

**Автономная некоммерческая организация высшего образования  
«Университет Иннополис»  
(АНО ВО «Университет Иннополис»)**

**ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА  
(МАГИСТЕРСКАЯ ДИССЕРТАЦИЯ)  
по направлению подготовки  
09.04.01 – «Информатика и вычислительная техника»**

**GRADUATION THESIS  
(MASTER GRADUATE THESIS)  
Field of Study  
09.04.01 – «Computer Science»**

**Направленность (профиль) образовательной программы  
«Анализ данных и искусственный интеллект»  
Area of Specialization / Academic Program Title:  
«Data Analysis and Artificial Intelligence»**

**Тема /  
Topic**

**Обработка естественного языка для аудита соответствия  
нормативным требованиям. Автоматизированная проверка  
соответствия соглашений об обработке данных на основе  
обработки естественного языка / NLP for Regulatory  
Compliance Audit. NLP-based Automated Compliance Checking  
of Data Processing Agreements against General Data Protection  
Regulation**

**Работу выполнил /  
Thesis is executed by**

**Оконича Озиома Ненубари/  
Okonicha Ozioma Nenubari**

подпись / signature

**Руководитель  
выпускной  
квалификационной  
работы /  
Graduation Thesis  
Supervisor**

**Садовых Андрей  
Александрович/  
Sadovykh Andrey  
Aleksandrovich**

подпись / signature

**Консультанты /  
Consultants**

подпись / signature

Иннополис, Innopolis, 2024

# Contents

<b>1</b>	<b>Introduction</b>	<b>10</b>
1.1	Background . . . . .	12
1.2	Significance of NLP in GDPR Compliance . . . . .	12
1.3	Problem Statement . . . . .	13
1.4	Research Questions and Objectives . . . . .	15
1.5	Thesis Structure . . . . .	16
<b>2</b>	<b>Literature Review</b>	<b>17</b>
2.1	Overview . . . . .	18
2.2	Privacy Policy Analysis: Reduction of Legalization to Increase Transparency . . . . .	19
2.3	Compliance Checking with GDPR: Automating Regulatory Adherence . . . . .	20
2.4	Semantic Annotation and Legal Metadata Extraction: Creating Superior Legal Reasoning . . . . .	21
2.5	Challenges in Leveraging NLP for GDPR Compliance . . . . .	23
2.6	Future Directions for Research and Exploration . . . . .	24
<b>3</b>	<b>Methodology</b>	<b>26</b>

---

3.1	Introduction . . . . .	26
3.2	Dataset Acquisition . . . . .	27
3.2.1	OPP-115 Dataset . . . . .	28
3.2.2	ACL Coling Dataset . . . . .	30
3.3	Data Preprocessing . . . . .	32
3.4	Mapping GDPR Principles . . . . .	34
3.5	Granularity of Analysis . . . . .	35
3.5.1	Sentence Level Analysis . . . . .	36
3.5.2	Entire Policy Level Analysis . . . . .	37
3.6	Model Selection . . . . .	38
3.6.1	SBERT . . . . .	38
3.6.2	BERT . . . . .	39
3.6.3	GPT2 . . . . .	39
3.7	Multi-label Classification . . . . .	39
3.8	Compliance Reporting . . . . .	41
3.9	Performance Metrics . . . . .	43
<b>4</b>	<b>Implementation</b>	<b>45</b>
4.1	Data Preparation . . . . .	45
4.2	Model Training and Saving . . . . .	47
4.2.1	SBERT . . . . .	47
4.2.2	BERT . . . . .	51
4.2.3	GPT2 . . . . .	55
4.3	Testing on Unlabeled Data . . . . .	57
4.3.1	SBERT . . . . .	58
4.3.2	BERT . . . . .	60

---

4.3.3	GPT2 . . . . .	62
4.4	Single Principle Checker . . . . .	64
4.5	OpenAI Embedding Experiments . . . . .	66
4.6	Model Evaluation . . . . .	67
<b>5</b>	<b>Evaluation and Discussion</b>	<b>68</b>
5.1	Overview of Model Performance . . . . .	68
5.2	Analysis of Results . . . . .	75
5.3	Challenges Encountered . . . . .	85
5.4	Discussion on GPT-3 Embeddings . . . . .	89
5.4.1	Generating GPT-3 Embeddings . . . . .	89
5.5	Implications . . . . .	91
5.5.1	Addressing Research Questions . . . . .	92
5.5.2	Practicality of the Compliance Report . . . . .	93
5.5.3	Contributions to the Field of Data Privacy and Compliance	95
5.6	Future Work and Improvements . . . . .	95
5.6.1	Enhancing Model Performance . . . . .	96
5.6.2	Real-time Compliance Monitoring . . . . .	96
5.6.3	Broadening the Scope of Compliance Checks . . . . .	97
5.6.4	User-Centric Enhancements . . . . .	97
<b>6</b>	<b>Conclusion</b>	<b>98</b>
6.1	Summary of Findings . . . . .	98
6.2	Discussion . . . . .	100
6.3	Practicality . . . . .	101
6.4	Answering the Research Questions . . . . .	101
6.5	Contributions to the Field . . . . .	103

<b>CONTENTS</b>	<b>6</b>
-----------------	----------

---

<b>Bibliography cited</b>	<b>104</b>
---------------------------	------------

<b>A Full Policy Report</b>	<b>116</b>
-----------------------------	------------

# List of Tables

5.1	Model Evaluation Metrics . . . . .	75
5.2	Comparison of metrics for single principle before class resampling	76
5.3	Comparison of metrics for single principle after class resampling	77
5.4	Compliance analysis for privacy principles using NLP models .	77
5.5	Comparison of metrics for multi principle classification in sentence level . .	83
5.6	Comparison of metrics for multi principle classification in policy level . . .	85
5.7	Comparison of Legal vs. Base Model for BERT . . . . .	88
5.8	Threshold and Identified Principles from GPT-3 Embeddings . .	91

# List of Figures

3.1	Mapping of OPP-115 categories to GDPR article 5 principles . .	35
3.2	High level overview of the architecture . . . . .	37
5.1	Metrics for SBERT performance on sentence level . . . . .	69
5.2	Metrics for SBERT performance on policy level . . . . .	69
5.3	Metrics for BERT performance on sentence level . . . . .	70
5.4	Metrics for BERT performance on policy level . . . . .	71
5.5	Metrics for GPT2 performance on sentence level . . . . .	72
5.6	Metrics for GPT2 performance on policy level . . . . .	73
5.7	Metrics for OpenAI embeddings performance . . . . .	74
5.8	SBERT Compliance Report for Sentence Level . . . . .	83
5.9	BERT Compliance Report for Sentence Level . . . . .	84
5.10	GPT2 Compliance Report for Sentence Level . . . . .	84
5.11	Enter Caption . . . . .	86



## **Abstract**

As it stands in the contemporary world, compliance with regulations concerning data protection such as GDPR is central to organizations. The another important issue analysis identified is the fact that compliance is hampered by the fact that legal documents are often complex and that regulations are ever changing. This paper aims to describe the ways in which NLP aids in keeping GDPR compliance effortless through automated scanning for compliance, evaluating privacy policies, and increasing the level of transparency.

