

Remarks



Silesian
University
of Technology

SILESIA N UNIVERSITY OF TECHNOLOGY
FACULTY OF AUTOMATIC CONTROL, ELECTRONICS
AND COMPUTER SCIENCE

PROGRAMME: INFORMATICS

Final Project

Data anonymisation web platform

author: Szymon Pluta

supervisor: Krzysztof Simiński, PhD DSc

Gliwice, January 2022

Contents

Abstract	1
1 Introduction	3
2 Problem analysis	5
2.1 Data explosion	5
2.1.1 Technology advancement	5
2.1.2 Data analytics	6
2.1.3 Big data	6
2.2 Data privacy	8
2.2.1 Authorization to share data	9
2.2.2 General Data Protection Regulation	10
2.3 Data anonymisation	12
2.3.1 Background	12
2.3.2 Definition	13
2.3.3 Data classification	13
2.3.4 Utility vs privacy trade-off	14
2.4 Data masking techniques	14
2.4.1 Suppression	15
2.4.2 Generalisation	17
2.4.3 Perturbation	18
2.4.4 Pattern masking	21
2.4.5 Hashing	22
2.4.6 Randomisation	23
2.4.7 Artificial data	24

2.4.8	Shortening	26
2.5	Existing solutions	26
2.5.1	Open-source software	27
2.5.2	Enterprise-class software	28
2.6	Anonymisation as a platform	29
3	Requirements and tools	31
3.1	Requirements engineering	31
3.1.1	Functional requirements – domain	31
3.1.2	Functional requirements – infrastructure	38
3.1.3	Non-functional requirements	47
3.2	System engineering	49
3.2.1	Use cases	49
3.2.2	REST API	49
3.3	Software tools and technologies	53
3.3.1	Tools	53
3.3.2	Server	56
3.3.3	Client	59
3.4	Development methodology	61
4	External specification	63
4.1	Domain specific glossary	63
4.2	Installation	64
4.2.1	Environments — overview	64
4.2.2	Cloud setup	65
4.2.3	Semi-cloud setup	66
4.2.4	Local setup	67
4.3	Activation	68
4.4	Types of users	68
4.5	User manual	69
4.6	System administration	69
4.6.1	Environments	69
4.6.2	Database	71

4.6.3	Files	72
4.6.4	Mail service	72
4.6.5	Unprotected resources	73
4.6.6	Authorization	73
4.6.7	User account	74
4.6.8	Scheduler	74
4.7	Security issues	75
4.7.1	User account	75
4.7.2	Continuous integration	75
4.7.3	Containers isolation	76
4.7.4	Validation	76
4.8	Working scenarios	76
5	Internal specification	79
6	Verification and validation	81
7	Conclusions	83
	Bibliography	XII
	Index of abbreviations and symbols	XV
	List of additional files in electronic submission (if applicable)	XVII
	List of figures	XIX
	List of tables	XXI

Abstract

1. coraz więcej danych produkowanych 2. które są wykorzystywane przez firmy
3. co sprawia że rośnie istota prywatności danych 4. więc potrzebne są coraz lepsze
techniki ochrony danych 5. anonimizacja jest taką metodą ochrony danych 6. więc
ją zaimplementuje w sposób ultra generyczny 7. jako nowoczesna web platforma
8. dzięki czemu będzie można ją wykorzystać jako wszechstronne narzędzie dla
biznesu i R&D

Keywords: 2-5 keywords, separated by commas

Keywords: data anonymisation, privacy and data protection, GDPR compliance, platform as a service, web

Chapter 1

Introduction

Technological advancement being observed in the past years fundamentally changed the relevance of data in today's digitalized world. Information became an innovation stimulus in the area of research and development. The quantity of data that organizations produce, process, store and share is at a continuous growth. An enormous amount of 1.8 zettabytes ($1.8 \cdot 10^{21}$ bytes) of new data was produced only in the 2011, and every two consecutive years this number is doubling [73]. After decades of observed technological advancement and innovation, the global internet traffic finally entered the zettabyte-era, as it had reached a magnitude of one zettabyte in 2016, and in the calendar month being as early as September [16].

The vast quantities of processed information allowed for brand new research fields such as data science or big data analytics to form, which are used by organizations to derive new insights in a way that was previously impossible. Organizations collect and process the data to enhance the services they provide to the customers through statistical analysis or newly developed computer science processes including data mining and machine learning. The utility of delivered services is increased at a lower cost and improved efficiency through the insights extracted from the collected information about how the services are consumed [8].

Następne do zrobienia, w uproszczeniu i w podanej kolejności:

- Jak widać, dane są używane wszędzie. Przeciętna osoba sobie nawet nie zdaje sprawy.
- Istota prywatności danych, ochrona danych osobowych.

- Prawo wolno reaguje na uregulowanie ochrony danych osobowych: GDPR dopiero w 2018. Prawo różnie działa w różnych krajach.
- W czasach narastających danych, dane muszą być zanonimizowane, bo...
- ...ale:

Wniosek: generyczny sposób anonimizacji - nie ma takiego.

- Mogą istnieć dobre algorytmy anonimizacji, ale są optymalne jedynie w określonych kontekstach
- **Cel pracy:** stworzenie generycznego rozwiązania do anonimizacji baz danych w dowolnej formie danych (nieważne co dane reprezentują).
- Kontroler danych zna swoje dane i sam dostosowuje optymalny sposób anonimizacji.
- Rozwiązanie może być dostarczone jako free-access (research), B2C lub nawet B2B.
- Zakres pracy, opis rozdziałów - na samym końcu, po zrobieniu innych rozdziałów.

- introduction into the problem domain
- settling of the problem in the domain
- objective of the thesis
- scope of the thesis
- short description of chapters
- clear description of contribution of the thesis's author – in case of more authors table with enumeration of contribution of authors

Chapter 2

Problem analysis

An extensive analysis of data anonymisation subject is required to be conducted to comprehend the growing data protection needs of the society and to engineer a solution that will properly suit those needs.

2.1 Data explosion

2.1.1 Technology advancement

The continuous and rapid exponential growth of data being collected globally is further excited by improvements to the overall population's accessibility to digital technology [10]. Cisco Systems estimates within its annual report [15] that 66 percent of the world population will have access to the web by 2023, compared to 51 percent in 2018, whereas the number of devices that are connected to the web will reach a staggering value of three times as many as the entire population size – demonstrating a total of 60 percent expansion when compared to 2018. Even the area of mobile connection, which was established long ago, is still sustaining growth – by 2023, mobile connectivity will be a privilege for 70 percent of the world's population, compared to 66 percent in 2018. The global average mobile network speeds will be tripled through a rapid increase from 13.2 Mbps in 2018 to 43.9 Mbps in 2023.

2.1.2 Data analytics

The raw representation form of the data is not interpretable until it is put under a context and processed into practical information. Acquiring relevant insights and conclusions from the information can be achieved through a wide term of analytics, which encompasses the actions needed to be performed to produce new information, including analysis, filtering, processing, grouping, and contextualizing the information. Newly discovered knowledge is inferred from the produced information. Apart from the processes, analytics also includes the technologies, methodologies, and algorithms to use and could be divided into descriptive analytics, diagnostic analytics, prescriptive analytics, and predictive analytics [10].

2.1.3 Big data

Big data analytics deals with the difficulties of managing the observed exponentially increasing collected volumes of data. Its purpose is not only to handle the processing and analysis of the data through specialized software tools and frameworks but also to handle the means on how this enormous amount of data is collected and stored in the first place. It is in its nature that big data is all about massive volumes of information that require specialized hardware infrastructure to store it [10].

Services of enterprise organizations are running on all the collected data which can take various forms such as database entries, metrics, logs, or outgoing messages. New data streaming technologies working at a large scale needed to be engineered to handle the continuous flow of data between systems and databases. An example of such technology includes Apache Kafka which generates even more than a trillion of messages per day for individual large enterprise organizations taking advantage of it [35, 54].

This only proves that big data deals not only with massive volumes – it has also to deal with the high velocities of data generation, which is yet another characteristic of data [10]. According to DOMO report published back in 2020, 90% of the world's data was generated just in the preceding two years, and on average every person in the world created 1.7 megabytes of data per second – which yields 2.5 quintillions ($2.5 \cdot 10^{18}$) bytes of new data each day [28].

The momentum of the immense big data interest growth among organizations is not fading away yet, as more and more new businesses and researchers are drawn to this subject. The benefits of big data especially concern scientific organizations and large enterprises of which the financial domain and IT industry are the common consumers [30]. Organizations find interest in information analytics for remarkably diversified reasons. It is recognized as a field that will entirely alter all parts of civilization such as businesses or society as a whole [46].

Applications

Various types of organizations collect data to take advantage of the insights derived out of the data. The big data analytics applications impact can be observed already today in a broad spectrum of domains.

Leading technology companies, such as Google and Facebook, to name a few, sell anonymised collected user data access to their partner advertisers [46]. This is legally possible as the information that was anonymised, i.e., de-identified in a way that it is no longer bound to an individual, may flow from one system to a system.

The Large Hadron Collider (LHC) located in CERN, being the largest physical experiment, annually produces approximately 30 petabytes of data. LHC takes advantage of light sensors that monitor the collisions of hundreds of millions of particles accelerated nearly to the speed of light. The collisions create an enormous amount of data to be processed by computer algorithms in the hope of discovering new particles, e.g., a Nobel prize-awarding discovery of the Higgs boson had taken place in 2012 [1].

Enterprise stores such as Amazon or SAP Commerce Cloud collect information regarding the way how the visitors browse and interact with these stores. Collected information may involve behavioural data related to customer engagement, such as the pages we visit, event clicks, or the way we scroll the page. The insights derived from the collected information enable making future improvements of these services – for example by improving the digital marketing or performance improvements based on the metrics [41]. The customer experience is also improved as based on the collected data the advertisements or item recommendations can be tailored to the

specific user's preferences. The recommendation engine may also attempt to match your profile data to people of similar profiles to provide better recommendations. Services attempt to analyse the behavioural patterns such as time of day we browse the store or what circumstances caused our last visit to finish. Even the details such as the exact neighbourhood location we live in, combined with its estimated wealth, organizations may attempt to guess our potential income level [46]. These data analytics are performed to improve the possibility of customers buying yet another item.

2.2 Data privacy

Privacy endangerment is inherent to big data and it is its major drawback. The personal data we continuously give away to third parties is the big data fuel.

As technology evolves, concerns relating to the privacy of our personal information should also grow – and for a relevant reason. We tend not to give a second thought to whom the data is shared, how it may be used, and in what kind of circumstances. We don't wonder how our data may be exploited to alter our thinking, decisions, or even ideas – whether in an ethical manner or not.

Although the use of our personal information existed ever since the very first census was created, data privacy is a relatively new concept, as it did not exist prior to the global adoption of the internet. Granted that our data was used by the researchers even before the digitalized era of the internet, the motives for that usage were not commercial [67]. Nowadays, the data has most certainly become an asset – a resource like any others, and a rather precious one.

It is argued that data privacy should be centered only around data usage that has the potential to be a privacy breach. On the other hand, it is argued that merely a collection of the data is already privacy harm. The information that was collected is endangered by many threats, including data misuse, data breach, data leak, or even authorities access without legal obligations. Anonymisation is the best method to mitigate conflicts raised by big data with respect to data privacy and data protection [71].

Luckily the law had finally caught up to the circumstances of the increasing data usage and the associated risks. New European regulations were adopted in

2018 in the form of the General Data Protection Regulation (GDPR) to protect our data privacy in a refined fashion. The data subjects, i.e., the individuals represented by the data [8], now have better control over their personal information – we are now entitled to know what information concerning us is being processed and for what purpose. We are also entitled to withdraw at any time the consent for the processing of our data. In case of violations, we have the right to complain to authorities and seek justice against both the data controllers, i.e., the entities that determine the intention and means of processing the personal information, and the data processors, i.e., the entities that process the personal information on behalf of the data controllers [74].

Overall awareness of the data privacy significance had improved in the society, and the means to achieve data privacy through data protection had also improved as organizations needed to adapt to the new situation by applying enhanced measures to their protection of data – anonymisation and pseudonymisation being the notable examples of such measures.

2.2.1 Authorization to share data

Majority of data privacy regulations are based on a consent of an individual, i.e., it is lawful to process and use the information for secondary purposes only if an individual explicitly acknowledge their consent for that [77]. This may appear easier said than done due to the unobvious difficulties data controllers face when trying to obtain such a broad authorization consent that will take into account all possible secondary purpose usages.

Consider a patient entering medical facility for an ordinary appointment. The patient would likely find it unusual, disturbing or even shocking if upon his entrance to the facility he was to receive an overwhelming form that included dozens of independent consent authorization requests. The consents could give the impression of being seemingly unrelated to his visit in the first place, e.g., a consent to share the data with researchers of an university located on another continent. In the end that could destroy the data subject's trust – in this case patient's trust.

This theoretical scenario may not easily be implemented in the real world counterpart, as it could be even impossible to know or predict all possible secondary

purpose usages in the first place, and consent-based authorization is all about knowing the usages.

Consider a newly discovered purpose to process personal information of an already existing database. Getting consent after the data had already been collected, i.e., backward in time, would be impossible to accomplish as the data controller would need to contact potentially hundreds of thousands of people for their explicit consent. New purposes can be discovered years after the data collection.

Having that in mind, no consent is required when processing the data that is already anonymised. Data that was stripped from personal identifying or identifiable information data can be used in any way and can be shared with third parties without previously agreed consent. Data controllers now face a realistic to solve the problem of information anonymisation rather than an unrealistic problem of consents collection [29].

It is worth mentioning that there exist cases which are defined under Article 6 of General Data Protection Regulation (GDPR) [22] when consent is not required to process the data, e.g., if the processing is required to defend the data subject's interests or in order not to break the compliance with legal obligations as a data controller. GDPR is a new European law that replaces the preceding Data Protection Directive regulation adopted by European Community in 1995 [77].

