinks in each code snippet using the symbols below.

- $O$  ranges over sources
- $I$  ranges over sinks
- $I^O$  ranges over source-specific sinks

We write  $\bar{O}$  as shorthand for a possibly empty sequence  $O_1, \dots, O_n$ . Here the underscore  $_$  represents a placeholder for an expression that is insignificant in terms of privacy - it is neither a source nor sink nor contains a value from a source.

Below is a list of the common flow abstracts between sources and sinks that we observed in each code snippet. Each abstract represents a typical flow, for example, ① to ③ show that there are values passing through a sink to a source, from a non-privacy sensitive value (①) or from another source (②) or from innovating a sink inside another source object (③).

- |                          |                       |
|--------------------------|-----------------------|
| ① $O = \_I(-)$           | ⑥ $\_O.I(-)$          |
| ② $O_2 = \_I(O_1, -)$    | ⑦ $\_O.I(-, \bar{O})$ |
| ③ $O_2 = \_O_1.I(-)$     | ⑧ $\_I^O(-)$          |
| ④ $\_ = \_O.I(-)$        | ⑨ $\_I^O(-, \bar{O})$ |
| ⑤ $\_ = \_I(\bar{O}, -)$ | ⑩ $\_I(\bar{O}, -)$   |

For each identified code snippet, we label them with 22 labels (9 types of source, 5 types of sink, 5 types of source-specific sink, and 3 types of change in the sensitivity level), which are listed in Table 1. Besides the definition of source and sinks, we also introduce an important label: sensitivity. The sensitivity level can increase, decrease, and stay the same in one identified code snippet.

Table 1: Labels to be assigned to each code snippet.

- |       |  |
|-------|--|
| $O$   | Nine types of source: $\{O^1, O^2, \dots, O^9\}$                           |
| $I$   | Five types of sink: $\{I^1, I^2, \dots, I^5\}$                             |
| $I^O$ | Five types of source-specific sink: $\{I^{O^1}, I^{O^2}, \dots, I^{O^5}\}$ |
| $S$   | Sensitivity level change: {up, down, equal}                                |

**Sensitivity Level.** Not every result shares the same level of sensitivity regarding personal data processing. After processing, the data from the source might remain at a similar sensitivity level, become more sensitive, or become less sensitive.

- $S = \text{up}$ : ①, ④, ⑤
- $S = \text{equal}$ : ②, ③, ⑥, ⑦, ⑧, ⑨
- $S = \text{down}$ : ⑩

## 4.4 Result Presentation

Johnson et al. (Johnson et al., 2013) pointed out that “because the results are dumped onto a code reviewer’s screen with no distinct structure causing him to spend a lot of time trying to figure out what needs to be done”. This indicates that developers and code reviewers may not benefit from ungrouped code snippets from static analysis tools if they are not presented in a sensible manner.

To tackle this issue, we present a two-phase technique to process the findings from Semgrep and present them to code reviewers in a smart way.

After each code snippet is labeled, we start to group them for presentation using their labels and other criteria. Criteria for grouping include not only the labels but also other properties:

- neighboring results will be combined (same file and within a line number threshold);

- same or similar source/sink name;
- same API usage (e.g. every code snippet that is related to the same API MongoDB).

Figure 4 following provides a straightforward illustration of how we present our results. The results are presented in two separate sections: plain text results and flow results. Users have the flexibility to select any label or label combination to filter the results.

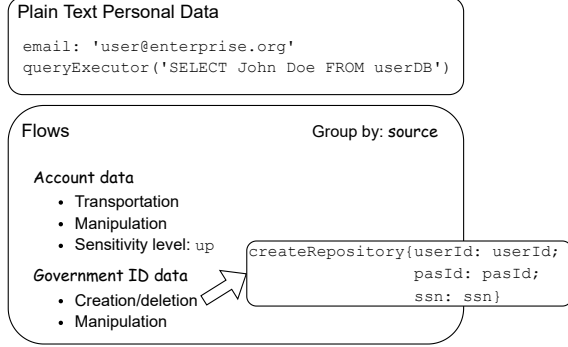


Figure 4: Example presentation of the result. *Personal data occurrences* is at the top and *personal data processing code* is at the bottom.

## 5 DEMONSTRATION

We created rules in Semgrep trying to capture as many useful findings for our analysis. The software we analyzed here is ToolJet<sup>5</sup>, an open-source low-code framework for building React-based web applications. ToolJet’s implementation is mostly in JavaScript and TypeScript. Users can build internal tools using ToolJet’s prebuilt UI widgets to connect to data sources like databases, API endpoints, and external services. This means ToolJet has many parts that process personal data, which makes it a good starting example.