The work does not only limit to exploring the application of NLP for dealing with the privacy policies and facilitate better understandings of the third-party data sharing but also proceed to perform the experiments to evaluate the difference of several NLP models. They implement and execute the models to distinguish the one that performs the best based on the efficiency and speed at which it automates the process of compliance verification and analyzing the privacy policy.

Furthermore, some issues covered in the study refer to automated solutions to GDPR, such as generating of the machine readable models, which make the process of compliance evaluation more efficient.

Therefore, this paper emphasizes the importance of NLP to help organizations overcome the difficulties of GDPR compliance, create a roadmap to a more client-oriented data protection regime. In this regard, by comparing experiments done in the test and showing the performance of the better model, it helps enhance the measures taken in compliance and fosters the defense of individual rights in the cyberspace.

# Chapter 1

## Introduction

Given the fact that people produce and share many personal data in the today's digital environment, protection of individuals' privacy has become more relevant. Leading this effort is the General Data Protection Regulation (GDPR), a suitable legal instrument and body of laws the EU has enacted to mitigate the processing of the personal data and protect privacy and data protection rights of the citizens. Altogether, the GDPR solves the problem of protection of trust by imposing strict legislation on the organization that processes the personal data and regulates data processing activities by principles of transparency, accountability, and consent.

Some definitions of the keywords that will show up throughout the thesis paper are as follows:

- **Natural Language Processing (NLP):** A branch of computer science that deals with making computer applications understand and analyze written or spoken human language [1].
- **Data Processing Agreements (DPAs):** Legally binding documents to be

entered into between the controller and the processor in writing or electronic form [2].

- **GDPR (General Data Protection Regulation):** A comprehensive legislative framework enacted by the European Union (EU) to regulate the processing of personal data and uphold individuals' rights to privacy and data protection [3].
- **Privacy:** Empowering users to make their own decisions about who can process their data and for what purpose [4].
- **Compliance:** Meeting the requirements for properly handling personal data as defined in the law [5].
- **Method for compliance verification:** The process of ensuring that an organization is compliant with the GDPR.
- **Compliance Audit:** A systematic and independent assessment of an organization's compliance with the GDPR [6].
- **Text Mining:** Deriving high-quality information from text [7].
- **Transfer Learning:** A technique in machine learning (ML) in which knowledge learned from a task is re-used to boost performance on a related task [8].
- **BERT:** An "encoder-only" transformer architecture and a baseline model in NLP [9].

## 1.1 Background

In this paper, the Background section would give the information on why the General Data Protection Regulation (GDPR) came about. It is rooted back to some problems that took their genesis in middle of the twentieth century about the management of personal information. With the emergence of internet as one of the communication technologies the above concerns were realized fully. The collection, usage and sharing of personal data was a growing concern to the public due to data breaches that were instances in almost all organization [10].

In connection with the increase in the use of Internet applications and various online services, the volume of personal data processed also increased. This clearly created a major concern on the issues of privacy and security. They had been cases of leaking sensitive information and personal identification details being given to the wrong hands, information being used in the wrong way, and never being told how our information was being used [11]. This paper seeks to isolate and explain these challenges and show how the GDPR sought to overcome them by providing direction to organizations when handling personal data.

## 1.2 Significance of NLP in GDPR Compliance

The understanding the Roles of NLP More About GDPR section identifies the roles of Natural Language Processing in guaranteeing that organizations follow the GDPR policies. It begins with understanding that legal texts, such as privacy policies which are significant for every Internet user, can hardly be easily comprehensible as they contain numerous terms foreign to an ordinary person. Typically, compliance checking has long been done manually, and it is quite

laborious as well as involving chances of error.

NLP fulfills this need through the usage of intelligent algorithms that help to analyze and understand human language. Thus, using NLP methods, organizations can improve the effectiveness of their activities related to analyzing and understanding legal texts. This ranges from analysis of clauses that include them to the extraction of certain information and determination of compliance on various aspects related to GDPR. The application of NLP is beneficial not only in increasing the efficiency of the compliance process but also in enhancing precision and cohesion in assessing organizations' compliance with the regulations [12], [13].

Further, with the help of NLP, enterprises turn more capable to respond to the changes in the legislation on privacy protection. While working with NLP models legal needs can be fulfilled and as the legal interpretations change and new requirements are added the models can be trained again. This ensures that organizations adhere to the GDPR mandates throughout while at the same time keeping them informed on issues to do with data protection and privacy regulations. Finally, NLP enables organizations to address the challenges of GDPR compliance more effectively so that individuals' rights in the context of digital transformation could be preserved.

## 1.3 Problem Statement

The issues described in this section are connected with the problems that organizations face when attempting to adhere to the provisions of the GDPR and decipher ambiguous legal texts. It highlights typical conditions such as how to understand what rules really say, laws in development, and how to implement the

concept of privacy policies to non-legal audiences? The objective of this part of the thesis is to explain these problems so that by the end of this paper, we can attempt to answer the following questions.

First of all, it is crucial to note that dealing with the legal standards provided for in the GDPR framework is very problematic for many organizations that work in the context of digital environments. One of them is the complexity which arises from the necessity to comprehend the language used in legislation that might be hard to understand without legal background. It is with the situation that the GDPR has integrated many rules and regulations [14], which are sometimes hard to understand, let alone meet.

However, the GDPR is not something frozen in time; rather, when it gets implemented it changes with time, with the amendments, added clarification regulations and the new ones being made constantly. This adds yet another level of dynamism to the process, because the organization has to be informed of corresponding change and adjust its actions to the change [15].

Further, introducing data protection principles and users' privacy policies also pose a great challenge in terms of how these policy statements are relayed to the users. Most written documents, especially privacy policies, use legal terms that many people find hard to understand, and this make such people unsure of their rights and how the data they have will be used.

These challenges present the narrative of why solutions that are flexible, out-of-the-box, easy to use and which are able to produce GDPR stipulations in simple language, with little effort can go a long way in easing the comprehension of legal texts laid down under the GDPR as well as aiding in improving transparency in data processing practices. Solving these challenges is possible to help organizations protect people's rights to privacy and data protection and to properly address the

current requirements of the legislation.

## **1.4 Research Questions and Objectives**

In the following part of the paper, we present research questions, which guide our study, and objectives, which we aim at achieving. These questions help us to move towards a more profound analysis and recommendations in the cases of GDPR compliance and Natural Language Processing (NLP).

- How effective are NLP models, including GPT-3, in automating the identification of GDPR compliance issues within organizational data privacy policies?
- What are the limitations of current NLP technology in interpreting and enforcing GDPR compliance, and how can these limitations be addressed?
- What role can NLP-powered tools play in supporting compliance officers and legal experts in maintaining GDPR compliance?