2.2.2 General Data Protection Regulation

One of the primary objectives of GDPR is personal data privacy protection which is a fundamental right and freedom of people as defined under the Recital 1 of GDPR [18] and the Charter of Fundamental Rights of the European Union [17]. Newly discovered challenges for the protection of personal data arising from the ongoing globalization and quick development of digital technologies. This in turn vastly had increased both the scope of the gathering of the data and the sharing of thereof. General Data Protection Regulation (GDPR) is a data protection law that came into force on May 25, 2018, to addresses these data privacy-related issues in a strict manner [20].

Compliance with GDPR law is critical for organizations in given the significant administrative fines they face. Violations of the data processor and data controller

obligations defined in GPDR are subject to costly penalties that are imposed by European authorities. Non-compliance with technical rules implies a penalty of 2 percent of the total annual turnover of the previous year, or €10 million, whichever one is higher, whereas non-compliance with basic data protection principles imply an even higher penalty of 4 percent of the total annual turnover of the previous year, or €20 million, whichever one is higher [21][50].

Implementation of this law had immediately increased the significance of data anonymisation as an information sanitization process in today's world [55]. Anonymisation is a specific form of data masking that suddenly became more relevant in today's world for the reason that the strict regulations, and therefore administrative fines, defined in GDPR do not apply to the anonymised information. Data protection principles covered throughout GDPR concern only the processing of information that is already identified to a natural person, or that is identifiable to a natural person, i.e., an individual is yet to be identified. Given the fact that anonymised information is by definition not relating to a person, hence it can be completely exempted from falling under GDPR requirements, which apply only to personal data, as stated under Recital 26 [19]:

The principles of data protection should apply to any information concerning an identified or identifiable natural person. [...] The principles of data protection should therefore not apply to anonymous information, namely information which does not relate to an identified or identifiable natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable. This Regulation does not therefore concern the processing of such anonymous information, including for statistical or research purposes.

GDPR distinguishes personal data, anonymised data, and pseudonymised data as distinct variations of data. The information that had gone merely through the pseudonymisation process would still fall under the regulations of GDPR, due to the existing relevant possibility of future re-identification of the data subject, whereas in the case of the anonymised data, such re-identification is by definition either impossible or extremely impractical, and the anonymisation is irreversible by definition. Anonymised data is completely exempted from being governed by

GDPR. Nevertheless, pseudonymisation is still one of many possibilities for the data controllers and data processors to be GDPR compliant [77]. The point of anonymisation in the context of GDPR is to be completely exempted from being governed by this regulation.

2.3 Data anonymisation

The data released by organizations exclude identity-related information such as names, addresses, telephone numbers, or credit card numbers. Personal data is stripped from disseminated data sets through anonymisation to protect the anonymity of the individuals, i.e., data subjects [68].

2.3.1 Background

Disseminating the medical data volumes is crucial for the evolution of the world's healthcare services. Consider a collection of medical data concerning patients' clinical information. Medical researchers and doctors take an advantage of the collected data sets to improve their comprehension of diseases and explore new possibilities to treat these diseases, and hence both the overall capability to treat the diseases and the general efficiency of health services are improved. At last, it is the patients who benefit from the research conducted on their data since the services they are offered continuously improve. Nevertheless, it is known that medical data is exceptionally sensitive by its nature due to the details that include e.g., patient data, laboratory tests results, diagnosis details, prescribed medications, and history of diseases [38].

Having understood how sensitive by its essence is the patient information and the vital needs to share this data, this is where data anonymisation plays an indeed crucial role. It would be impossible to disseminate patient information without prior anonymisation of thereof.

2.3.2 Definition

Anonymisation is a statistical disclosure control method of particular importance. The ultimate anonymisation goal is to de-identify the data by pruning the personal information in such a way that the relationship between the data subject and corresponding records is blurred.

The data anonymisation is considered to be an effective one only if both of the following criteria hold true [40]:

- performed anonymisation operations are irreversible,
- data subject re-identification is impossible or impractical.

2.3.3 Data classification

The prerequisite for anonymisation is to understand the context of what does the data represent and by whom – in addition to how – will it be processed [29].

The commonly adopted approach when attempting to anonymise the data is to first classify the data attributes as direct identifiers (i.e., personal identifying confidential attributes unique to an individual, such as address or tax identification number), indirect identifiers (i.e., so-called quasi-identifiers that when combined can discover the individual's identity – such as sex or date of birth) and non-identifiers [27, 29].

It is context-dependent to determine which attributes are quasi-identifiers – and finally which techniques need to be combined to achieve anonymisation. Virtually, even the most unexpected attributes could be quasi-identifiers. Ignoring such fields when preparing for anonymisation may eventually lead to future re-identification of the data subjects, as successfully demonstrated, e.g., when breaking the anonymity of the Netflix Prize dataset back in 2006 [7]. Since then significant advancements in the field of de-identification were yielded – and in various domains [9].

On the other hand, too much anonymisation of quasi-identifiers strips the data out of utility.

2.3.4 Utility vs privacy trade-off

The organizations need to compromise between the utility and privacy when releasing the data. As anonymisation level increases when more and more restrictive techniques are applied, the utility of the same data decreases.

Consider the two theoretical boundary cases of either releasing data in its original form which maximizes the utility at the cost of discarding the protection of the data subjects' privacy altogether, and the second case of not releasing any data at all which completely preserves the privacy [80].

Clearly the data must not be released in its original state – strict GDPR that imposes large administrative fines [21] exists to mitigate such violations. On the other hand, data that completely suppresses all the information has no value. The data controllers therefore need to strike a balance between the two cases. An appropriate anonymisation strategy maximizing the utility without endangering the privacy must be adopted.

For this reason the data controllers must select appropriate data masking techniques.

2.4 Data masking techniques

The analysis of the domain concludes that the designed system needs to be enriched with diverse masking techniques to meet the fundamental requirement of offering a generic context flexible anonymisation tool. This is achievable on the account of the breadth of the available data masking techniques being genuinely remarkable. The methods that are supported by the anonymisation platform include suppression, generalisation, perturbation, pattern masking, hashing, attribute randomisation, record randomisation, substitution, tokenisation, shortening, or random number generation. All of the available methods are subject to customisations.

An individual method can be classified either as a non-perturbative method if it reduces the data detail by removing some of the information (e.g., suppression, generalisation, pattern masking or shortening) or as a perturbative method if it alters the data by creating some level of uncertainty (e.g. by adding noise through

perturbation or randomisation) [38]. Perturbative techniques are the best for maintaining a high analytical value of the data as they preserve the data truthfulness, while non-perturbative methods are the best for protecting privacy as they prune the information.

It is not possible to predetermine – in a general fashion and without an extra context – the data masking techniques that need to be combined to achieve actual anonymisation of the previously identifiable information. Instead, a context consisting of:

- exact representation form of the data being processed (*What is the data?*)
- data processors who use the data (*Who will use the data?*)
- processing purpose, e.g., research objective (*How will the data be used?*)

is always required when considering an effective way to anonymise data under that context [29].

2.4.1 Suppression

Suppression is the strongest anonymisation technique which completely removes all the values associated with the given attribute, hence rendering the complete protection of the data and therefore enforcing the best data subjects privacy. It is guaranteed that no further implication attacks can be performed against this data [68]. Nevertheless, suppression being the strongest privacy-preserving method implies the highest (i.e., worst) data utility loss as an inseparable consequence – the ultimate cost of complete anonymisation that is to be paid.

Consider the raw data of the three attributes presented in Tab. 2.1 to be put under the process of suppression.

Exemplary results of performed suppression are depicted in Tab. 2.2. All of the three attributes values were suppressed to the value of the suppression token supplied by the data controller.

Phone number suppression deserves a brief discussion. Consider a context in which an attacker successfully obtained the information of exactly one individual data subject of his interest. Consider that the attacker decided to call the anonymised phone number. This number could be suppressed in a way to specifically

Table 2.1: Suppression – input data.

Sex	PIN codes	Phone number
F	3248	1212 – 345345
F	8090	4000 – 303030
M	1337	5191 – 915100

Table 2.2: Suppression – masked data.

Sex	PIN codes	Phone number
F/M	####	3000 – 123123
F/M	####	3000 – 123123
F/M	####	3000 – 123123

protect against this attack and the attacker could call someone intentionally most definitely did not, i.e., data protection officer or authorities. In this scenario, it could be trivial to trace the call and find the potential attacker. As demonstrated, suppression could be used in unobvious ways.

Suppression is typically applied on the column-level of an attribute, however, it is also applicable to the individual records. When compared to plain generalisation, suppression yields a higher information loss and can be thought of as a particular case of generalisation [38].

The designed software offers column level suppression with a user-specified value (i.e., suppression token). The values are transformed to this value using a simple algorithm shown in Fig. 2.1.

```

1 procedure suppression(values, token)
2   for (value in values)
3     value := token
4   return values

```

Figure 2.1: Pseudocode – suppression.

Table 2.3: Generalisation – input data.

Age	Salary	Location
27	36 000	Poland
52	54 000	Canada
30	180 000	Poland
68	128 000	Switzerland

Table 2.4: Generalised – masked data.

Age	Salary	Location
26 – 30	1 – 60000	Europe
51 – 55	1 – 60000	North America
26 – 30	120001 – 180000	Europe
66 – 70	120001 – 180000	Europe

2.4.2 Generalisation

Generalisation is the second most common non-perturbative anonymisation technique [27]. This technique processes the raw data by aggregating and substituting the initial values of a given attribute to the values that are more general [6]. The process of generalisation constitutes a conversion of any value to a more general scope.

The proposed software solution includes two generalisation strategies, namely:

- based on the size of distributions,
- based on the number of distributions.

Consider the raw data of the three attributes presented in Tab. 2.3 to be put under the process of generalisation.

Results shown in Tab. 2.4 include exemplary possible values after being processed by the generalisation method. As depicted in the table, all three attributes have been grouped into broader value ranges, hence undergoing the generalisation, yet every column was generalised with a different generalisation strategy.

The age attribute was generalised with a strategy based on the size of distribution. In this case, the size of distribution was defined as 5, as every interval

aggregates five distinct increasing integer values.

On the other hand, salary data was generalised with a strategy based on the number of distributions – the values were aggregated to three even distributions.

Finally, the generalisation does not necessarily concern only numerical values as observed in the generalisation of location attribute. To better understand this type of generalisation, consider that the age attribute could be generalised into the values of e.g., *young adult*, *adult*, *senior*. This type of generalisation is not supported by the developed software, however, an illustration of the existence of such a method is necessary. The fundamental assumption of the developed software is to be generic, i.e., data context-agnostic, and this generalisation method is too specific in its essence to be implemented generically.

Both supported generalisation strategies were further enhanced with customizations. The data processor is allowed to predefine minimum and maximum boundary values for the generalised intervals. The computed starting value for the distribution concerning the lowest interval is either minimum raw value or user predefined minimum value, whichever is lower. The maximal value for the highest interval is computed in an analogical way.

The generalisation strategies can be summarized with the pseudocodes shown in Fig. 2.2 and Fig. 2.3.

Generalisation may be applied on a global level or local level [68]. In terms of database-related vocabulary, as this paper primarily concerns an anonymisation of databases, we say generalisation on a column or record level, respectively.

2.4.3 Perturbation

Perturbation is a data masking technique that encompasses the process of swapping the original values with artificial ones in a way that the statistical factors of the data are similar to the initial data. Perturbation is typically achieved by adding noise to the information so that the values are slightly different [34].

Consider weight and height data shown in Tab. 2.5.

Tab. 2.6 shows the data after it was perturbed using the two different perturbation methods that are supported by the engineered software.

The supported methods include adding the noise based on the strategies of:

```

1 procedure generalise(values, configuration)
2   computedMin := min(values.min(), configuration.userMin())
3   computedMax := max(values.max(), configuration.userMax())
4   distributionSize := configuration.distributionSize()
5
6   // Phase 1: Generate empty intervals.
7   intervals := <Interval, Values>
8   i := computedMin
9   while (i < computedMax)
10    interval := asEmptyInterval(i, i + distributionSize)
11    insert interval into intervals
12    i += distributionSize
13
14   // Phase 2: Populate intervals.
15   for (value in values)
16     for (interval in intervals)
17       if (value is within interval)
18         interval := get interval from intervals
19         add value to interval
20   return intervals

```

Figure 2.2: Pseudocode – generalisation based on distribution size.

Table 2.5: Perturbation – input data.

Height	Weight
166	58
170	66
194	91

- fixed value,
- percentage value.

Height attribute values could have been perturbed using the former (i.e., fixed noise of ± 3), whereas weight attribute values could have been perturbed using the latter (i.e., percentage value of $\pm 5\%$). In either case, the linkage between individual records is blurred with this technique.

```

1 procedure generalise(values, configuration)
2   computedMin := min(values.min(), configuration.userMin())
3   computedMax := max(values.max(), configuration.userMax())
4   numberOfDistributions = configuration.numberOfDistributions()
5
6   // Phase 1: Generate empty intervals.
7   intervals := <Interval, Values>
8   i := computedMin
9
10  distributionSize := (computedMax - computedMin) / numberOfDistributions
11  j := 0
12  while (j < numberOfDistributions)
13    interval := asEmptyInterval(i, i + distributionSize)
14    insert interval into intervals
15    i += distributionSize
16    j += 1
17
18  // Phase 2: Populate intervals.
19  for (value in values)
20    for (interval in intervals)
21      if (value is within interval)
22        interval := get interval from intervals
23        add value to interval
24  return intervals

```

Figure 2.3: Pseudocode – generalisation based on number of distributions.

```

1 procedure perturbation(values, noise)
2   for (value in values)
3     random := pick random from [value - noise, value + noise] interval
4     value := random
5     value := fit value within boundaries
6   return values

```

Figure 2.4: Pseudocode – perturbation based on fixed noise.