Our Semgrep rules produce a total of 1,589 results from ToolJet’s source code. We manually reviewed each of the results and calculated the precision for each category. If a single result can clearly demonstrate the processing of personal data, we consider it relevant and it could be beneficial for privacy code review. Surprisingly, most false positives come from the personal data occurrence detector (with a precision of only 46.6%), while most personal data processing results are relevant (with an average of 90.9% precision for categories that have more than 50 code snippets identified).

Detailed statistics are listed in Tables 2 and Table 3.

Table 2: The code snippet count for each identified source and sink identified, ‘-’ marks labels for which our approach detected no code snippet. Sink types are: *data manipulation (M)*, *data transportation (T)*, *data creation/deletion (C/D)*, *database (DB)*, *encryption (E)* and *log (L)*.

	M	T	C/D	DB	E	L
Account	66	171	84	24	-	21
Contact	89	175	36	3	-	3
Personal ID	56	133	41	7	1	4
Online ID	6	26	1	-	-	1
Location	1	2	-	-	-	-

Table 3: The precision of code snippet relevance (in %) for each identified type of source and sink, ‘-’ marks the labels for which our approach did not detect any code snippet, ‘\*’ marks the labels for which our approach detected less than 10 results. Sink types are: *data manipulation (M)*, *data transportation (T)*, *data creation/deletion (C/D)*, *database (DB)*, *encryption (E)* and *log (L)*.

	M	T	C/D	DB	E	L
Account	90.9	90.6	95.2	91.67	-	95.2
Contact	89.9	94.9	80.6	*	-	*
Personal ID	92.9	81.9	85.4	*	*	*
Online ID	*	84.6	*	-	-	*
Location	*	*	-	-	-	-

Figure 5 shows a simple interesting example of a grouped result showing how personal data `userId` is retrieved from a local repository in `app_users.service.ts` and then utilized to generate many data structures, such as the app object in `app_service.ts`.

```

1  async create(user: User): Promise<App> {
2    const app = await this.appsRepository.save(
3      this.appsRepository.create({
4        name: 'Untitled app',
5        createdAt: new Date(),
6        updatedAt: new Date(),
7        organizationId: user.organizationId,
8        userId: user.id,
9      })
10   );
    app_service.ts

```

```

1  async create(user: User, appId: string, organizationUserId: string,
2    role: string): Promise<AppUser> {
3    const organizationUser = await this.organizationUsersRepository.
4      findOne({ where: { id: organizationUserId } });
5
6    return await this.appUsersRepository.save(
7      this.appUsersRepository.create({
8        appId,
9        userId: organizationUser.userId,
10       role,
11       createdAt: new Date(),
12       updatedAt: new Date(),
13     })
14   );
15 }
    app_users.service.ts

```

Figure 5: Grouped example results showing how `organizationUserId` flows between functions.

<sup>5</sup><https://github.com/ToolJet/ToolJet>

## 5.1 Future Work

Since our objective is to identify all relevant processing of personal data in source code, reducing false negatives is our next primary priority. However, in our case, false positives are not a major concern. Due to the subtlety of personal data processing, determining relevance without human assistance is particularly challenging. Specifying the analysis to certain specific patterns would ease manual analysis. This necessitates the implementation of a privacy taxonomy. Using Ethyca's taxonomy (Ethyca, 2022) as an example, we may modify our labels to match the technique with the taxonomy.

As an extension of this article, we propose an automated mapping of personal data in an unpublished (under review) manuscript (Tang et al., 2023) to assist developers and code reviewers in identifying privacy-related code. The mapping based on static analysis automatically detects personal data and the code that processes it, and we offer semantics of personal data flows.

## 6 CONCLUSIONS

This short paper presented ongoing work on a novel, customizable approach to identify personal data processing for code review. This three-phase technique first uses Semgrep to match patterns in the code based on rules for sources and sinks, then associates code snippets generated from pattern matching with a set of behavioral labels, and finally groups results to reduce code reviewer workload. Our demonstration shows the utility and feasibility of this method for gathering and presenting code snippets related to personal data processing from a codebase.

Along with the continued development of the approach architecture (refined rules for source and sink, more meaningful labels, and additional criteria for grouping), future work will focus on expanding the case study to include a larger set of open-source software from various domains and conducting a thorough user evaluation of the resulting platform.