The objectives are clear: to gain a deeper insight into quite the process of NLP in improving the compliance and to design more of the effective guidance tools for the organisations to follow the data protection laws. Thus, the goal is to conduct a profound research and, if possible, an investigation, to help people and organizations adhere to the rules established by the GDPR by expanding the knowledge base in the field of data privacy and compliance.

## 1.5 Thesis Structure

This part of the paper will endeavour explain how this entire thesis is presented and what is covered in each of the sections. It gives a general outlook of what to expect.

First, there is introduction, as it provides the premise for what is to follow, stating the general issue considered and its significance. Lastly, we proceed to the literature review part of the study in order to determine what others have done in this line. Next is the method where we discuss how we undertook the study after the literature review and theoretical formulation of the issues at hand.

The experimental results which we proceeded with are given in the next section of the work. At last, the conclusions are presented where the outcome of the research study, and implications are highlighted.



## Chapter 2

# Literature Review

The topic of compliance in the past few years, and more specifically in relation to the GDPR as an example of data protection legislation, one cannot fail to notice the seemingly ‘*revolutionary*’ changes that have been brought about by the incorporation of NLP capabilities [12], [16], [17]. As a result of this evolution, NLP has emerged into the upcoming area for automating diverse processes in accomplishing the data checking and privacy policies based on the GDPR and other data protection laws [18]–[20].

Combining NLP possibilities with legal compliance activities is a revolution in comprehending and dealing with regulatory requirements and legal documents. Typically, data protection compliance has been a time-consuming and expensive affair that has entailed paper-based sifting and analysis of vast legal texts [21]–[23]. Though there is not much information about its use in handling compliance, the use of NLP has brought efficiency and scalability in the organization’s use of Machine learning and natural language understanding to handle compliance work [24]–[26].

## 2.1 Overview

Thus, this extensive literature review aims at discussing key publications in this field and highlights various methodologies, models, and approaches that rely on NLP to analyze the complexity of GDPR compliance and other legal documents. Looking at the specifics of data protection and privacy legislations, data protection officers and academic pioneers have developed brand new approaches that concern the difficulties in corresponding compliance procedures and the clarification of data practices in policy documents such as privacy policies, data processing agreements, and regulation requirements [27]–[29].

Key areas of focus within this literature review include:

- **Compliance Checking:** There are well thought and advanced models and frameworks that describe how compliance checking of the GDPR requirements can be automated. These models act as legal advisors that employ NLP strategies to analyze legal texts and determine whether an organization conforms to the requirements of legal precedents [3], [30], [31].
- **Privacy Policy Analysis:** Corpora as well as tools like PrivaSeer have emerged as powerful NLP solutions of analyzing policies and have greatly helped in large scale collection of data extraction as well as classification [32]–[34]. These performances may be pursued to increase or improve the current state of transparency together with improving the ability of the users in making proper choices regarding their right to privacy and protection of their data.
- **Semantic Annotation:** Previous attempts at indexing and annotating legal texts with semantics have made it possible to achieve enhanced perform

for search as effectively as for details mining [35]–[37]. Through adding metadata and semantic tags into legal texts, the scholars have created opportunities for practically working in the future to improve the availability and understanding of legal regulations.

Thus, it is within this framework of utilizing these pioneers in the field of NLP that this literature review aims at uncovering the possibility of the NLP methodology in the GDPR context and data protection. Thus, by integrating various methodologies and approaches, the scholars try to enhance innovation, increase the roles of methods' transparency, and help organizations to manage the challenges of the constantly changing legal environment effectively [21], [28], [38].

## 2.2 Privacy Policy Analysis: Reduction of Legalization to Increase Transparency

A privacy policy is one of the critical legal instruments that provide information about an organization's data management from numerous users. Nevertheless, such policies are written in legal jargon, and as such, users struggle to understand what the policies state. In response to this problem, scholars have adopted unique methods based on NLP to work on the automation of privacy policy evaluation in order to increase transparency [25], [39], [40].

[25] presents an innovative project known as PrivaSeer which is focused on the issue of complexity of the privacy policies and the possible ways to help users to easier and better comprehend them. About one million English language website privacy policies are included in the dataset referred as PrivaSeer, which is the

## **2.3 Compliance Checking with GDPR: Automating Regulatory Adherence**

---

largest dataset of this nature. Corpus accumulation is done through elaborate web crawling and filtering processes that comprise the creation pipeline and guarantee the corpus's comprehensiveness and relevance. As it offers the researchers the list of links to the privacy policies that contain both strengths and weakness for comparison, information analysis with the help of PrivaSeer and identification of areas for the improvement of the privacy policy formulation can be effectively achieved.

In the same manner, [33] deals with the identification and categorization of third-party entities from privacy policies that are connected with mobile applications. In this regard, the study acknowledges the need for expounding data recipients in relation to user's perception about the apps' privacy. Employing state-of-art Named Entity Recognition (NER) models, the research seeks to improve users' awareness of third-party entities present in the privacy policies. Accomplishing feature-specific annotations and detailed models, this study aims at producing effective and efficient tools for extracting a proper information set from privacy policies and thus enabling the stakeholders make rational decisions on their data safeguarding.

## **2.3 Compliance Checking with GDPR: Automating Regulatory Adherence**

The GDPR remains a cause of immense concern to organizations especially in terms of the level of compliance to the data protection laws. Due to the complexity and activity of GDPR mandates, the concept of checking has had to evolve, and one of the potential solutions includes Natural Language Processing

[20], [31], [41].

[20] follows a model driven engineering approach to propose a conceptual model that can be analyzed by a machine for capturing the characteristics of DPAs under GDPR. Thus in enhancing the specific criteria for the evaluation of DPA compliance and suggesting the automated tool for checking the compliance, the research aims to relieve the compliance pressure on data controllers and processors. The DPAs have therefore derived a conceptual model that forms the basis of assessing the compliance of the organizations with GDPR to help in the early detection and prevention of compliance issues.

In the same vein, [18] discusses the SPECIAL H2020 project through which tools enabling organizations to perform automated compliance checks will be provided. This is based on the key theme of the project which is the coming up with of a policy language that will be capable of expressing consent, business policies as well as regulatory obligations in such a way that a machine will be able to understand them. The project provides two essentially different variants of automated compliance verification; that allows offering the organizations flexible opportunities to check their GDPR compliance. The SPECIAL H2020 is an innovative project involving NLP-automated solutions aimed at improving compliance with requirements and standards in various fields.

## 2.4 Semantic Annotation and Legal Metadata Extraction: Creating Superior Legal Reasoning

Semantic annotation of legal texts is a crucial prerequisite in implementing such elements of Searcher's intent, supporting the process of legal argumentation

and interpretative scalability. These notable advancements in this area make this paper [42] an important source of information since it provides a mixed-method approach in extracting semantic legal metadata with the help of Natural Language Processing (NLP) [17], [28], [43]. Thus, by suggesting the uniform conceptual model of semantic metadata types related to the analysis of legal requirements and the list of automated extraction rules based on the constituency and dependency parsing, the research provides the basis for the systematic analysis of the legal provisions. This way, the so-called conceptual model implies the differentiation of phrase-level and statement-level concepts, thus rounding off the framework for the semantics of legal statements. The study done with tokenization, sentence splitting, POS tagging, NER, and both constituency and dependency parsers shows the fine extraction of Semantic legal metadata. This way, the validation of the extraction rules has merely positive results, while precision ranges from 87.4% to 97.2%, recall ranges from 85.5% and 94.9%, which averts to the efficiency of the proposed approach in achieving meager legal semantics [43].