Similarly to generalisation, both perturbation techniques allow the user to explicitly specify the minimum and maximum accepted boundary values.

Perturbation technique works the best with continuous values – it is an effective method for de-identification of quasi-identifiers such as dates and numbers [69]. Scientifically advanced perturbation approaches are further divided into linear and non-linear perturbation models [64].

Table 2.6: Perturbation – masked data.

Height	Weight
169	57
168	67
193	91

```

1 procedure perturbation(values, noise)
2   for (value in values)
3     coeff := pick random from [1 - noise, 1 + noise] interval
4     value := value * coeff
5     value := fit value within boundaries
6   return values

```

Figure 2.5: Pseudocode – perturbation based on percentage noise.

2.4.4 Pattern masking

Pattern masking is a data distortion technique designed to be used when only some parts of the data should be masked. This technique is typically used to mask e.g., codes of various forms, phone numbers, credit card numbers, and other structured data.

Consider the data shown in Tab. 2.7 to be masked using the patterns presented in Tab. 2.8.

Exemplary masked output is visible in Tab. 2.9.

Analysis of pattern masking possibilities has yielded the need to create a versatile pattern configuration. This technique is therefore highly configurable, i.e., available pattern tokens include: *O* – preserve the character, *X* - mask the character, *U* – mask with random uppercase letter, *L* – mask with random lowercase

Table 2.7: Pattern masking – input data.

PIN code	Software version	Product code
54850185	2.7.1	BAR/service/1
03013844	2.4.0-rc.3	FOO/service/7
76590209	1.0.1-alpha	QUX/utility/0

Table 2.8: Pattern masking – exemplary patterns.

PIN code	Software version	Product code
OOXXXXXO	OOXOX	UUUOOOOOOOOON

Table 2.9: Pattern masking – masked data.

PIN code	Software version	Product code
54#####5	2.0.0	RAN/service/0
03#####4	2.7.1	DOM/service/7
76#####9	0.1.9	IZE/utility/1

letter, N – mask with random digit, A – mask with random alphabetic character and C – mask with random alphanumeric character.

Users can specify the masking character, e.g., $\#$, and can also decide whether the result is truncated to the pattern length (as shown in the masking of software version in Tab. 2.9).

2.4.5 Hashing

Hashing is a deterministic masking technique [55] that transforms data to the fixed-length output. This technique is best used for unstructured data. Available implementations include SHA2 and SHA3.

Consider the server logs depicted in Tab. 2.10. Exemplary hashed values are located in Tab. 2.11.

Although a mere inspection of the generated hash value does not immediately trace back to the original value, an attack to de-anonymise the data can be easily performed when the attacker knows what the hashed data represents, e.g. a struc-

Table 2.10: Hashing – input data.

Server logs
185.184.2.198 — 200 POST: /api/v1/auth/refresh-token
185.184.2.198 — 200 POST: /api/v1/worksheets
255.7.141.233 — 200 POST: /api/v1/outcomes/generate

Table 2.11: Hashing – masked data.

Server logs
b27ffd54e5b05a538f333157363f18df0a2aaae5754dfd9ec9daad9cc4ccd7a2 477784538ed600c38f586079a7d5e99aac4af97d1cb322888de54edeb600b14d 2cd3e1912285c765f1746d5b68b1fdbbfff6be9460e305acc18a1d9d777d89b5e

Table 2.12: Randomisation – attribute – input data.

Identity	Virus
John	Influenza A
Marc	Pneumonia
Stephen	Bronchitis

tured birth date [72]. Like in the case of all techniques, the decision to use this technique is context-dependent.

2.4.6 Randomisation

Randomisation methods offer the best middle ground between maintaining the data utility and preserving privacy [65]. This method can be performed on a column or a row level.

Column shuffle

Column shuffle is an exceptionally effective method for breaking strong links of data belonging to individual data subjects. All of that is achieved while preserving the data utility [67], hence this technique is great for maintaining the analytical value of the data [65].

Consider the medical data from Tab. 2.12

An exemplary randomised data is shown in Tab. 2.13.

This technique can be used in two ways:

- shuffle without repetitions which preserves the data distribution,
- shuffle with repetitions which distorts the data distribution.

Table 2.13: Randomisation – attribute – masked data.

Identity	Virus
Stephen	Pneumonia
Marc	Influenza A
John	Bronchitis

Table 2.14: Randomisation – record – input data.

Hex triplet	Set of decisions
FF00FF	1101
54E7CD	1010
E5E5E5	0000

Row shuffle

The data may also be needed to be shuffled on a record level – similarly, with or without preserving the distribution.

Consider shuffling the characters of an RGB colour described in hexadecimal format. Furthermore, consider a concept of storing a set of abstract decisions (e.g., willingness to participate in the election voting) as a sequence of bits where each bit represents an individual boolean decision. Tab. 2.14 shows the data.

The table 2.15 shows the data after the masking process. The colours were processed without preserving the original characters distribution.

2.4.7 Artificial data

Replacing the original data with an artificial data can be achieved through substitution, tokenisation or random number generation.

Table 2.15: Randomisation – record – masked data.

Hex triplet	Set of decisions
0000F0	0111
E57DC4	1100
5E5555	0000

Figure 2.6: Substitution.

(a) Input data.		(b) Masked data.	
Name	Surname	Name	Surname
Jan	Gold	Lucius	Lucci
Bob	Ng	Decimus	Rector
Bob	Xi	Decimus	Lucci
Maria	Robin	Amanda	Rector

Substitution

Substitution is of particular relevance when the masked data needs to be realistic [61]. Consider the data from Tab. 2.6a undergoing the substitution process. Consider that the data controller provides the following data for the name parameter: *Lucius*, *Decimus*, *Amanda*, and the following data for the surname parameter: *Lucci*, *Rector*. The table 2.6b presents an exemplary possible masked result. The name attribute is masked without memorylessness, and the masking occurs in a circular manner. The method is designed to be used with or without memorylessness in mind. Duplicate entries of the same unmasked data values may always be altered to the same masked data values. This technique is also a great extension point for data controller customisations, as the external data set needs to be provided, therefore the data controller can have full control over the data.

Tokenisation

In comparison to the substitution technique, this operation is relevant when the masked data does not need to be realistic [61], i.e., it could be useful to perform the analytics on the frequencies of data duplicates.

Consider the data from a survey question shown in Tab. 2.7a.

The tokenised data is shown in Tab. 2.7b. The *agree* value was tokenised to the same token value hence conserving the frequency characteristics of the data.

Figure 2.7: Tokenisation.

(a) Input data.	(b) Masked data.
Survey response	Survey response
Agree	1
Not sure	2
Agree	1
Strongly disagree	3

Figure 2.8: Random number.

(a) Input data.	(b) Masked data.
Customer satisfaction	Customer satisfaction
2	5
1	5
3	4

Random number

The fundamental requirement of anonymisation platform is to be a versatile software to anonymise the data, hence even a simple random number technique needs to be supported to possibly cover a broader range of the business needs.

Consider the data of customer satisfaction described as a value ranging from 1 to 5 shown in Tab. 2.8a. The values from Tab. 2.8b show the new masked values.

2.4.8 Shortening

Consider the surnames from Tab. 2.9a to be shortened to a length of 5 characters. The shortened data is visible in Tab. 2.9b. It is possible to specify whether or not each data record should be terminated with a dot.

2.5 Existing solutions

All of the tools have a shared goal of being GDPR compliant in mind. As this section shows, the existing solutions are greatly similar to the engineered software

Figure 2.9: Shortening.

(a) Input data.	(b) Masked data.
Surname	Surname
Kowalski	Kowal.
Kowalewski	Kowal.
Nowak	Nowak

and often offer *less* versatility than the designed anonymisation web platform.

Of course, as data privacy concerns are increasing in society, the need for data protection also increase. The relevance of anonymisation will only increase in the future. The market is just not steady yet.

2.5.1 Open-source software

Anonimatron

Anonimatron [63] is a Java tool designed for anonymisation of databases specifically for development purposes. The idea behind the tool is that the software developers may not be authorized to access the production data, yet may be required to reproduce, e.g., data-related performance problems. Internally, anonimatron uses the concept of synonyms.

Anonymizer

Similarly to Anonimatron – Anonymizer [25] is a tool that anonymises the databases specifically for development. This tool was written in Ruby.

ARX

ARX [59] is a Java based desktop solution offering a significant variety of advanced data masking techniques [60]. The tool offers the functionality of analysing the usefulness of the generated outcomes.

Amnesia

Amnesia [56] is a software system that allows uploading the data, configuring the anonymisation parameters, and downloading anonymised outcomes. However, it supports merely the text files.

2.5.2 Enterprise-class software

SAP Data Services

SAP Data Services is a data management software that internally uses Data Mask component of the Transformation module to provide the masking functionalities [66]. Data Mask component offers highly configurable functionalities like perturbation, generalisation, and pattern masking.

Oracle Data Masking and Subsetting

This is yet another enterprise-class data masking software, however, it is solely focused on masking. This software is optimised for Oracle database, however other databases are also partially supported. The available operations include shuffling, randomisation, pattern masking [23].

Microsoft Azure ecosystem

The anonymisation techniques that are used throughout Azure Synapse Analytics, Azure SQL Managed Instance, or Azure SQL Database include pattern masking, random number generation, and text substitution [2].

TurboCat.io

This tool primarily oscillates between the encryption method to protect sensitive data [45].

2.6 Anonymisation as a platform

The people who were software engineers back in the 1990s remember how the web had changed everything. The market market was rapidly dominated and the way business systems are now engineered [47].

The companies started to massively migrate their services to a cloud with the quite recent rise of containerization and microservice architecture. Trending software concepts such as serverless or event-storming and tools such as Kubernetes — all of them are oscillating around the web, which is the only clear and technologically available choice to build a new system. All this technological evolvement is enabled by the open-source which is also continuously growing together with the web-based software systems, in a feedback loop manner [12]. The start of the decline in the web's growth is yet to be discovered and likely will not emerge in the foreseeable future.

With that being said, the designed software system is categorized and delivered as Platform as a Service (PaaS) and runs in a cloud. The fundamental and the most significant objective of the system is to design an innovative platform that will offer anonymisation capabilities through a broad variety of highly configurable anonymisation techniques in such a way that the data controllers can fit their unique business needs easily.

Chapter 3

Requirements and tools

3.1 Requirements engineering

3.1.1 Functional requirements – domain

Data masking

There should exist data transforming services that will offer the functionalities of the following techniques:

- suppression
- generalisation
- perturbation
- randomisation
- tokenisation
- random number replacement
- hashing
- column shuffle
- row shuffle

- shortening
- pattern masking

There should be a validation standard established and enforced that will determine if a particular anonymisation technique is compatible with another technique. The techniques should be configurable.

Starting template generation

The user should be able to start the generation of a new template based on the database dump file he provides. The acceptable dump files should include a compressed archive or plain script file. The user must successfully upload the dump before starting the template generation.

Template generation progress

The template generation should involve the sequence of three steps: persisting the dump file, restoring a new database from the dump, extracting metadata about the restored database. The user should be able to view the template generation progress.

This should be presented as a progress panel containing a visualization of three steps. Completion of a step either marks it as a success or an error. Currently processing step should be shown. Errors informing about e.g. invalid databases must be reported to the user. The generation progress should internally update status of the template and should be performed as a sequence of events.

Restore database

The system should be able to restore a new database from a dump file of either a compressed archive or script file.

Mirror database

The system should be able to create a mirror of an existing user database. The mirror should be used for writing the anonymisation script into it and then dumping it into the final outcome result.

Extract metadata

The system should extract the metadata from the uploaded database when creating a new template. The metadata must include information concerning the structure of the database, e.g., the tables, columns, and primary keys the database contains. Extracted metadata should be persisted in the server database as JSON type.

Inspecting metadata

The user should be able to inspect the extracted metadata.

Download metadata

The user should be able to download the extracted metadata.

Viewing templates

The user should be able to view his templates.

Producing worksheet

The user should be able to produce a new worksheet from the templates he created earlier on. The worksheet should be used for building up the anonymisation setup. The worksheet must not be accessible to other users.

Viewing worksheets

There should be a table giving an overview of the created worksheets. Only worksheets belonging to the user should be shown. The table should embed buttons to go to the summary view or outcome generation view.

Viewing summary

There should be a summary view which is an overview for the given worksheet. It should contain four distinct sections of information of the worksheet: template section, tables section, operations section, and outcomes section.

Viewing tables

The tables included in the template should be enlisted in a worksheet summary.

Adding column operations

Users must have the possibility of adding new anonymisation operations to the anonymisation concerned columns.

Viewing column operations

The so far accumulated column operations should be enlisted for each individual database table. Different operations should use different colors to provide better visualization. Apart from accumulated operations, the view should also show information such as column name, type, nullability, or whether the column is a primary or foreign key.

PK and FK detection

Primary and foreign keys should be disallowed to undergo anonymisation. An appropriate information should exist to indicate to the user that a given column is a primary or foreign key.

Starting outcome generation

The user should be able to start the generation of an outcome based on the worksheet he prepared. It should be possible to select a compressed archive or script file as a target output for the outcome. Starting the outcome generation must happen in an asynchronous non-blocking way.

On a high level, the outcome generation should include the complex sequential steps of:

- creating a writeable mirror database of the read-only template database
- creating a new empty anonymisation script
- populating the anonymisation script based on the user-defined anonymisation setup

- executing the anonymisation script against the mirror database
- dumping the anonymised database to the compressed archive or script file

There should be an outcome database entity created with a status field. Each consecutive step should update this status. The total time it took to process an outcome should be measured.

Browsing outcomes

The user has to have the possibility to browse the outcomes he created in a form of a dedicated table. Each outcome should include the status, anonymisation script download button, outcome dump download button, processing time, template name, etc.

Anonymisation script download

There should be a way to download an anonymisation script.