## ACKNOWLEDGEMENTS

This work is part of the Privacy Matters (PriMa) project. The PriMa project has received funding from European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 860315.

## REFERENCES

- Bambauer, D. E. (2013). Privacy versus security. *J. Crim. L. & Criminology*, 103:667.
- Berčić, B. and George, C. (2009). Identifying personal data using relational database design principles. *International Journal of Law and Information Technology*, 17(3):233–251.
- Bertino, E. (2016). Data security and privacy: Concepts, approaches, and research directions. In *2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC)*, volume 1, pages 400–407. IEEE.
- Blume, P. (2016). Impact of the EU General Data Protection Regulation on the public sector. *Journal of Data Protection & Privacy*, 1(1):53–63.
- Braz, L. and Bacchelli, A. (2022). Software security during modern code review: The developer's perspective. *arXiv preprint arXiv:2208.04261*.
- Buse, R. P. and Zimmermann, T. (2012). Information needs for software development analytics. In *2012 34th International Conference on Software Engineering (ICSE)*, pages 987–996. IEEE.
- Ethyca (2022). Fides language. <https://ethyca.github.io/fideslang/>. (Accessed on 11/15/2022).
- Finck, M. and Pallas, F. (2020). They who must not be identified—distinguishing personal from non-personal data under the GDPR. *International Data Privacy Law*, 10(1):11–36.
- Fugkeaw, S., Chaturasrivilai, A., Tasungnoen, P., and Techaudomthaworn, W. (2021). AP2I: Adaptive PII scanning and consent discovery system. In *2021 13th International Conference on Knowledge and Smart Technology (KST)*, pages 231–236. IEEE.
- Hadar, I., Hasson, T., Ayalon, O., Toch, E., Birnhack, M., Sherman, S., and Balissa, A. (2018). Privacy by designers: software developers' privacy mindset. *Empirical Software Engineering*, 23(1):259–289.
- Jain, P., Gyanchandani, M., and Khare, N. (2016). Big data privacy: a technological perspective and review. *Journal of Big Data*, 3(1):1–25.
- Johnson, B., Song, Y., Murphy-Hill, E., and Bowdidge, R. (2013). Why don't software developers use static analysis tools to find bugs? In *2013 35th International Conference on Software Engineering (ICSE)*, pages 672–681. IEEE.
- Lenhard, J., Fritsch, L., and Herold, S. (2017). A literature study on privacy patterns research. In *2017 43rd Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*, pages 194–201. IEEE.
- McGraw, G. (2008). Automated code review tools for security. *Computer*, 41(12):108–111.
- McIntosh, S., Kamei, Y., Adams, B., and Hassan, A. E. (2014). The impact of code review coverage and code review participation on software quality: A case study of the qt, vtk, and itk projects. In *Proceedings of the 11th working conference on mining software repositories*, pages 192–201.
- Notario, N., Crespo, A., Martín, Y.-S., Del Alamo, J. M., Le Métayer, D., Antignac, T., Kung, A., Kroener, I., and Wright, D. (2015). PRIPARE: integrating privacy

- best practices into a privacy engineering methodology. In *2015 IEEE Security and Privacy Workshops*, pages 151–158. IEEE.
- Pandit, H. J., Polleres, A., Bos, B., Brennan, R., Bruegger, B., Ekaputra, F. J., Fernández, J. D., Hamed, R. G., Kiesling, E., Lizar, M., et al. (2019). Creating a vocabulary for data privacy. In *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*, pages 714–730. Springer.
- Pormeister, K. (2017). Informed consent to sensitive personal data processing for the performance of digital consumer contracts on the example of “23andMe”. *Journal of European Consumer and Market Law*, 6(1).
- r2c (2022). Semgrep. <https://semgrep.dev/>. (Accessed on 11/15/2022).
- Ren, J., Rao, A., Lindorfer, M., Legout, A., and Choffnes, D. (2016). Recon: Revealing and controlling pii leaks in mobile network traffic. In *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*, pages 361–374.
- Tang, F., Østvold, B. M., and Bruntink, M. (2023). Mapping personal data in source code for GDPR compliance.
- van der Plas, N. (2022). Detecting PII in Git commits. *TU Delft Master’s thesis*.
- Voss, W. G. and Houser, K. A. (2019). Personal data and the GDPR: providing a competitive advantage for US companies. *American Business Law Journal*, 56(2):287–344.