Likewise, [17], [44] offers a work on the Pragmatic Approach for Semantic Annotation and the design methodology of CLAL in the XML context with a proper formalization. In this respect, the study refers to the improvement of legal provisions by semantically annotating texts, which merely adds interpretation-neutral information as metadata to the legal texts and thus enhances the search as well as a more effective legal analysis of the provisions. CLAL plays the role of annotation language, and it is the common language used for legal texts' annotation and involves quite a number of semantic characteristics that are indispensable for legal reasoning and understanding. This is coordinated via measures of inter annotator agreement which is used to confirm the reliability and security of the annotation language in the course of the process of annotation [17]. This

annotated corpus combined with outlines of the CLAL schema can be useful for the development of further research and for the practical usage aimed at search and analysis of legal texts with enhanced criteria for modern studying.

## 2.5 Challenges in Leveraging NLP for GDPR Compliance

Nevertheless, there are some issues that should be noted, proving that further research and development in the field of NLP for GDPR, as well as legal text analysis in general, should continue. These issues must be further solved for the purpose of providing dedicated solutions that are reliable and can easily address various aspects related to regulation and legal documentation [12], [20], [45].

**1. Dataset Limitations:** One of the issues is the fact that these models need to be currently trained on larger and more diverse datasets for best results – hence accessibility. Currently produced datasets may not be sufficient in providing the needed scope and depth in coverage to reflect the variations of the legal language use and the associated compliance situations. There is, therefore, a dire need to assemble relevant collections of overall legal works, legal codes, and linguistic differences [21], [28], [46].

**2. Adaptability to Evolving Legal Frameworks:** The legal frameworks in case of NLP-based compliance solutions are GDPR and other data protection laws which emerge as the biggest problem due to how volatile they are. Legal texts change often, get interpolated, amended, or simply reinterpreted quite regularly, which requires models to be at par with these changes as and when they happen. NLP models that are flexible and dynamic enough to follow any changes to the

legal specifications for compliance are relevant are needed [36], [47], [48].

**3. Exploration of Newer NLP Architectures:** Although the current NLP architectures have shown a good level of achievement in multiple utilities such as compliance checks and legal textual analysis, researchers are still trying to discover new architectures or pre-trained models. The Generative models including GPT-3 and T5 can be used to produce natural responses and also for context-based legal documents comprehension. It would be interesting to try out these neural architectures and incorporate these architectures into tasks that deal with compliance; there could be much more potential for enhanced accuracy and efficiency here [21], [49], [50].

## 2.6 Future Directions for Research and Exploration

On the basis of these research gaps and the limitations mentioned in the present work to further develop the NLP for GDPR compliance field, the following strategies can be outlined for future research:

**1. Experimentation with Generative Models:** The method of using generative models should be extended in future studies with up-to-date generative models such as GPT-3 and T5. These models have enhanced abilities in NL processing and generation that are especially appropriate for complex applications such as compliance, legal document review, and similar processes [24], [25], [37].

**2. Extensive Evaluations:** It is critical to conduct more comprehensive investigations with respect to NLP-based compliance solutions and applications as to how well they work, how resilient they are, and how easily they can be scaled up. The evaluation framework should contain multiple datasets, metrics, and applications to give recommendations for utilizing the approach and information



about its effectiveness [3], [18], [45].

**3. Interdisciplinary Approaches:** Another possible method for further research involve the focus on both legal knowledge incorporated into the engineering of compliance solutions and NLP methods to improve such solutions. Multidisciplinary collaborations with law enforcement, IT, compliance, and NLP specialists can work towards improving the current solutions by creating tools that are more aware of the specific issues within the complex sphere of compliance [27], [39], [43].

If researchers employ these challenges and future directions towards enhancing NLP for GDPR compliance, they will be able to progress this field in terms of making potent, expandable, and responsive models which can help organisations in gaining the skills and tools to cope with these unpredictable easy guidance services.

In summary, the literature reviewed supports the centrality of NLP for the complex problems concerning GDPR regulation, privacy policy comprehension, and legal documents deciphering. Further research and experimentation have to be conducted to fill this gap and enhance the state of the art in the applied approaches for automated compliance checking and legal text analysis.

# Chapter 3

## Methodology

### 3.1 Introduction

This chapter explains how the methods of this thesis were used to detect automated compliance with General Data Protection Regulation (GDPR) using Natural Language Processing (NLP) techniques. The fact that legal texts are specific and legal standards change quite often makes the systematic approach of developing, training, and evaluating the NLP models relevant. In this study, two different data sets were used: OPP-115 and ACL Coling. They were utilized to train and verify the performance of a couple of NLP models, including SBERT, BERT, and GPT2, which are one of the most powerful languages processing technologies.

First in the methodology we begin with getting and pre-processing the dataset that the models will be trained on. Next, there is a thorough process of mapping the categories from OPP-115 dataset to principles in the GDPR 5. This will later on assist in the training the models for multi-labelled classification. Furthermore, we analyze different granularities to pursue. There will be analysis both at the

sentence and whole policy levels to find out the effectiveness of the models in different situations.

Also in this methodology there will be a comprehensive explanation of the capabilities and limitations of the chosen models used in the detection of compliance with the GDPR. Specifically, the pros and cons of the mentioned models are examined thoroughly as each emphasize their own unique capabilities and reason for being selected. So I will highlight their suitability for the task.

The core of the methodology lies in a multi-label classification training process for GDPR principles, the production of compliance reports, and the evaluation metrics specifically created to assess models' performance and precision. Finally, I explain the fine-tuning techniques taken to improve the performance of the models.

Hence, in the following sections, I will outline both the technical and theoretical approaches used to achieve the research objectives shown in chapter 1.

## 3.2 Dataset Acquisition

I made use of two datasets gotten from [Usable Privacy Policies \[51\]](#). The combination of the two datasets allows for more reliable results from training and testing the models. While one of them allows a targeted model training as it has privacy policies that are properly annotated, the other dataset challenges the models. With no targets or annotation, I was able to test how adaptable and generalizable the models are.

This way of selecting the datasets was strategic and it ensures that the trained models don't specialize only in checking GDPR compliance in the specifically defined boundaries. But that they are also capable of performing well in other

legal contexts.

### 3.2.1 OPP-115 Dataset

OPP-115 which comes from Online Privacy Policy Project is a dataset of 115 policy policies gathered in 2016 from various websites. It provided a good starting point for the training on NLP models on GDPR compliance. The policies in this dataset are all annotated in detail by a set of defined categories that do align with the principles provided under the GDPR article 5. Each policy was cleaned and has a pretty print in html format and it's respective annotation in csv format. They were all looked at and labeled with annotations that characterize specific data elements. Hence this dataset is very good for training models to understand privacy policy nuances.