Anonymised dump download

There should be a way to download an anonymised dump result.

Account creation

To create a new account, an anonymous user needs to provide valid information including e-mail address, password and confirmation, name, surname, and optionally the purpose for anonymisation. Creating an account with an e-mail that is already in use must be prohibited. A new account must be verified to access the broad functionalities.

Account verification

Upon successful account creation, the user must receive an e-mail with a link to verify the account. User role should change from *unverified* to *verified* due to confirmation.

Admin verification

Admin has the means to verify the specified user on his behalf from the users panel.

Verification expiration

The verification link that is sent to the user after creating a new account should contain a token that expires after a set amount of time. The system administrator has to have a possibility of runtime configuration of the expiration time.

User profile

The user should be able to view and edit his profile details. The user must not be able to change his e-mail address.

Users panel

There should be a users panel for the admin to present an overview and manage the users' accounts.

Block user

Admin should have the means to block the user from the users panel, effectively disabling the possibility of logging in.

Unblock user

Admin will have the means to unblock the blocked user from the users panel, restoring the login capabilities.

Expiring accounts

Users that do not login to the platform for a configured amount of time should be deleted due to the account expiration. The user should receive a notification e-mail prompting him to login to prevent the removal action.

Mark account for removal

The user must be able to initiate the process involving the removal of his account along with all the associated data, hence ensuring GDPR compliance.

The process is started by an explicit user's requests to do so. The system should warn the user with a detailed message and a prompt to type in both the confirmation and the password. If successful, the account is marked for future removal – the account is removed by a dedicated cron-based data pruning service.

Forcing account removal

Admin must have the possibility to force the account to be removed in case of urgent needs e.g. due to reasons concerning system security or performance.

Revert marking account for removal

A scenario in which a user changes his mind to delete his account needs to be supported. The user can revoke his request to remove the account by accessing the link that was earlier sent to his e-mail address.

The revert action should be possible only if all three conditions are met:

- the link has not yet expired
- the admin did not explicitly force the account removal
- the account has not already been removed by the pruning service

Revert account removal expiration

The undo removal link sent to the user after requesting his account removal contains a token that expires after a set amount of time. The time of exactly how long is the link valid should be determined by the runtime configurable by the system administrator.

Forgot password

There should be a way to reset the password that was forgotten. The link should be sent to the user's e-mail address and should be valid only for a configured duration.

Tasks panel

There should be a tasks panel for the admin to retrieve and monitor the tasks. The panel should permit a non-blocking on-demand execution of the selected individual tasks.

3.1.2 Functional requirements – infrastructure**Login**

The user should be able to login to the system and access his account's resources.

User context

The client should know the details of the user who is logged in. There should exist a user context functionality. Upon successful login the information about the principal user should be fetched and stored on the client's side in the user context.

Access token

A successful authentication should generate an access token that encompasses all the information necessary for user authentication and authorization. The token should be valid for a brief duration, e.g., 5 minutes – which should be a matter of runtime configuration conducted by the system administrator. The token should be implemented using JSON Web Token (JWT) standard and should be stored on the client's side.

Refresh Token

The refresh token which is generated together with the access token during authentication should be used to periodically generate a new pair of the access and refresh tokens. There should exist a REST endpoint that will accept the currently valid refresh token and will output new access and refresh token - effectively prolonging the user's session. This token should be valid for a very long duration, e.g., 6 months – which should be a matter of runtime configuration.

Logout

Discarding the access and refresh tokens permits user logout.

Remember me

The user should not have to login to his account when starting a new visit unless

- the refresh token has expired after e.g., 6 months of inactivity
- the user has cleared browser's cache
- the user is logging from a new device

Restarting the anonymisation platform should not logout the user.

Seamless tokens regeneration

Axios HTTP client should be configured in such a manner that it will detect a state of the expired access token and will seamlessly generate new tokens. When the access token expires, the client should invoke refresh token service to generate new token pairings in such a manner that the user does not even notice this implementation-related behaviour. It is required that no unnecessary page refreshes, errors, logouts, or prompts to login show up to the user.

Role-based authorization

All non-whitelisted REST resources should be protected with a declarative role-based authorization.

Resources whitelist

The system administrator should be able to configure the REST resources that are whitelisted (i.e., unprotected). These endpoints should not require to be authenticated (i.e., logged in) to access them.

Redirection

The platform should redirect the user to the login page when accessing routes that are non existing or forbidden to access for them.

Password encryption

User password should be stored in the database in an encrypted form using e.g., bcrypt hashing.

Navigation menu

The client application should have a navigation menu. The panel should be located on the left side of the screen. Each navigation item should consist of a large icon and a text label. The entire panel should be hideable – the entire page content should seamlessly adapt to this. Optionally, smooth transition animation should be configurable during TypeScript compilation time.

The items should be easily defined, ordered, and configured by adopting a standard. Visibility of a particular item must be role-based, e.g. only admin can see *users* item.

Timestamp data collection

The system should collect timestamp data about various resources such as users or templates. Some of the data should be client renderable, and some of the data should be merely collected.

Exemplary dates concern:

- account registration, block, request for removal, removal
- template generation

Requests data collection

All HTTP requests to the anonymisation server should be logged in the system. The logged information should include HTTP method, HTTP response status, and the request URI.

Request validation – client

The system should disallow sending HTTP requests that fail to meet the validation criteria set with Yup library.

In-depth validation – client

Validation rules set by Yup library for the client forms should be enforced by React Hook Forms.

The client should have the possibility to specify required and optional fields and use custom validation rules such as e-mail or file presence validation.

Display validation errors

The client displays validation errors below concerning HTML inputs.

Error notifications

The frontend view uses toast pop-ups to display client and server HTTP errors [33] – for example, 400 Bad Request. The notifications along with the associated message are displayed in the upper-right corner. The pop-ups are red to underline its purpose and the purpose should be further emphasized - by adding more red color – when displaying the errors of extreme importance e.g. database restore failure.

Success notifications

Client-sent HTTP requests that finished with a successful response [33] – for example, 201 Accepted - should be displayed in a green coloured pop-up along with the associated message.

Request validation – server

The system must validate the HTTP requests on the server to protect the system from passing malformed requests sent by malicious external HTTP clients. The validations should be declaratively used on the Data Transfer Object (DTO) classes.

Exemplary validation constraints could include checking for values that meet any of the following criteria:

- not null
- not blank
- not empty
- for integers – of a value in an allowed range
- for strings – of size in an allowed range
- non-standard such as e-mail

Maximum file size

The system administrator should be in control of the configuration of the maximum accepted sizes of the files sent to the system (i.e. database dumps). The file sizes that are too large must not be accepted.

In-depth validation – server

The server should perform numerous other advanced validations specific to individual scenarios.

Exemplary infrastructure related cases should check if:

- the system is running in cloud environment
- an entity non compliant with the database schema is trying to be persisted
- database connection is established

Exemplary authorization related cases should check if:

- user is trying to access another user's resource
- user is unauthorized to access the resource due to lack of privileges
- long lasting refresh token has expired

Exemplary domain related cases should check if:

- anonymisation operation is applicable the particular column
- column is a primary or foreign key
- a given task is manually executable

Client tables settings

The tables presented used extensively by the client application should offer the functionalities of:

- pagination to a set page size
- filtering by attribute
- sorting by attribute
- hiding and showing of columns

File selector

There should be a custom file selector input that will allow uploading files i.e., database dumps. The component should show the upload progress bar and success or error message upon completion.

Uniform design

The user must not be surprised with a particular page looking in an unexpected way. There should be established a uniform design for the client pages. The design should uniformise, e.g., colors, borders, paddings, margins, shadows, or components.

Theme

The client application should allow an easy way to change the theme colors of the entire application by changing it in one configuration.

Skeletons and spinners

The client should use skeletons and spinners to provide a seamless user experience.

Mail notifications

The system has the capabilities of sending e-mails notifying about various actions such as expiring accounts. This functionality should be runtime configurable by the system administrator.

Database

The system needs to use a database. The system should not be coupled to a particular database provider.

Database preloader

There should be a functionality to preload the system with prepared data. A new system instance must include at least one administrator account.

Dynamic database connection

There should exist a factory for data sources that will establish a connection to the dynamically created databases restored from dump files uploaded by the users.

Prune connection pool

There should be a functionality to drop all existing database connections, effectively allowing the creation of the database mirrors dynamically.

Containerization

Software components of the system should be containerized. All services should be running as containers when running the system in the production environment.

Multiple environments

The designed system has to be runnable in multiple environments, i.e.:

- production
- development
- development with the *server* profile enabled
- locally on the host with a database container
- fully locally on the host

Environment based configuration

Different environments, i.e., production, develop, and local should use different configurations of the environment and anonymisation server.

Storage configuration

The location of the files managed by the system, i.e., uploaded templates, anonymisation scripts and anonymised dumps (script based or compressed archive-based), should be configurable.

Running processes from Java

There should be the possibility of running external processes such as *psql* directly from the Java code.

Asynchronous events

There should be a configuration for handling asynchronous events. The core functionalities of *uploading* and *processing* server modules should be initiated with an asynchronous non-blocking event.

Scheduling

There should be a functionality for defining schedulable task services. The services should rely on the cron expression-based triggers. Individual tasks should be runtime configurable in terms of their:

- cron expression
- whether the task is enabled for automated scheduling
- whether the task is enabled for manual execution
- displayed description in the panel

Secrets

There should be a functionality prepared to load the secrets (i.e., password for the database, JWT secret key, mail service password) from an external file non-published to the public git repository. Continuous Integration related secrets should be configured directly in GitHub.

SonarQube

Both the client and the server should be scanned by the SonarQube static analysis tool. This software should be integrated into the Quality CI workflow.

CodeQL

The application should be scanned by CodeQL for security issues as a part of the Security CI workflow.

Swagger

The REST API of the anonymisation server is to be documented with Swagger UI.

File storage

There should be a way to store the files, i.e., user input dumps, anonymisation scripts, anonymised output dumps (archives, scripts). The files should be stored in grouped directories whose paths are runtime configurable. The file names should be randomly picked e.g. as a universally unique identifier (UUID).

Prune user data

Service responsible for deleting accounts that were marked for removal should be executed in a periodic manner. Such a service should prune the user data. The frequency of how often is the service executed is described as a cron expression. It should be a system administrator runtime configurable property.

Prune system data

There should be a script pruning the Docker data from the host system.

3.1.3 Non-functional requirements**Modular**

The anonymisation server should be divided into modular architecture where each module has its own responsibility. Modules have to be further divided into layers or packages. The system implements client-server architecture.

Secure

Achieved through, e.g., vulnerability scanning, token-based authentication enhanced with refresh token functionality, client and server validation, containerization, role-based endpoint protection.

Configurable

The platform should be highly configurable by the system administrator. The software can be separately configured on the server, client and containers side.

Confidential

Data sent by the users should not be accessible by the admins.

Extensible

The system should be engineered in a way to be easily extensible with new features such as supporting new database types or new data masking techniques.

Adaptable

The system should be able to easily adapt to changing requirements by the means of clean code practices such as modularity and testability.

Performant

Achieved through, e.g., using lazy loading of the relations and individual optimisations.

Responsive

Achieved through, e.g., performant client routing.

Engaging

Achieved by displaying skeletons, spinners, and notifications that catch the user's attention.

Usable

The system features should be easily understandable – the functionalities must not be confusing to be used.

Testable

The system code should be written in a testable manner, e.g. by ensuring high quality of the code and separation of concerns.

Reliable

Achieved through continuous integration testing and using well-known secure libraries.

User experience

Achieved through client functionalities such as intuitiveness, responsive design, and uniform styling of the system.

Measurable

The system should collect metrics about e.g. processing time of outcomes.

Modern

The system must not use obsolete technologies or software engineering design techniques.

3.2 System engineering

3.2.1 Use cases

The use cases were split into two distinct diagrams which group the related use cases concerning:

- management – shown in Fig. 3.1
- anonymisation – shown in Fig. 3.2

3.2.2 REST API

The communication with the anonymisation server is driven by the RESTful APIs. An overview of the available endpoints is listed in Tab. 3.1. The particular endpoint's purpose can be inferred from the surrounding context i.e., requirements, use cases, controller name, and the endpoint URI.

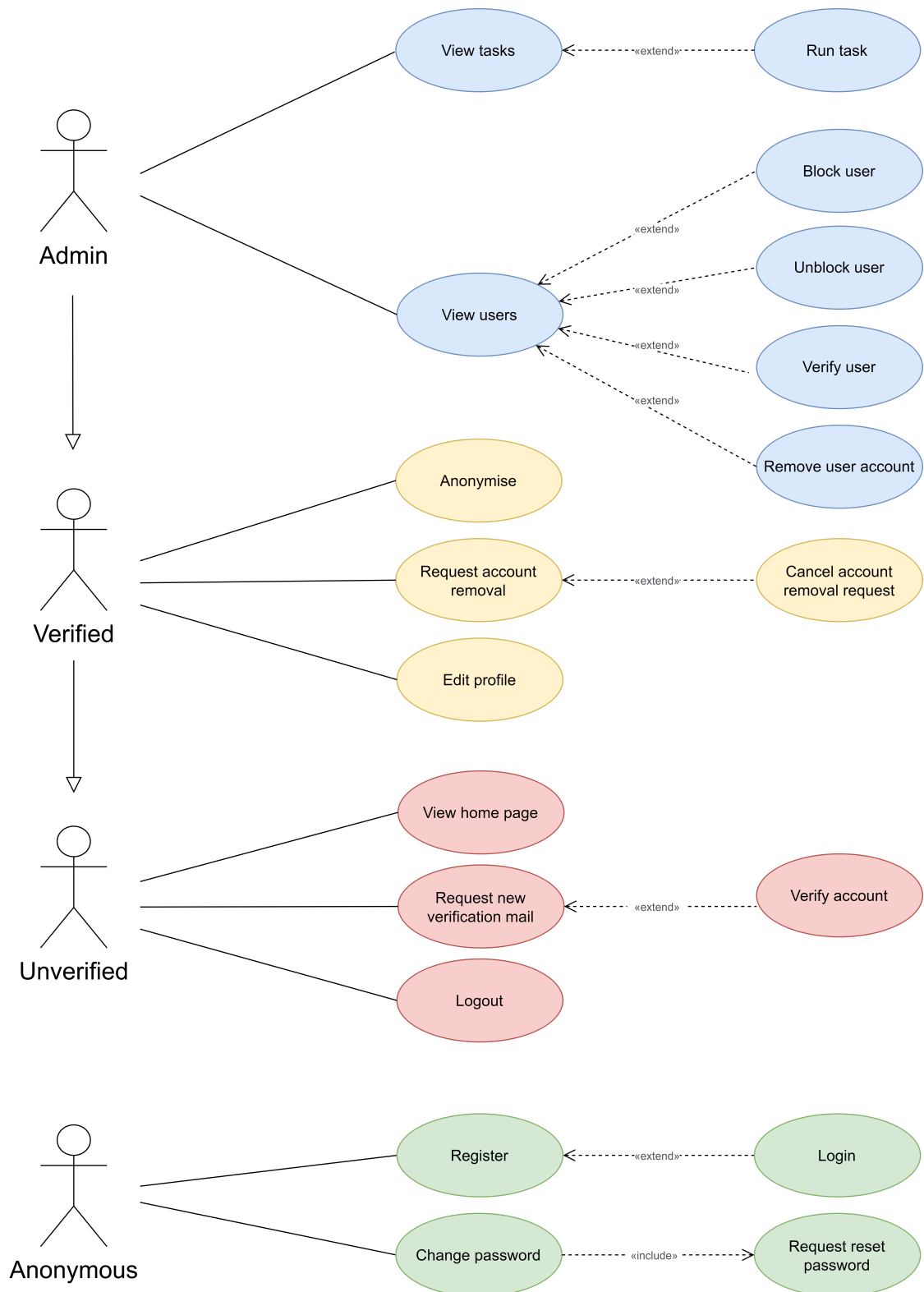


Figure 3.1: Use cases – core domain

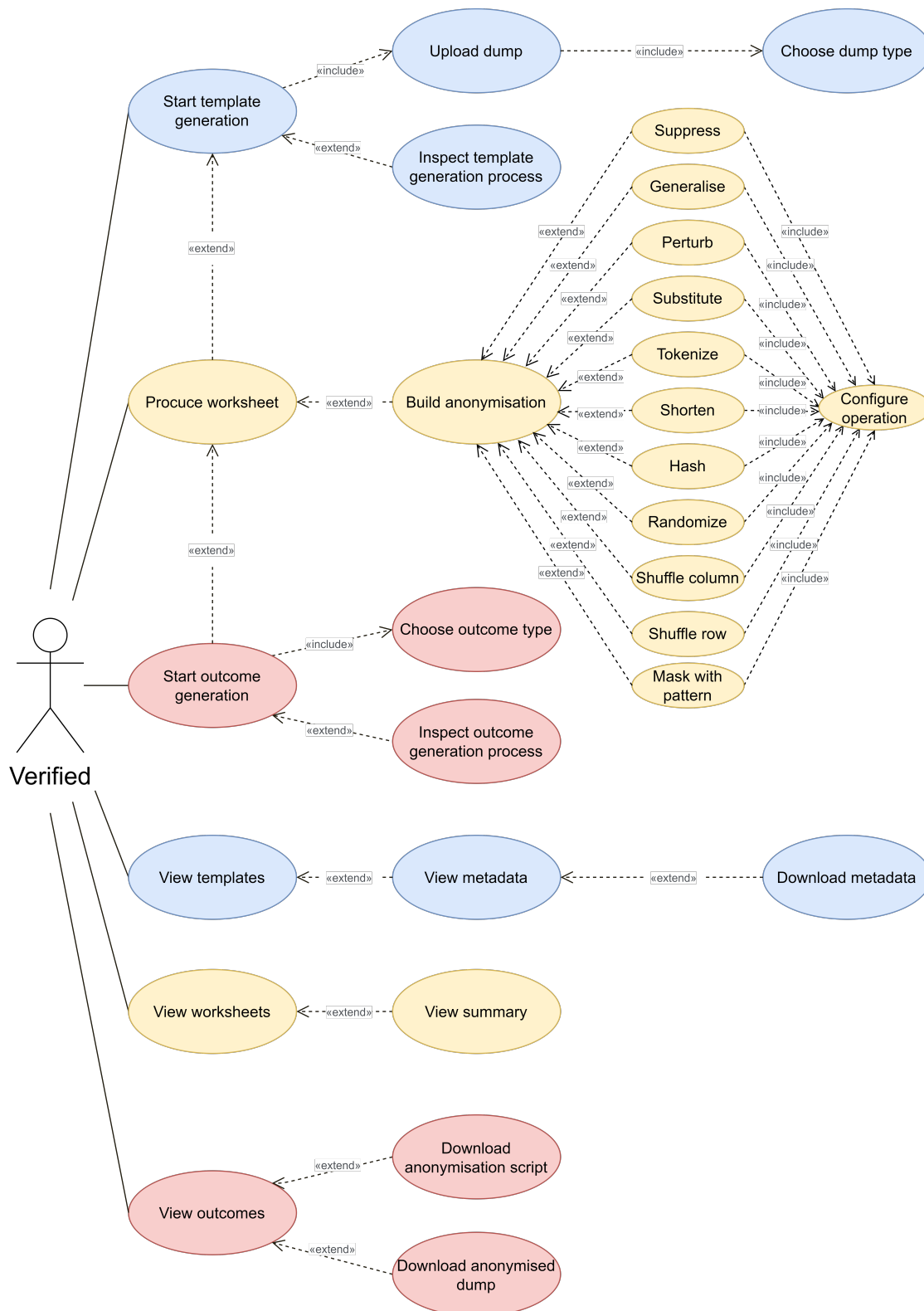


Figure 3.2: Use cases – anonymisation domain

Table 3.1: REST resources – summary.

Controller	Type	Endpoint
Auth	POST	/api/v1/auth/login
	POST	/api/v1/auth/refresh-token
User	GET	/api/v1/users
	POST	/api/v1/users/register
	PUT	/api/v1/users/id/block
	PUT	/api/v1/users/id/unblock
	PUT	/api/v1/users/id/force-removal
Verify	POST	/api/v1/users/verify-mail
	POST	/api/v1/users/verify-mail/send-again
	POST	/api/v1/users/{id}/confirm-mail-verification
ResetPassword	POST	/api/v1/reset-password/request-reset
	POST	/api/v1/reset-password/change-password
	GET	/api/v1/reset-password/show-change-password-form
Task	GET	/api/v1/tasks
	POST	/api/v1/tasks/{task}/execute
Template	POST	/api/v1/templates
	GET	/api/v1/templates/me
	GET	/api/v1/templates/{id}/status
	GET	/api/v1/templates/{id}/metadata
	GET	/api/v1/templates/{id}/dump/download
Worksheet	POST	/api/v1/worksheets
	GET	/api/v1/worksheets/me
	GET	/api/v1/worksheets/me/{id}/summary
TableOperation	GET	/api/v1/worksheet/{id}/table-operations/{table}
ColumnOperation	PUT	/api/v1/worksheet/{id}/column-operations/add-suppression.
	PUT	/api/v1/worksheet/{id}/column-operations/add-generalisati.
	PUT	/api/v1/worksheet/{id}/column-operations/add-perturbation
	PUT	/api/v1/worksheet/{id}/column-operations/add-column-shuff.
	PUT	/api/v1/worksheet/{id}/column-operations/add-hashing
	PUT	/api/v1/worksheet/{id}/column-operations/add-pattern-mas.
	PUT	/api/v1/worksheet/{id}/column-operations/add-random-num
	PUT	/api/v1/worksheet/{id}/column-operations/add-row-shuffle
	PUT	/api/v1/worksheet/{id}/column-operations/add-shortening
	PUT	/api/v1/worksheet/{id}/column-operations/add-substitution
	PUT	/api/v1/worksheet/{id}/column-operations/add-tokenization
Outcome	POST	/api/v1/outcomes/generate
	GET	/api/v1/outcomes/me
	GET	/api/v1/outcomes/{id}/status
	GET	/api/v1/outcomes/{id}/anonymisation-script/download
	GET	/api/v1/outcomes/{id}/dump/download

3.3 Software tools and technologies

The following listings are presented to give an overview of the designed system in terms of tools and technologies involved to engineer it without diving into the implementation details. This is a brief summary that merely presents the technological landscape of the designed system – should a particular technology receive a greater detail of explanation, e.g., exemplary usage, it will be described in another dedicated unit.

3.3.1 Tools

This section encompasses a summary of tools needed to develop, build, and maintain the system.

Docker

Docker is a virtualization software used to develop, run and manage applications in the form of containers taking advantage of Linux kernel features namespaces and cgroups. It consists of three components: a runtime, a daemon engine, and an orchestrator.

To run the production environment of anonymisation platform, only prior installation of Docker on the host system is required, because containerization technologies like Docker enable a separation of the software from the infrastructure the software is running on [26, 57].

Apache Maven

Apache Maven is an open-source build tool that automates the build process of foremostly Java based projects. Primary Maven objectives include ensuring an easy build process and an uniformly established build system which is achieved through Maven's build lifecycle (i.e., predefined build steps, e.g., *validate*, *compile*, *test* or *install*) [36].

Anonymisation platform uses this tool to build the backend server. Maven is put into service in one of the two ways – depending on the selected environment – namely:

- internally from the Docker container where Maven can be considered as an abstraction that the system administrator does not have to be aware of,
- on the host system itself.

Maven is also used to run the tests. Gradle could be a great alternative to Maven, the choice being preference based in usual scenarios.

Node.js

Node.js is a Javascript runtime engine encapsulating V8 engine developed by Google which implements ECMAScript [43]. While node.js may function as a dedicated backend server, in the case of an anonymisation platform it is used merely to build the client application with yarn.

yarn

Yet another resource negotiator (yarn for short) is a package manager developed by Facebook, Google and others [53], whose primary purpose is – similarly to *npm* – installation and management of dependencies used by Node.js runtime environment based projects. Yarn also focuses on ensuring high level of security, performance and reliability [32].

While the purpose of Maven is to build and run the backend server, yarn's purpose is to build and run the frontend client application. Consequently, yarn is used in the same way as Maven – either internally in Docker or on the host system – depending on the chosen environment.

nginx

Nginx is an open-source that works as a reverse proxy, web server and HTTP load balancer with the first one being the primary function for an anonymisation platform. Reverse proxies improve resiliency, security, scalability and performance of the software systems. This is particularly relevant when designing systems that scale horizontally [24].

git

It's been a while since most developers had already started using git for development, which is the leading software for version control, however not everyone knows it was originally invented by Linus Torvalds specifically for Linux kernel development needs [75]. It is clearly necessary to use git during software development due to the essential needs to publish, stash or revert the code changes [14].

GitHub Actions

GitHub Actions is a continuous integration (CI) and continuous delivery (CD) tool that enable creating workflows which automate development processes such as building, testing and deployment of the system.

Developed software uses three distinct workflows to ensure an automation of testing, quality assurance and security resilience.

SonarQube

SonarQube is a large scale open-source tool to conduct the code inspection through static analysis methods. This tool platform can reliably detect vulnerabilities, security hotspots, code smells and bugs.

In the context of designed anonymisation platform, this software is integrated together with Github Actions to ensure quality and security of developed software. Two distinct instances of SonarQube are configured, independently for the client and the server.

ESLint

ESLint is yet another tool for static code analysis for JavaScript applications.

It had been necessary to enhance this tool with TypeScript extension in order to successfully integrate it with client application and Github Actions workflow concerning quality.

Prettier

For an unified code formatting standard to be established for the client's application, prettier tool which is a modern code formatter was configured within JetBrains Webstorm IDE.

3.3.2 Server

This section encompasses a summary of technologies used for the implementation of the modular and REST-driven anonymisation server.

Java

Anonymisation server uses the newest long-term support JDK 17 version released in September 2021, although the server is backward compatible with JDK 11.

Java is still globally the most widely adopted development language [13], even after the 25 years that had passed after its first release – it is a truly mature language, accompanied by all types of frameworks and libraries that one can imagine.

Spring Framework

Spring projects can be viewed together as a modular toolkit that addresses the concerns of modern software development [78].

Spring driven software system may include the following Spring projects [31]:

- Spring Framework - the core supporting e.g., dependency injection and web applications
- Spring Boot – to achieve an easier setup of the Spring driven system
- Spring Data – to retrieve information from databases in a consistent, database-agnostic way
- Spring Security – to secure the application with authentication and authorization

Java-based systems typically employ Spring technologies, and so does the anonymisation server, which extensively uses all of the listed above projects.

Hibernate

Hibernate – a JPA implementation – which is a well-tested and high performant [49] object-relational mapper was a natural choice for the engineered system given the fact that Spring has built-in support for this library. It is used to map the database entities to the corresponding database tables [62].

Hibernate functionality was additionally enhanced with Hibernate Types library [48] to provide the support for JSON column type – which is not supported by the plain Hibernate but is greatly supported by the underlying PostgreSQL database [39].

PostgreSQL – Server

The choice of an appropriate database was of particular importance due to the domain of the designed system.

PostgreSQL was the fourth most popular database management system worldwide back in 2017 as shown in the annual Stack Overflow Developer Survey report [51]. Each consecutive report demonstrated further growth of PostgreSQL. As PostgreSQL is at the continuous growth, MySQL (i.e., the most popular database) is at the continuous decline. The latest report published in 2021 [52] placed PostgreSQL in second place and it is believed that it will at last beat MySQL this year.