To dive deeper, a brief description of OPP-115 categories [52] is as follows:

1. First Party Collection/Use: What, why and how information is collected by the service provider
2. Third Party Sharing/Collection: What, why and how information shared with or collected by third parties
3. Data Security: Protection measures for user information
4. Data Retention: How long user information will be stored
5. User Choice/Control: Control options available to users
6. User Access, Edit and Deletion: If/how users can access, edit or delete information

7. Policy Change: Informing users if policy information has been changed
8. International and Specific Audiences: Practices pertaining to a specific group of users
9. Other: General text, contact information or practices not covered by other categories.

A sample of how an annotation looks like can be seen below. The text highlighted in yellow shows the data type under the category and the text highlighted in blue shows the extracted sentence that pertains to the category.

#### Sample 3.2.1: Snippet of annotated data from OPP-115 dataset

Third Party Sharing/Collection

```
{ "Third Party Entity": "endIndexInSegment": 256, "startIndexInSegment": 0, "selectedText":  
"We reserve the right to share information in order to investigate, prevent, or take action regard-  
ing illegal activities, suspected fraud, violations of the Geekdo Terms of Service, situations  
involving potential threats to the physical safety of any person", "value": "Unnamed third  
party", ... }
```

<https://boardgamegeek.com/privacy>

It is visible for the category "Third Party Sharing/Collection" the start and end index for the sentence extracted from the entire policy. Because of its structured format, the detail of annotations and labels that can be used in training, it was a given to go ahead with the OPP-115 dataset.

### 3.2.2 ACL Coling Dataset

The ACL Coling dataset [53] is chosen to be used mainly for the evaluation purposes. It has a different range of 1,010 legal documents collected in 2014 in xml format. The corpus itself was made for the Computational Linguistics Conference and it is made up of different types of texts. It acts as the perfect collection of policies for testing the strength of the models trained on the OPP-115 dataset that is more structured, to generalize and be resilient to unseen data. A sample of how the data looks like is shown below:

#### Sample 3.2.2: Sample of data structure from ACL-Coling dataset

```
<POLICY      modification_date="March      01,      2013"      pol-
icy_url="http://www.zendesk.com/company/privacy"      website_category="Computers"
website_index="096" website_url="zendesk.com">
```

```
<SECTION>
```

```
<SUBTITLE />
```

```
<SUBTEXT>Privacy Policy
```

Effective as of March 1, 2013. For the prior version of our Privacy Policy, click here.

At Zendesk, we respect and protect the privacy of visitors to our website, [www.zendesk.com](http://www.zendesk.com) (together with the other websites we own and control, the “Zendesk Websites”), and our customers who use our on-demand customer service support platform, tools and services offered on the Zendesk Websites (together with the Zendesk Websites, the “Service”). This Privacy Policy (“Policy”) explains how we collect and use visitors’ and customers’ information, particularly personal information, as part of the Service. The information Zendesk collects and uses is limited to the purpose for which customers engage Zendesk and other purposes expressly described in this Policy. Any discussion of your use of the Service in this Policy is meant to include your visits and other interactions with the Zendesk Websites, whether or

not you are a user of Zendesk’s on-demand customer service support platform.

</SUBTEXT>

</SECTION>

<SECTION>

<SUBTITLE>Privacy Certifications </SUBTITLE>

<SUBTEXT> </SUBTEXT>

</SECTION>

</POLICY>

From the figure, 3.2.2, the text highlighted in pink is the root of the given policy, inside it are sections, highlighted in orange. Each section has a subtitle, highlighted in green, and a subtext, highlighted in purple.

Attribute	OPP-115	ACL Coling
<b>Focus</b>	Privacy Policy Analysis	Computational Linguistics Research
<b>Content Type</b>	Website Privacy Policies	Website Privacy Policies
<b>Annotations</b>	Data Practices (e.g., collection, sharing)	None
<b>Use Case</b>	Training models to analyze privacy policies	Broad NLP and linguistic research
<b>Data Format</b>	Annotated Text	Plain Text, PDF
<b>Accessibility</b>	Restricted Access	Open Access
<b>Size</b>	115 policies	1,010
<b>Year of Release</b>	2016	2013 & 2014

---

Attribute	OPP-115			ACL Coling
Maintained by	Usable	Privacy	Policy	Association for Computa-
	Project			tional Linguistics

---

### 3.3 Data Preprocessing

For any text analysis how effective the data preprocessing goes a long way which affects the performance of the models. Especially in legal texts, it is very important as text representation accuracy plays a significant role. This section discusses the preprocessing techniques implemented on both the OPP-115 and ACL Coling datasets to make them suitable for continuing the experiments.

The OPP-115 dataset required extensive text cleaning especially because to get the policy with its equivalent label, we needed to go through the annotated version and put sentences together to get the entire privacy policy. The preprocessing steps included:

- **Text Cleaning:** Essentially removing some leftover HTML tags and also non-alphanumeric characters. Here we want to get rid of elements that could potentially affect the text analysis.
- **Tokenization:** Breaking down the text into tokens. For natural language, this step helps parse text into a form that the model can better process.
- **Lemmatization:** Shortening words to their root form. Lemmatization tends to preserve the context and this is vital for maintaining the semantic integrity of the policies.



The ACL Coling dataset, requires a slightly different approach to preprocessing. Given the nature of XML structure, the process begins with parsing the XML documents to extract properly the textual content. So code was written to identify and read the specific elements highlighted in 3.2.2 in the XML hierarchy. These are the section subtexts, because they are where the relevant legal text is stored.

Once the text is extracted, it undergoes some further cleaning steps:

- **Normalization:** All text is converted to lowercase as it helps to reduce how complex the rest of the processing will be.
- **Whitespace and Newline Handling:** If there are extra whitespaces or newline characters, these are removed. It prevents any error in parsing that could occur during the tokenization phase.
- **Tokenization:** Then text is split into tokens to turn the text into a form friendlier for the model to process.
- **Lemmatization:** Each token is reduced to its base form at the same time maintaining the semantic meaning of the text.
- **Punctuation Removal:** Punctuation, are removed to focus on the words themselves. If not they could add noise into NLP models.

These preprocessing steps are done in order to maximize how efficient the NLP models used are in this experiment. And that is achieved by providing clean, consistent, and meaningful text data.

## 3.4 Mapping GDPR Principles

The mapping of GDPR principles to the categories of the OPP-115 dataset is the next step in the methodology. Because we want to classify policies as compliant or not based on if they follow the GDPR principles. This section outlines how and why each category from the OPP-115 dataset map to corresponding GDPR principles from Article 5.

The OPP-115 dataset categorizes listed in [3.2](#) map to the GDPR principles as demonstrated by [54]. The mapping is as follows:

- **First Party Collection/Use:** This category is mapped to Lawfulness, Fairness, Transparency; Purpose Limitation; and Data Minimization.
- **Third Party Sharing/Collection:** Similarly, this involves Lawfulness, Fairness, Transparency; Purpose Limitation; and Data Minimization.
- **User Choice/Control:** Mapped to Lawfulness, Fairness, Transparency.
- **User Access, Edit, and Deletion:** Linked to Lawfulness, Fairness, Transparency; and Accuracy.
- **Data Retention:** Corresponds to Storage Limitation.
- **Data Security:** Maps to Integrity and Confidentiality.
- **Policy Change:** Also mapped to Lawfulness, Fairness, Transparency.
- **International and Specific Audiences:** Maps to Lawfulness, Fairness, Transparency.
- **Do Not Track and Other:** These categories do not have a direct mapping to specific GDPR principles.

This mapping process not only guides the training of NLP models but also helps in structuring the compliance checks. We can also see this mapping visually as shown below:

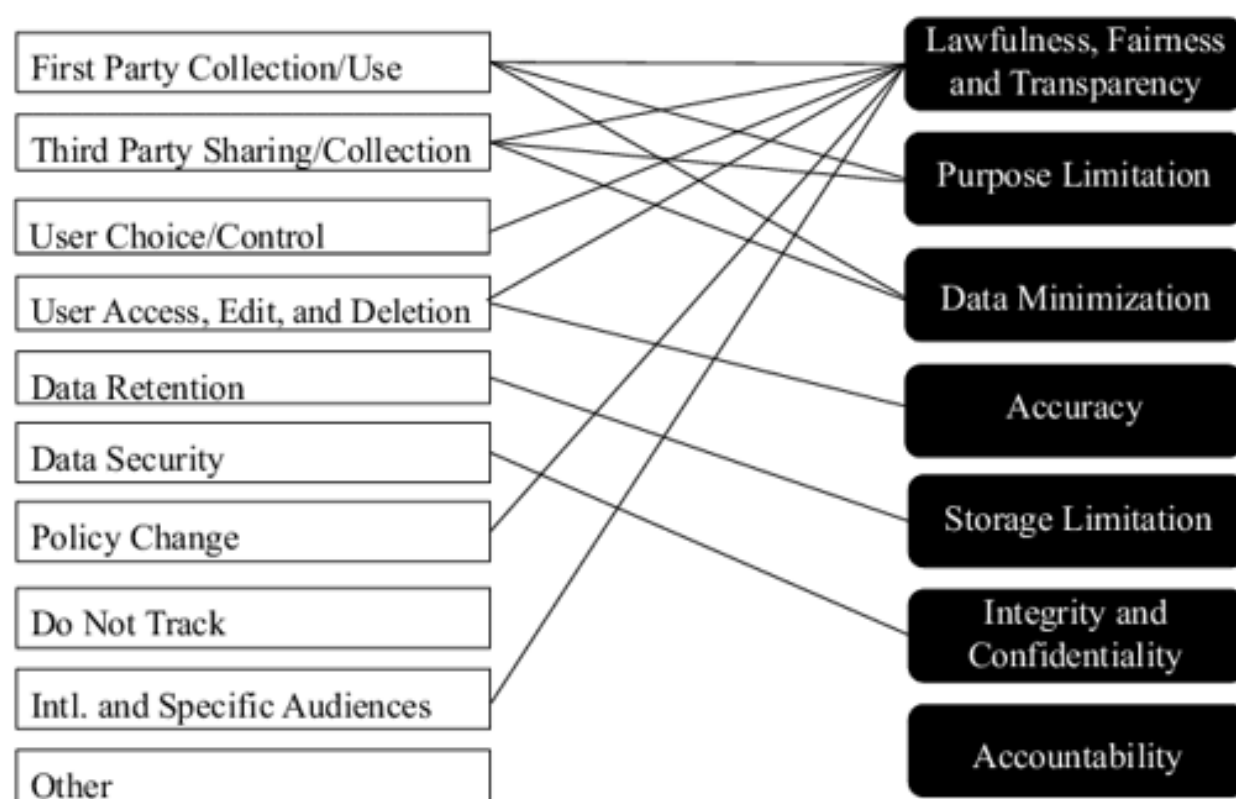


Fig. 3.1. Mapping of OPP-115 categories to GDPR article 5 principles

## 3.5 Granularity of Analysis

In the context of applying NLP to GDPR compliance, the granularity of the analysis is important as we apply different models. That is why we use and compare two levels of granularity: sentence level and entire policy level. Each level has its own benefits and challenges as we will cover now.

### 3.5.1 Sentence Level Analysis

**Benefits:**

- **Precision in Detection:** The identification of compliance-related issues at the sentence level is slightly more accurate than the identification at the document level or at the word level. It allows the model to concentrate on specific phrases and clauses in the text that relate to GDPR principles thus it can be very useful for compliance screening purposes.
- **Contextual Relevance:** Moving the analysis to the sentence level ensures that the surrounding contextual co-text is considered and that however the Anselmian relation is working, it is done so on a per-sentence basis.

**Challenges:**

- **Loss of Broader Context:** While considering small sections of text, the reader gets highly specific and effective results; However, the strategy fails to consider the general idea of the policy, which could lead to overlooking certain patterns or goals that came through a number of sentences or sections.
- **Fragmentation:** This granularity may cause problems in the form of fragmentation where the analysis is done on a mere sentence basis, they may be limited ideas on the overall contextual layout of the policy.