For this reason, PostgreSQL was chosen as the database for the server system. This database was also chosen as the supported database type for the data dumps uploaded by data controllers for the purpose of anonymisation.

PostgreSQL – Client

PostgreSQL client is installed in the containerized anonymisation server to execute *psql* processes throughout the application's runtime, hence providing the communication capabilities with the potentially remote PostgreSQL server.

The client is used by the ZT Processor Executor library – directly from the Java code.

ZT Processor Executor

ZT Process Executor provides capabilities to run external processes directly from the Java code.

This library found many domain-critical applications related to PostgreSQL resources management. The library executes PostgreSQL client processes including:

- creation of a new database,
- replication of an existing,
- database dump restoration,
- anonymisation script execution,
- database dump to script or to archive.

JUnit 5

JUnit 5 is the new and refined version of the worldwide recognized Java testing framework. It is used to drive the unit and integration tests [76].

Testcontainers

Testcontainers is a cutting-edge Java library to support the instantiation of lightweight Docker containers for the tests needs. Effectively – the tests run in a containerized environment just like the real application hence blurring the technical differences between real and test environments [5].

This state-of-the-art library is used by the integration tests to bootstrap the PostgreSQL database for the tests.

REST Assured

REST Assured is an integration tests library and an HTTP client which allows sending true HTTP requests to the test server (bootstrapped by Testcontainers library). This library also offers the means to validate received HTTP responses [3].

The library is used by the integration tests.

Swagger UI

One of many Swagger modules is Swagger UI which deals with the design, description, and documentation of REST APIs. The library allows visualization and interaction with the API [70].

This library is relevant due to the complexity of the designed system, as it provides the API documentation for the client.

JJWT

JSON Web Tokens are required to transport the implemented access and refresh tokens necessary for authentication and authorization purposes. This popular library deals with the creation and validation of the tokens [4].

Bouncy Castle

One of many cryptographic APIs offered by Bouncy Castle library is SHA 3 which internally implements Keccak-f[1600] algorithm [11].

This library is used by the hashing data masking technique.

Apache Commons Lang

The sole purpose of this library is a generation of random alphabetic and alphanumeric characters for the pattern masking anonymisation technique [37].

3.3.3 Client

This section encompasses a summary of technologies used for the implementation of React-based client.

TypeScript

TypeScript is a transcompiled superset of JavaScript that the developers want to work with the most unless already doing so [52]. The language was designed by Microsoft to extend JavaScript with types, interfaces, and generics.

The language is powering the frontend client application and React.

React.js

React.js was the most popular web framework as of 2021 [52] as it has finally surpassed the obsolete jQuery library. This is an open-source library developed by Facebook to build user interfaces for the web. As a high-level overview – React centers around the concept of virtual DOM [79].

React ecosystem

Independent authors create many external libraries to drive React with extended functionalities. React ecosystem is rich in various libraries which typically include *React* word in their names.

Anonymisation client code extensively uses the following popular libraries:

- React Router – to handle the URL related concepts including routing, navigation and query parameters
- React Query – as an abstraction and improvement layer for fetching the data
- React Hook Form – to provide user input validation in the forms based on Yup defined rules

Yup

Yup is a validation schema builder. In the case of engineered software, the schemas created by Yup are consumed by React Hook Form library to validate the forms.

Material UI

Material UI is yet another React library and a large one. This library is responsible for the appearance of the application.

Axios

Axios is an asynchronous HTTP client used for sending requests to the anonymisation server.

3.4 Development methodology

The development was driven by git like in the case of any software project that is of non-trivial size. A total number of 6XX **todo** commits¹ were pushed to the master branch. Conventional Commits specification was followed in a simplified way. The branches were deleted following the merge.

There were three possible issue types:

- feature – describing a new functionality
- improvement – describing an upgrade to the existing functionality
- epic – describing a big user story that aggregates other issues [42]

Note that an improvement to the functionality that is still under development is considered a *feature* and not an *improvement*.

Furthermore, the labels with the increasing values of XS, S, M, L, and XL were being added to all the tasks upon creation to estimate the time-to-finish complexities. The distribution of the estimations is shown in Fig. 3.3. **todo updated pie chart**

The non-epic tasks were marked either as resolved or as canceled when closing the tasks. The vast majority of tasks were successfully resolved as shown in Fig. 3.4. Exemplary reasons to close a given task include changing requirements or task duplication.

¹the data as of 21.01.2022

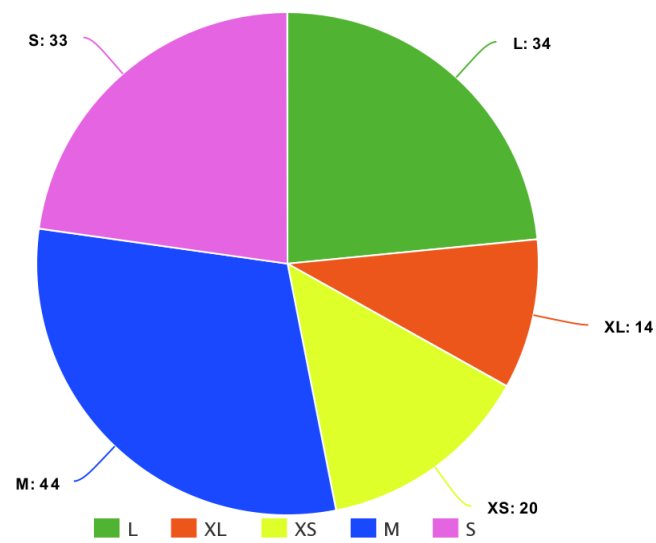


Figure 3.3: GitHub – distribution of tasks estimations.

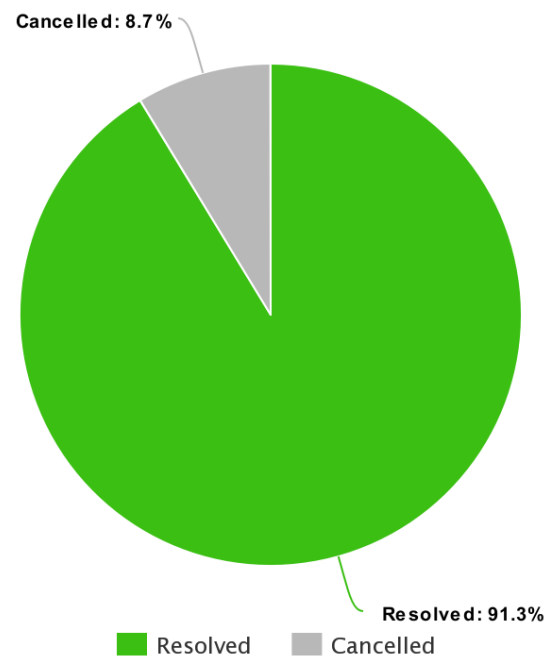


Figure 3.4: GitHub – distribution of resolved tasks

Chapter 4

External specification

4.1 Domain specific glossary

The application defines several business names related to the anonymisation domain to concisely introduce complex technical problems. These names include:

- **Template** – the user restores the database from the dump he provides to create a read-only replicable template that is used to produce worksheets.
- **Metadata** – represents the information about the template database, i.e., information concerning columns, tables, primary keys, foreign keys, number of tables and records, etc.
- **Worksheet** – a concept that represents the user’s working area for configuring the anonymisation to match the exact business requirements of his interest. The worksheet is produced from the template and is used to generate the outcomes.
- **Summary** – this is the worksheet’s representation – a summary layer of the configured anonymisation that presents in one place all the information concerning the template, worksheet, operations, and outcomes.
- **Operation** – an individual anonymisation technique applied to the given attribute.

- **Outcome** – the outcomes are generated from the worksheet after processing the anonymisation request and encapsulate the anonymisation script and anonymised dumps.

4.2 Installation

One of the primary reasons why the containerization concept was so remarkably well-received [52] and exhibited by the massive global adoptions of the Docker technology is the ease it offers regarding software installation. Decoupling the software from the hardware it runs on allows a portable way to install the software in a simple manner. Instead of making it a responsibility for the end-users to manually install the dependent software (e.g., database server, JDK, Node.js, or yarn), all the dependencies come encapsulated within the image, effectively transferring the difficulties of managing the dependencies to the application developers.

4.2.1 Environments — overview

The software solution is designed to be runnable in multiple deployment environments hence reproducing the deployment characteristics of a simple modern business-class software system.

As depicted in Fig. 4.1 the available environments include:

- production
- development

The environments are further divided into three available setup options, i.e.:

- cloud setup
- semi-cloud setup
- local setup

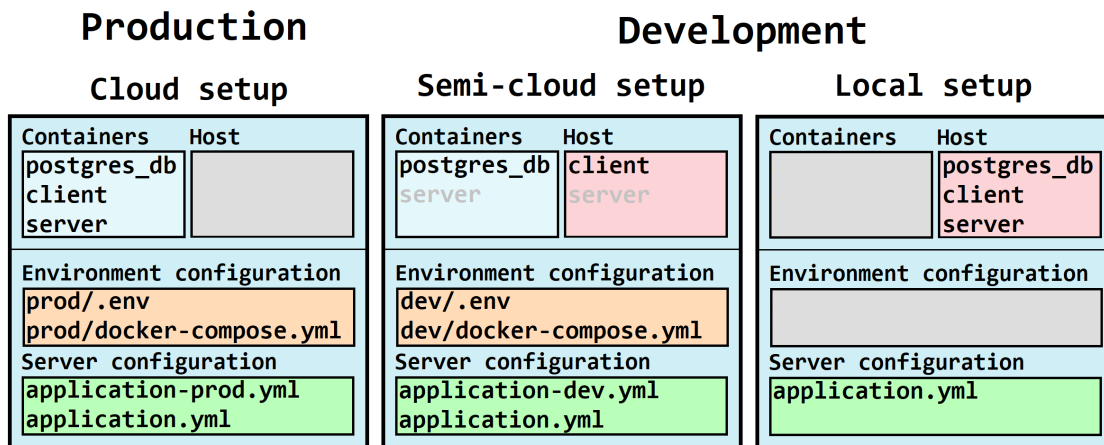


Figure 4.1: Environments

The shared prerequisite for the cloud and semi-cloud setup is to have Docker up and running on the host system. Furthermore, the development environments Local setup does require the installation of additional dependencies.

4.2.2 Cloud setup

The cloud-ready production environment is the easiest one to install.

Prerequisite

Docker up and running.

Installation procedure

Start up everything at once:

```
cd docker/prod
docker compose up
```

If successful the application should be accessible at `http://localhost:3000`.

The setup is customized from:

- `prod/.env`
- `application.yml`

- `application-prod.yml`

4.2.3 Semi-cloud setup

This is the preferred setup for development. PostgreSQL *server* is abstracted away from the host thanks to this setup. Optionally, the anonymisation server can also run in the container.

The semi-cloud setup requires the installation of additional software.

Prerequisite

JDK 17, Node.js, and yarn need to be installed on the host system — apart from the Docker. Note that the Apache Maven which is the build tool for the anonymisation server does not need to be installed neither in the host nor in the container because Maven wrapper is provided and used.

The semi-cloud setup requires an additional installation of the PostgreSQL *client* in the host system if the setup is run without the server profile. The minimum supported version is version 13.5 released on November 11, 2021.

The setup is customized from:

- `dev/.env`
- `application.yml`
- `application-dev.yml`

Verify if the required software is installed:

```
java --version
psql -V # if running without server profile
node -v
yarn -v
```

Installation procedure

Start up the database server:

```
cd docker/dev
docker compose up
```

Start the anonymisation server:

```
cd backend
./mvnw spring-boot:run
```

Alternatively, the server could be started up in the container together with the database server by specifying the server profile:

```
cd docker/dev
docker compose --profile server up
```

Start the client:

```
cd frontend
yarn install
yarn start
```

4.2.4 Local setup

While the preferred setup for development is the semi-cloud, the local on-premise setup is the preferred option if the Docker is not installed.

The setup is customized from:

- `application.yml`

Prerequisites

Apart from Docker, local setup requires everything that the semi-cloud setup requires, and an additional installation of the PostgreSQL *server*.

Verify if the required software is installed:

```
java --version
postgres -V
psql -V
node -v
yarn -v
```

Installation procedure

```
cd backend
./mvnw spring-boot:run
```

Start the anonymisation server:

```
cd backend
./mvnw spring-boot:run
```

Start the client:

```
cd frontend
yarn install
yarn start
```

If all steps are successful the application is fully accessible.

Note that the services are decoupled from each other, e.g. the client service (frontend) does not necessarily need to be started to access the anonymisation server capabilities, i.e., the server would still be accessible through the REST API.

4.3 Activation

The database of a new instance of the anonymisation platform is preloaded with at least one admin account. The credentials of this initial account can be configured from the properties and must be overridden in production:

```
core:
  preloader.admin:
    login: admin@admin.com
    password: admin
```

4.4 Types of users

Four types of users can interact with the system:

- admin user
- verified user

- unverified user
- anonymous user

It is very important to note that while the admin user can manage other users' accounts, he is not authorized to access the users' confidential data such as anonymisation scripts. The admin user who may be a non-technical business person should not be confused with the system administrator.

The relation between the users and their use-cases are summarized in Fig. 3.1 and 3.2.

4.5 User manual

4.6 System administration

Anonymisation server is highly configurable from `application.yml`, `application-<env>.yml`, and `<env>/.env` files. The idea is that concrete environments may need to behave differently, therefore the possibility to override only the selected local properties must be supported and is achievable through the creation of extra overriding files. A property from `application.yml` may be overridden from the `application-<env>.yml`.

The system administrator must be provided with a possibility to configure the platform in an externalized way. The configuration changes merely require to re-deploy the platform – it is not required to re-build it.

System administrator and admin user must not be confused. The former is a technical user responsible for technical software processes such as configuration and deploy, whereas the latter is a non-technical business user responsible for managing the users.