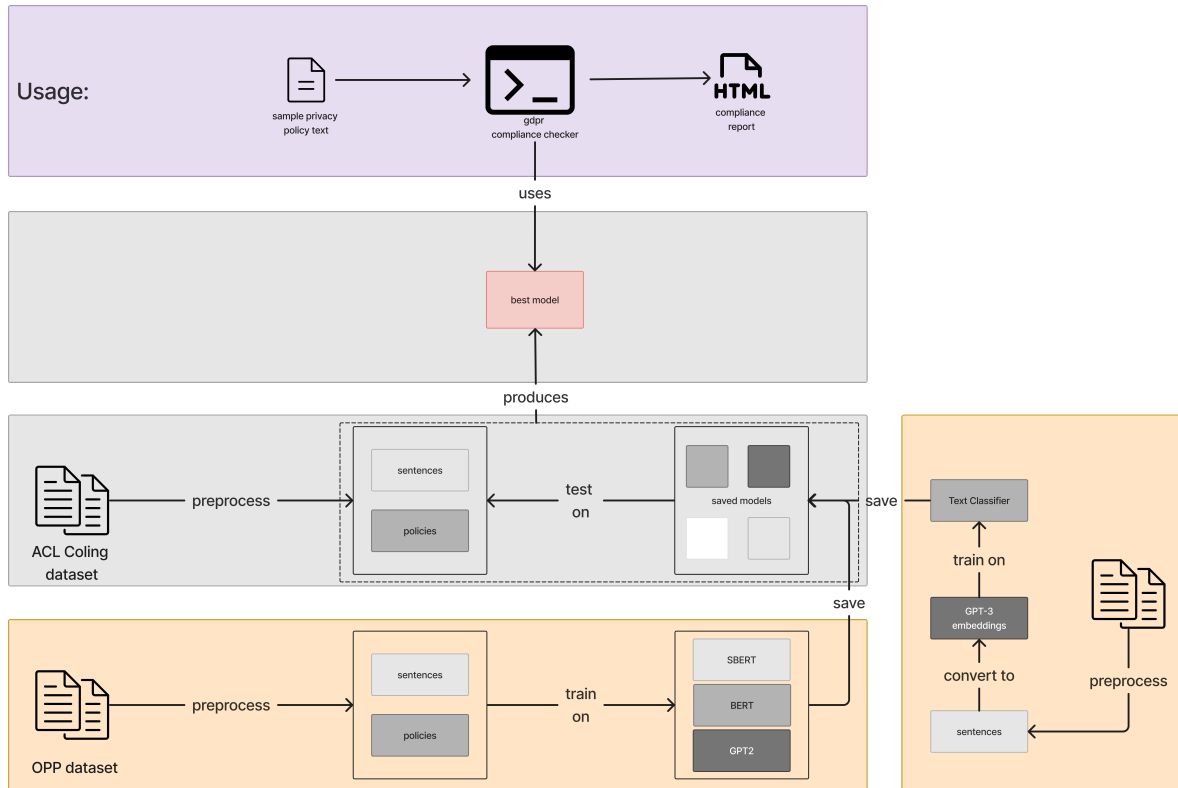


Fig. 3.2. High level overview of the architecture

### 3.5.2 Entire Policy Level Analysis

#### Benefits:

- **Holistic Understanding:** Evaluating a policy as a single text allows to consider the whole idea of the document and compliance of the policy to the requirements of GDPR. Especially helpful is the outline of general compliance with the rules of data protection.
- **Efficiency:** However, in large data sets, it is more convenient to assess entire documents, albeit to a less extent than in a more detailed option of analyzing the results of the work on a sentence level.

**Challenges:**

- **Potential for Overgeneralization:** Such an approach of analysis might result in overlooking particular non-conforming words, phrases or even paragraphs in the entire document while on the surface, seems to be absolutely compliant.
- **Complexity in Handling Large Texts:** Full policies, particularly long and complicated ones, are tricky to dissect at once from a computational standpoint and due to the model's precision and stability when faced with large texts.

This way, the intention is to have the benefits of both and subsequently, go for the level that has the best outcome. This will increase the effectiveness of the checking process giving more reliable results on compliance checks.

## 3.6 Model Selection

The selection of right NLP models is the most important step for the success of any machine learning task. Especially in the field of legal text analysis where accuracy and underlying tones of interpretation matter. We look at three different models, and now each of them are examined closely to understand why they were chosen.

### 3.6.1 SBERT

SBERT performs well when getting the sentence-level embeddings which are also more semantically valid than ordinary BERT. This is particularly useful for some exercises such as sentence similarity and clustering; therefore, is

very relevant in classifying differences in GDPR compliance in various privacy policies.

### **3.6.2 BERT**

One of the biggest strengths of BERT is that it is proficient in capturing language context and structural relations hence performs exceedingly well in relation recognition assignments. What makes BERT suitable for assessing and deriving exact, legally related information from intricate records such as the policies of privacy are.

### **3.6.3 GPT2**

It is worth noticing that GPT2 is built mainly for creating text that is both grammatically correct and contextually appropriate. It is beneficial for the activities that imply the creation of content, like privacy policies or GDPR-based explanatory texts, which improve users' interaction with compliance tools through the automated, easily understandable responses.

## **3.7 Multi-label Classification**

For this thesis, the multi-label classification method was employed to define the measure of noncompliance of each chunk of text, single and numerous SENTs, or the entire POLICY with the seven principles outlined under Article 5 of the GDPR.

Nevertheless, within the framework of GDPR, every single statement or even the entire policy is assessed not only in terms of compliance or non-compliance

but is considered as to how many principles of the GDPR are simultaneously violated by it. This is a real-world situation of a single privacy policy seeming to address some of the GDPR principles and disregard the others. Hence, each text unit is connected with a numeric vector of length seven, which captures the admittance to the seven principles of the GDPR. Every component of the vector is binary, equal to either 0 or 1; small number 1 means compliance with the specific GDPR principle while number 0 refers to violations.

**Labels** The labels for this classification task are derived from the GDPR principles [55], which are:

- Lawfulness, Fairness, and Transparency
- Purpose Limitation
- Data Minimization
- Accuracy
- Storage Limitation
- Integrity and Confidentiality
- Accountability

Each of these principles forms a dimension in the label vector for each data point. A sentence or policy is annotated with '1' for a principle if it aligns with the GDPR requirements for that principle, and '0' otherwise.



## 3.8 Compliance Reporting

Effective GDPR compliance requires not only identifying potential non-compliance but also clearly reporting these findings. This section outlines the process used in this thesis to transform model outputs into comprehensive compliance reports.

The compliance report generation process is as follows:

- **Interpretation of Model Outputs:** All the NLP models described in this thesis provide probabilities for all GDPR principles regarding if the given text, which may be a sentence or a policy, complies with those principles. These probabilities are, therefore, transformed to binary states, either compliant or non-compliant depending on the defined probability limit. For example, if the probability that the text meets the first ‘data minimisation’ principle is higher than zero. 8, it is categorized under compliant. This threshold was optimized, as was mentioned in the previous section concerning the discussion of the precision-recall curve.
- **Aggregation of Results:** After classifying individual sentences or sections, the last step is to combine these results from the classification of individual sentences or sections of the policy to provide an overall policy level compliance report. This encompasses a summary of all the result of checked GDPR principles with the overall compliance picture being presented. Whenever there is a violation with any of the principles stipulated, then the whole document is marked indicating the areas to be considered as non-compliant with the principles.
- **Report Format:** The final compliance report is structured to provide clear

and actionable insights. Each report includes:

- **Summary Section:** An overview indicating overall compliance status.
- **Detailed Findings:** Example excerpts from the text highlighting compliant areas.
- **Recommendations:** Based on the non-compliance findings, advice on which principles need to be added to met GDPR standards.

### Sample 3.8.1: Example of a Compliance Report

GDPR Compliance Report for Policy XYZ

Summary: Non-compliant with 2 out of 7 principles evaluated.

Detailed Findings:

1. Data Minimisation: Compliant.
2. Integrity and Confidentiality: Non-compliant.

Example: "User data may be stored indefinitely for analytics."

Recommendations:

- Review the data retention policy to align with the 'storage limitation' principle.

Thus, the presented compliance reporting framework guarantees that the results of the models are not only understandable but also practical, which makes them significant for organizations that need proper GDPR compliance.

## 3.9 Performance Metrics

One of the most important questions which have to be solved in order to ensure proper working and reliability of the NLP model, and which is considered to be particularly important at the tasks of compliance with GDPR – as this issue is rather critical – is a correct evaluation of the latter. This section details the metrics used to assess the performance of the models deployed in this thesis: As discussed the performance metrics includes: loss, accuracy, precision, recall, and F1-score. All of them are selected according to the opinion that the given measures can provide a better understanding of the areas of different models and their potential.