4.6.1 Environments

Consider the docker compose visible in Fig. 4.3 which is further configured by the environment file visible in Fig. 4.2.

```

1 POSTGRES_TAG=13.1
2 POSTGRES_HOST_PORT=5007
3 POSTGRES_CONTAINER_PORT=5432
4 POSTGRES_PORTS_MAPPING=${POSTGRES_HOST_PORT}:${POSTGRES_CONTAINER_PORT}
5
6 POSTGRES_DB=anonymisation_db
7 POSTGRES_USER=postgres
8 POSTGRES_PASSWORD=postgres
9 POSTGRES_DATA_PATH=/var/lib/postgresql/data
10
11 SERVER_PORTS_MAPPING=8080:8080
12 CLIENT_PORTS_MAPPING=3000:3000

```

Figure 4.2: Environment file – development.

Note that the environment files are typically not uploaded to the repository. The system administrator must configure the docker compose files from the environment files. Most of the configuration is related to the database connection.

The anonymisation server can override the individual properties in `application-<env>.yaml`. The configuration visible in Fig. 4.4 shows an example of how the production environment may want to override certain properties.

```

1 # Miscellaneous
2 server.environment.cloud: true
3
4 # Uploading module configuration
5 uploading.templates.path: /var/lib/anonymisation/templates
6
7 # Processing module configuration
8 processing:
9   anonymisations:
10     scripts.path: /var/lib/anonymisation_platform/anonymisations/scripts
11   dumps:
12     scripts.path: /var/lib/anonymisation_platform/dumps/scripts
13     archives.path: /var/lib/anonymisation_platform/dumps/archives

```

Figure 4.4: application-prod.yaml – example.

Prune containers data

The system administrator may want to execute `docker/prune.sh` script to prune all Docker-related data. The prepared script may be useful for development and must not be used for production.

To run the script after the prior assessment of necessity to do:

```
1 cd docker
2 sh prune.sh
```

The script's internals:

```
1 #!/bin/sh
2 docker system prune -f
3 docker container stop $(docker container ls -aq)
4 docker container rm $(docker container ls -aq)
5 docker rmi $(docker images -aq)
6 docker volume prune -f
```

4.6.2 Database

The primary database configuration can get complex.

Schema

By default, the database schema is dropped and re-created when the server restarts. This is configured through `spring.jpa.hibernate.ddl-auto` and must be overridden in production. Typically, production environments need to execute the schema-altering SQL scripts manually and tools like Flyway [44] are designed to support the versioning of the database.

```
1 spring:
2   jpa.hibernate:
3     ddl-auto: create
4     show-sql: true
5     generate-ddl: false
```

When `spring.jpa.hibernate.show-sql` property is enabled, all SQL queries are logged and enabling the `spring.jpa.hibernate.generate-ddl` will generate the log of the initial database schema. The script could be used by Flyway as the first initial version of the database.

Connection

The database url, username and password are all configurable and individual deployment setups override this property.

Local setup:

```
1  spring:
2    datasource:
3      url: jdbc:postgresql://localhost:5007/anonymisation_db
4      username: postgres
5      password: postgres
```

For the containerized setups, this property is overridden in docker compose:

```
1  services:
2    anonymisation_server:
3      SPRING_DATASOURCE_URL: jdbc:postgresql://anonymisation_postgres_db
   ↪   :${POSTGRES_CONTAINER_PORT}/${POSTGRES_DB}
```

4.6.3 Files

Files are the fuel for anonymisation platform and hence the need for customizations.

Maximum accepted file size of the user uploaded dumps:

```
1  spring.servlet.multipart:
2    max-file-size: 100MB
3    max-request-size: 100MB
```

The location configuration of the files managed throughout the anonymisation platform: processing:

```
1  anonymisations.scripts.path: stored_files/anonymisations/scripts
2  dumps:
3    scripts.path: stored_files/dumps/scripts
4    archives.path: stored_files/dumps/archives
```

Note that those are relative paths. Containerized environments override these properties to the full paths as depicted in Fig. 4.4.

4.6.4 Mail service

Mail service is also configurable:


```
1  spring:
2    mail:
3      host: smtp.gmail.com
4      port: 587
5      username: data.anonymisation
6      # password: <Specified in secrets.properties>
7      properties:
8        smtp:
9          auth: true
10         starttls.enable: true
```

In this case, the password for the mail service is provided in separate `secrets.properties` file that is not the part of the git repository.

4.6.5 Unprotected resources

By default, all REST resources require authentication and may require authorization. REST resources that do not require the user to authenticate must be configured:

```
1  core:
2    api:
3      unprotectedResources: /api/v1/auth/**,\
4                           /api/v1/users/register/**,\
5                           /api/v1/users/verify-mail/**,\
6                           /api/v1/reset-password/**,\
7                           /api/v1/me/restore-account/**
```

4.6.6 Authorization

The access and refresh tokens are used for authorization and are generated after the user authenticates himself by logging in. The configuration is necessary to control the behaviour of this:

```
1  core:
2    jwt:
3      algorithm: HS512
4      secretKey: /EOAqvJOzlPxRdTJ05ib1CYSCGMXsVaCU47Bpvxv019L87
5      ↪ /pVpyoab9sy2rDzHS5vNwHzu8rX/E8FJGqx5oCkA==
6      accessToken.expireTimeInSeconds: 1800
7      refreshToken.expireTimeInSeconds: 2419200
```

Note that the access token expires quickly while the refresh token is to stay for a long time. The secret key needs to be both valid for the selected hashing

algorithm and also base64 encoded. Currently, the only supported algorithm is HS512.

4.6.7 User account

The settings for the expiring time duration of the user account related services:

```
1 core:
2   resetPasswordToken.expireTimeInSeconds: 3600
3   verifyMailToken.expireTimeInSeconds: 604800
4   undoRemoveAccountToken.expireTimeInSeconds: 604800
```

4.6.8 Scheduler

The existing schedulers which are triggered based on the cron expression must be highly customizable, as there are quite a few of them, namely the schedulers are responsible for:

- notifying the users about approaching expiration of their account expiration due to inactivity,
- removing the inactive accounts after the configured duration of being inactive,
- removing the unverified accounts after the configured duration of being unverified,
- pruning the user data who wish to delete their account.

For brevity, the configuration of only one scheduler is shown:

```
1 scheduler:
2   notifyExpiringAccounts:
3     scheduled: true
4     executable: true
5     cron: '0 0 0 * * *'
6   notifyAfterTimeInSeconds: 300
7   description: Send mail notifications to expiring accounts
```

The expression from the example is triggered daily and an easier form of `@daily` could also be used [58]. There is a functionality that is able to compute the exact date of the next scheduled execution and shows it in the tasks panel.

The available parameters include:

- `scheduled` – controls if the given scheduler is periodically launched
- `executable` – controls if the given scheduler is manually executable from the tasks panel
- `cron` – the trigger expression.
- `description` – the description for the tasks panel to identify the scheduler
- `notifyAfterTimeInseconds` – this particular parameter is specific to the scheduler and specifies a way of how long the user needs to be inactive to receive a notification e-mail with a prompt to login

4.7 Security issues

4.7.1 User account

The system requires the user to verify his e-mail address to access the platform capabilities. The application allows the blocking of potentially malicious users. To ensure the security of users' data through data protection measures, the users can request to remove their accounts.

4.7.2 Continuous integration

Multiple advanced continuous integrity workflows were established to follow the state-of-the-art software development practices of ensuring quality and security. The plans verify the state of quality, security, and tests. The workflows are triggered daily at a set time, upon creating a pull request or upon committing to the master branch.

4.7.3 Containers isolation

The platform runs as a secure network of isolated Docker containers hence minimising the possible system vulnerabilities that could be taken advantage of.

Authentication and authorization

Platform endpoints are gated behind the authorization logic, i.e., to access the platform resources the user first needs to authenticate himself by logging in to his account, and then his account is also required to have the appropriate role authority to access the particular resources such as the templates or worksheets.

The authentication and the authorization is implemented using JWT access tokens. To increase the safety of the system, the functionality is further enhanced by the functionality of refresh tokens.

4.7.4 Validation

The application pays special attention to the validation which is performed on the client's side and the server's side in a multi-level fashion. Taking the proper validation measures prevents the potentially devastating effects of putting the application in a breaking, unsupported state.

4.8 Working scenarios

```
1  version: '3.9'
2
3  services:
4    anonymisation_postgres_db:
5      image: postgres:${POSTGRES_TAG}
6      ports:
7        - '${POSTGRES_PORTS_MAPPING}'
8      environment:
9        POSTGRES_DB: ${POSTGRES_DB}
10       POSTGRES_USER: ${POSTGRES_USER}
11       POSTGRES_PASSWORD: ${POSTGRES_PASSWORD}
12      volumes:
13        - postgres_data:${POSTGRES_DATA_PATH}
14      tty: true
15
16    anonymisation_server:
17      ports:
18        - '${SERVER_PORTS_MAPPING}'
19      build:
20        context: ../../backend
21      depends_on:
22        - anonymisation_postgres_db
23      environment:
24        SPRING_PROFILES_ACTIVE: prod
25        POSTGRES_IP_ADDRESS: anonymisation_postgres_db
26        SPRING_DATASOURCE_URL: jdbc:postgresql://anonymisation_postgres_db
27                               ↪ :${POSTGRES_CONTAINER_PORT}/${POSTGRES_DB}
28      tty: true
29
30    anonymisation_client:
31      ports:
32        - '${CLIENT_PORTS_MAPPING}'
33      build:
34        context: ../../frontend
35      depends_on:
36        - anonymisation_server
37
38  volumes:
39    postgres_data:
```

Figure 4.3: Docker compose – production.

Chapter 5

Internal specification

- concept of the system
- system architecture
- description of data structures (and data bases)
- components, modules, libraries, resume of important classes (if used)
- resume of important algorithms (if used)
- details of implementation of selected parts
- applied design patterns
- UML diagrams

Use special environment for inline code, eg **descriptor** or **descriptor_gaussian**.

Chapter 6

Verification and validation

- testing paradigm (eg V model)
- test cases, testing scope (full / partial)
- detected and fixed bugs
- results of experiments (optional)

Portability tests testcontainers, rest assured, junit, mockito, extendable CI working manual bugs

Chapter 7

Conclusions

Wyniki są b. zadow. bo soft może konkurować z innymi.

Bibliography

- [1] G. Aad et al. ‘Observation of a new particle in the search for the Standard Model Higgs boson with the ATLAS detector at the LHC’. In: *Physics Letter B* 716.1 (2012), pp. 1–29.
- [2] David Trigano et al. *Oracle Data Masking and Subsetting*. URL: <https://docs.microsoft.com/en-us/azure/azure-sql/database/dynamic-data-masking-overview> (visited on 21/01/2022).
- [3] Johan Haleby et al. *REST Assured – Documentation*. URL: <https://rest-assured.io> (visited on 21/01/2022).
- [4] Les Hazlewood et al. *jwt – GitHub Repository*. URL: <https://github.com/jwtk/jwt> (visited on 21/01/2022).
- [5] Richard North et al. *Testcontainers*. URL: <https://www.testcontainers.org> (visited on 21/01/2022).
- [6] Jens Albrecht, Marc Fiedler and Tim Kiefer. ‘A Rule-Based Approach to Local Anonymization for Exclusivity Handling in Statistical Databases’. In: *Privacy in Statistical Databases: UNESCO Chair in Data Privacy International Conference*. Ed. by Josep Domingo-Ferrer and Mirjana Pejić-Bach. Springer Publishing, 2016, pp. 81–93. ISBN: 978-3319453804.
- [7] Jens Albrecht, Marc Fiedler and Tim Kiefer. ‘Robust De-anonymization of Large Sparse Datasets’. In: *Privacy in Statistical Databases: UNESCO Chair in Data Privacy International Conference*. Ed. by Josep Domingo-Ferrer and Mirjana Pejić-Bach. Springer Publishing, 2016, pp. 81–93. ISBN: 978-3319453804.

- [8] Luk Arbuckle and Khaled El Emam. *Building an Anonymization Pipeline: Creating Safe Data*. O'Reilly Media, Inc, 2020. ISBN: 978-1492053439.
- [9] Narayanan Arvind and Shmatikov Vitaly. *Robust de-anonymization of large sparse datasets: a decade later*. URL: <https://www.cs.princeton.edu/~arvindn/publications/de-anonymization-retrospective.pdf> (visited on 21/01/2022).
- [10] Arshdeep Bahga and Vijay Madisetti. *Big Data Science and Analytics: A Hands-On Approach*. Vpt, 2016. ISBN: 978-0996025539.
- [11] Legion of the Bouncy Castle Inc. *The Legion of the Bouncy Castle Java Cryptography APIs*. URL: <https://www.bouncycastle.org/java.html> (visited on 21/01/2022).
- [12] Brendan Burns. *Designing Distributed Systems: Patterns and Paradigms for Scalable, Reliable Services*. O'Reilly Media, Inc., 2018. ISBN: 978-1491983645.
- [13] Damien Carey. *Celebrating 25 Years of Java*. URL: <https://blogs.oracle.com/oracleuniversity/post/celebrating-25-years-of-java> (visited on 21/01/2022).
- [14] Scott Chacon and Ben Straub. *Pro Git*. Apress, 2014. ISBN: 978-1484200773.
- [15] Inc. Cisco Systems. *Cisco Annual Internet Report (2018–2023)*. Tech. rep. URL: <https://www.cisco.com/c/en/us/solutions/collateral/executive-perspectives/annual-internet-report/white-paper-c11-741490.pdf> (visited on 21/01/2022).
- [16] Inc. Cisco Systems. *The Zettabyte Era Officially Begins*. URL: <https://blogs.cisco.com/sp/the-zettabyte-era-officially-begins-how-much-is-that> (visited on 21/01/2022).
- [17] European Commission. *Charter of Fundamental Rights of the European Union*. URL: <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:12012P/TXT&from=EN> (visited on 21/01/2022).
- [18] European Commission. *General Data Protection Regulation - Recital 1*. URL: <https://gdpr-info.eu/recitals/no-1> (visited on 21/01/2022).

-
- [19] European Comission. *General Data Protection Regulation - Recital 26*. URL: <https://www.privacy-regulation.eu/en/recital-26-GDPR.htm> (visited on 21/01/2022).
 - [20] European Comission. *General Data Protection Regulation - Recital 6*. URL: <https://gdpr-info.eu/recitals/no-6> (visited on 21/01/2022).
 - [21] European Comission. *General Data Protection Regulation – Art. 83 – General conditions for imposing administrative fines*. URL: <https://gdpr-info.eu/art-83-gdpr/> (visited on 21/01/2022).
 - [22] European Comission. *General Data Protection Regulation – Art. 83 – Lawfulness of processing*. URL: <https://gdpr-info.eu/art-6-gdpr/> (visited on 21/01/2022).
 - [23] Oracle Corporation. *Oracle Data Masking and Subsetting*. URL: <https://www.oracle.com/security/database-security/data-masking> (visited on 21/01/2022).
 - [24] Derek DeJonghe. *NGINX Cookbook: Advanced Recipes for High-Performance Load Balancing*. O'Reilly Media, Inc., 2020. ISBN: 978-1492078487.
 - [25] Divante. *Anonymizer – Documentation*. URL: <https://github.com/DivanteLtd/anonymizer> (visited on 21/01/2022).
 - [26] Inc. Docker. *Docker overview*. URL: <https://docs.docker.com/get-started/overview> (visited on 21/01/2022).
 - [27] Josep Domingo-Ferrer and Jordi Soria-Comas. ‘Anonymization in the Time of Big Data’. In: *Privacy in Statistical Databases: UNESCO Chair in Data Privacy International Conference*. Ed. by Josep Domingo-Ferrer and Mirjana Pejić-Bach. Springer Publishing, 2016, pp. 57–68. ISBN: 978-3319453804.
 - [28] Inc. Domo. *Data Never Sleeps 6.0*. Tech. rep. URL: https://www.domo.com/assets/downloads/18_domo_data-never-sleeps-6+verticals.pdf (visited on 21/01/2022).
 - [29] Khaled El Emam and Luk Arbuckle. *Anonymizing Health Data: Case Studies and Methods to Get You Started*. O'Reilly Media, Inc, 2020. ISBN: 978-1449363079.

- [30] Can Eyupoglu et al. ‘An Efficient Big Data Anonymization Algorithm Based on Chaos and Perturbation Techniques’. In: *Entropy* 20.5 (2018), p. 373.
- [31] Facebook. *Spring – Projects*. URL: <https://spring.io/projects> (visited on 21/01/2022).
- [32] Facebook. *Yarn – Introduction*. URL: <https://yarnpkg.com/getting-started> (visited on 21/01/2022).
- [33] Roy T. Fielding and Julian F. Reschke. *Hypertext Transfer Protocol (HTTP/1.1): Semantics and Content*. URL: <https://datatracker.ietf.org/doc/html/rfc7231> (visited on 21/01/2022).
- [34] Carl Folkesson. ‘Anonymization of directory-structured sensitive data’. MA thesis. Linköping University, 2019.
- [35] Apache Software Foundation. *Apache Kafka – uses*. URL: <https://kafka.apache.org/uses> (visited on 21/01/2022).
- [36] Apache Software Foundation. *Maven – Introduction*. URL: <https://maven.apache.org/what-is-maven.html> (visited on 21/01/2022).
- [37] The Apache Software Foundation. *Apache Commons Lang – Documentation*. URL: <https://commons.apache.org/proper/commons-lang> (visited on 21/01/2022).
- [38] Aris Gkoulalas-Divanis and Grigorios Loukides. *Anonymization of Electronic Medical Records to Support Clinical Analysis*. Springer Publishing, 2012. ISBN: 978-1461456698.
- [39] PostgreSQL Global Development Group. *PostgreSQL Documentation – JSON types*. URL: <https://www.postgresql.org/docs/14/datatype-json.html> (visited on 21/01/2022).
- [40] Richie Koch. *Data anonymization and GDPR compliance: the case of Taxa 4×35*. URL: <https://gdpr.eu/data-anonymization-taxa-4x35> (visited on 21/01/2022).

- [41] Chris Lee. *Integrating Web Analytics with your SAP Commerce Cloud Storefront - An Overview*. URL: https://www.sap.com/cxworks/article/2589633935/integrating_web_analytics_with_your_sap_commerce_cloud_storefront_an_overview (visited on 21/01/2022).
- [42] Patrick Li. *JIRA Agile Essentials*. Packt Publishing, 2015. ISBN: 978-1784394912.
- [43] Google LLC. *V8 JavaScript Engine*. URL: <https://v8.dev/docs> (visited on 21/01/2022).
- [44] Red Gate Software Ltd. *Flyway – Documentation*. URL: <https://flywaydb.org/documentation> (visited on 21/01/2022).
- [45] TeskaLabs Ltd. *TurboCat.io*. URL: <https://teskalabs.com/products/turbocat.io> (visited on 21/01/2022).
- [46] Bernard Marr. *Big Data in Practice: How 45 Successful Companies Used Big Data Analytics to Deliver Extraordinary Results*. John Wiley and Sons Ltd., 2016. ISBN: 978-1119231387.
- [47] Robert C. Martin. *Clean Architecture: A Craftsman’s Guide to Software Structure and Design*. Pearson, 2017. ISBN: 978-0134494166.
- [48] Vlad Mihalcea. *Hibernate Types – GitHub repository*. URL: <https://github.com/vladmihalcea/hibernate-types> (visited on 21/01/2022).
- [49] Vlad Mihalcea. *High-Performance Java Persistence*. Independently published, 2016. ISBN: 978-9730228236.
- [50] Andrew Moore. *The GDPR and Managing Data Risk*. John Wiley and Sons, Ltd., 2018. ISBN: 978-1119487746.
- [51] Prosus N.V. *Stack Overflow Developer Survey 2017*. Tech. rep. URL: <https://insights.stackoverflow.com/survey/2017> (visited on 21/01/2022).
- [52] Prosus N.V. *Stack Overflow Developer Survey 2021*. Tech. rep. URL: <https://insights.stackoverflow.com/survey/2021> (visited on 21/01/2022).
- [53] Christoph Nakazawa, Sebastian McKenzie and Jamie Kyle. *Yarn: A new package manager for JavaScript*. URL: <https://engineering.fb.com/2016/10/11/web/yarn-a-new-package-manager-for-javascript> (visited on 21/01/2022).

- [54] Neha Narkhede, Gwen Shapira and Todd Palino. *Kafka: The Definitive Guide: Real-Time Data and Stream Processing at Scale*. O'Reilly Media, Inc., 2017. ISBN: 978-1491936160.
- [55] Petter Ødegård. 'Data Anonymization for Research'. MA thesis. Norwegian University of Science and Technology, 2019.
- [56] OpenAIRE. *Amnesia – Documentation*. URL: <https://amnesia.openaire.eu/about-documentation.html> (visited on 21/01/2022).
- [57] Nigel Poulton. *Docker Deep Dive: Zero to Docker in a single book*. Independently published, 2016. ISBN: 978-1521822807.
- [58] Arjen Poutsma. *New in Spring 5.3: Improved Cron Expressions*. URL: <https://spring.io/blog/2020/11/10/new-in-spring-5-3-improved-cron-expressions> (visited on 21/01/2022).
- [59] Fabian Prasser et al. *ARX – Data Anonymization Tool*. URL: <https://arx.deidentifier.org> (visited on 21/01/2022).
- [60] Fabian Prasser et al. 'Flexible data anonymization using ARX—Current status and challenges ahead'. In: *Software: Practice and Experience* 50.7 (2020), pp. 1277–1304.
- [61] Balaji Raghunathan. *The Complete Book of Data Anonymization: From Planning to Implementation*. Auerbach Publications, 2013. ISBN: 978-1439877302.
- [62] Inc. Red Hat. *Hibernate Documentation – ORM*. URL: <https://hibernate.org/orm> (visited on 21/01/2022).
- [63] Rolfje. *Anonimatron – Documentation*. URL: <https://realrolfje.github.io/anonimatron/documentation> (visited on 21/01/2022).
- [64] Rathindra Sarathy and Krish Muralidhar. 'Perturbation Methods for Protecting Numerical Data: Evolution and Evaluation'. In: ed. by Ranajit Chakraborty, C.R. Rao and Pranab Sen. Vol. 28. *Handbook of Statistics*. Elsevier, 2012, pp. 513–531. DOI: <https://doi.org/10.1016/B978-0-44-451875-0.00019-1>. URL: <https://www.sciencedirect.com/science/article/pii/B9780444518750000191>.

-
- [65] Rathindra Sarathy and Krish Muralidhar. ‘Protecting Numerical Confidential Data using Data Shuffling: A Demonstration of Effectiveness of Approach and Flexibility of Delivery’. In: *Federal Committee on Statistical Methodology Research Conference*. 2009.
- [66] SAP SE. *SAP Data Services – Data Mask – Documentation*. URL: <https://help.sap.com/viewer/e54136ab6a4a43e6a370265bf0a2d744/4.2.14/en-US/6073fc460d794b699e90ca8ce5bb095c.html> (visited on 21/01/2022).
- [67] Sanjay Sharma. *Data Privacy and GDPR Handbook*. John Wiley and Sons Ltd., 2019. ISBN: 978-1119594246.
- [68] M. S. Simi, Nayaki K. Sankara and Elayidom M. Sudheep. ‘An Extensive Study on Data Anonymization Algorithms Based on K-Anonymity’. In: *IOP Conference Series: Materials Science and Engineering* 225.1 (2017).
- [69] Personal Data Protection Commission of Singapore. *Guide to Basic Data Anonymisation Techniques*. 2018. URL: [https://www.pdpc.gov.sg/-/media/Files/PDPC/PDF-Files/Other-Guides/Guide-to-Anonymisation_v1-\(250118\).pd](https://www.pdpc.gov.sg/-/media/Files/PDPC/PDF-Files/Other-Guides/Guide-to-Anonymisation_v1-(250118).pd) (visited on 21/01/2022).
- [70] SmartBear Software. *REST API Documentation Tool – Swagger UI*. URL: <https://swagger.io/tools/swagger-ui> (visited on 21/01/2022).
- [71] Jordi Soria-Comas and Josep Domingo-Ferrer. ‘Big Data Privacy: Challenges to Privacy Principles and Model’. In: *Data Science and Engineering* 1.1 (2015), pp. 21–28.
- [72] European Data Protection Supervisor. *Introduction to the hash function as a personal data pseudonymisation technique*. URL: https://edps.europa.eu/data-protection/our-work/publications/papers/introduction-hash-function-personal-data_en (visited on 21/01/2022).
- [73] Colin Tankard. *Big data security*. URL: <https://www.sciencedirect.com/science/article/abs/pii/S1353485812700636> (visited on 21/01/2022).
- [74] IT Governance Privacy Team. *EU General Data Protection Regulation (GDPR): An Implementation and Compliance Guide*. IT Governance Publishing Ltd., 2020. ISBN: 978-1787782495.

- [75] Linus Torvalds. *Tech Talk: Linus Torvalds on git*. 2007. URL: <https://www.youtube.com/watch?v=4XpnKHJAok8> (visited on 21/01/2022).
- [76] Catalin Tudose. *JUnit in Action*. Manning, 2020.
- [77] Paul Voigt and Axel von dem Bussche. *The EU General Data Protection Regulation (GDPR): A Practical Guide*. Springer Publishing, 2017. ISBN: 978-3319579580.
- [78] Craig Walls. *Spring in Action*. Manning, 2018. ISBN: 978-1617294945.
- [79] Frank Zametti. *Modern Full-Stack Development: Using TypeScript, React, Node.js, Webpack, and Docker*. Apress, 2020. ISBN: 978-1484257371.
- [80] Sherali Zeadally and Mohamad Badra. *Privacy in a Digital, Networked World: Technologies, Implications and Solutions*. Springer Publishing, 2015. ISBN: 978-3319084718.

Appendices

Index of abbreviations and symbols

DNA deoxyribonucleic acid

List of additional files in electronic submission (if applicable)

Additional files uploaded to the system include:

- source code of the application,
- test data,
- a video file showing how software or hardware developed for thesis is used,
- etc.

List of Figures

2.1	Pseudocode – suppression.	16
2.2	Pseudocode – generalisation based on distribution size.	19
2.3	Pseudocode – generalisation based on number of distributions. . . .	20
2.4	Pseudocode – perturbation based on fixed noise.	20
2.5	Pseudocode – perturbation based on percentage noise.	21
2.6	Substitution.	25
2.7	Tokenisation.	26
2.8	Random number.	26
2.9	Shortening.	27
3.1	Use cases – core domain	50
3.2	Use cases – anonymisation domain	51
3.3	GitHub – distribution of tasks estimations.	62
3.4	GitHub – distribution of resolved tasks	62
4.1	Environments	65
4.2	Environment file – development.	70
4.4	application-prod.yml – example.	70
4.3	Docker compose – production.	77

List of Tables

2.1	Suppression – input data.	16
2.2	Suppression – masked data.	16
2.3	Generalisation – input data.	17
2.4	Generalised – masked data.	17
2.5	Perturbation – input data.	19
2.6	Perturbation – masked data.	21
2.7	Pattern masking – input data.	21
2.8	Pattern masking – exemplary patterns.	22
2.9	Pattern masking – masked data.	22
2.10	Hashing – input data.	22
2.11	Hashing – masked data.	23
2.12	Randomisation – attribute – input data.	23
2.13	Randomisation – attribute – masked data.	24
2.14	Randomisation – record – input data.	24
2.15	Randomisation – record – masked data.	24
3.1	REST resources – summary.	52