The models are evaluated using the following performance metrics: The models are evaluated using the following performance metrics:

- **Loss:** During training, it is important to track the loss function for better understanding of the model's learning and convergence. Talking about the specifics of this thesis, general loss metrics are essential while adjusting hyperparameters of the model architectures. Similar to this, reduction in loss is directly proportional to the improvement of the model's perception of GDPR specifics.
- **Accuracy:** After the datasets preparation, accuracy is useful to measure the general performance of the model in training. In each of them, accuracy is useful to assess the performance enhancement with the increase in the epochs and to compare simple models with the enhanced versions incorporated later.
- **Precision:** In the legal field, leaving a compliant text in the non-compliant

category entails a cost known as false positive for which organizations face reviews and audits. Therefore, precision was particularly focused on in the model training routines. And the changes such as class weighting and the introduction of regularization could improve the model's precision while maintaining other characteristics.

- **Recall:** Recall should be very high in order to capture all possible non-compliant issues. In my experiments, recall is a measure of interest especially when measuring the model's performance on unseen test sets. It is essential to guarantee the models are accurate enough in detecting as many instances as possible, so they do not leave out a lot of real positive cases.
- **F1-Score:** Due to the imbalance of classes after preprocessing steps such as stop words removal and stemming, the use of F1-score was inevitable. It gives a clear comparison of each model's precision and recall. This assisted in establishing which model was the best performer of all the models without overfitting the majority class.

Combined, all these metrics guarantee that the models developed for the specific use of checking for GDPR compliance are tested and validated which in return provides confidence in using the models for automated compliance checking.

# Chapter 4

## Implementation

This chapter describes the implementation details of the machine learning models used in this thesis, including SBERT, BERT, and GPT2. Each model was trained on both the sentence level and entire policy level using the labeled OPP-115 dataset. The models were then saved and subsequently tested on unlabeled policies from the ACL Coling dataset to assess their generalizability and performance. Additionally, the implementation included experimenting with GPT-3 embeddings followed by training these embeddings on a classifier.

### 4.1 Data Preparation

At first, sentences and the whole policies themselves were annotated with binary vectors which can describe the compliance concerning seven GDPR principles, which allowed for the separate models' training for the sentence- and the policy-level analysis. All the data passed through tokenization and encoding to fit the given model's needs; for example, SBERT, BERT, and GPT2 required different data formats to work on.

**Label Assignment and Distribution** Notably, each sentence or policy was described with a binary 7-vector where each component is 1 at the position relating to the corresponding GDPR principle. Each of them is brought to the vector; if the value of the element is 1, then the organization adheres to the corresponding principle, and if it is 0, then the organization does not adhere to this principle. The distribution of these labels across the dataset is as follows: The distribution of these labels across the dataset is as follows:

- **GDPR Principle 1 (Lawfulness, Fairness, and Transparency):** 8460 sentence instances and 115 policy instances marked as compliant.
- **GDPR Principle 2 (Purpose Limitation):** 6209 sentence instances and 115 policy instances marked as compliant.
- **GDPR Principle 3 (Data Minimization):** 6209 sentence instances and 115 policy instances marked as compliant.
- **GDPR Principle 4 (Accuracy):** 646 sentence instances and 90 policy instances marked as compliant.
- **GDPR Principle 5 (Storage Limitation):** 396 sentence instances 76 policy instances marked as compliant.
- **GDPR Principle 6 (Integrity and Confidentiality):** 1000 sentence instances 102 policy instances marked as compliant.
- **GDPR Principle 7 (Accountability):** 0 sentence and policy instances marked as compliant.

This label distribution shows that the principles are divided into subtopics of interest within the data and that it may be a problem to train models, particularly

in principles with relatively few positive samples. Therefore, the difference in the label frequency highlights the fact that models should be able to detect infrequent but essential compliance issues, including Accuracy and Storage Limitation.

Before training, the OPP-115 dataset was prepared at two different granularities:

1. **Sentence Level:** In the case of the privacy policies, each sentence was annotated with regard to how well it adhered to the presented GDPR principles, which comprised the training data for the sentence-level models.
2. **Entire Policy Level:** Full privacy policies were categorized based on the GDPR compliance, even or odd, as a whole to train the policy-level models.

## 4.2 Model Training and Saving

As part of model training, it was necessary to set up structures for neural networks since the documents' embeddings are preprocessed using SBERT, BERT, and GPT2 models. All of the above models were trained using OPP-115 dataset split in training and validation sets After training all the models were serialized in a standard way.

### 4.2.1 SBERT

#### Sentence Level

Hence, the SBERT model version [all-MiniLM-L6-v2](#) which was adopted to provide semantically matching embeddings of the sentences. This capability is important hence for reporting GDPR compliance at the level of sentences.

Firstly, with the help of SBERT, embeddings were generated from the OPP-115 dataset, which were then taken in by a classifier as features. The embeddings also retain the entire background about the sentence and its suitability to high performance of classification.

Further, to manage these embeddings as well as their associated GDPR compliance labels which are in form of a 7-element binary vector for each sentence, a new PyTorch custom dataset, named as `SBERTEmbeddingDataset` was designed. Here with this dataset during the actual training process data loading and batching becomes easier.

**Listing 4.1:** Code snippet of the custom dataset.

---

```
1 sentence_sbert_model = SentenceTransformer('all-MiniLM-L6-v2')
2
3 def sentence_sbert_generate_embeddings(texts):
4     embeddings = sentence_sbert_model.encode(texts, convert_to_tensor=True)
5     return embeddings
6
7 class SBERTEmbeddingDataset(Dataset):
8     def __init__(self, texts, labels):
9         self.embeddings = sentence_sbert_generate_embeddings(texts)
10        self.labels = labels
11
12    def __len__(self):
13        return len(self.embeddings)
14
15    def __getitem__(self, idx):
16        return {
17            'embeddings': self.embeddings[idx],
18            'labels': torch.tensor(self.labels[idx], dtype=torch.float)
19        }
```

---



As for the training, it used to establish the neural network classifier that includes LSTM and the attention layers that would allow using the sentence embeddings destined for the multi-label classification. The training process was controlled by a loop where each iteration called epoch entailed training on a training set, then on a validation set.

The following metrics were used to assess the effectiveness of the model's performance. These include accuracy, which is the ratio of the number of correct predictions to the number of samples, and precision, recall and F1 score, which, particularly in multi-label sample domains where the compliance with each principle is considered as a unique label, are more insightful measures of a model's performance.

Subsequently, the model was saved using PyTorch save mechanism where the parameters of the model are stored with a view of replicating the model and performance of additional tests depending on the model, or in the creation of a production model.

Such a detailed implementation of the SBERT model at the level of individual sentences guarantees that each sentence will be recognized as compliant with GDPR with the help of sufficiently effective NLP methods concerning the analysis of legal texts.

### Entire Policy Level

At the policy level, the same Sentence-BERT model, that is the [all-!\[\]\(003082e50e3009141f59bd5df831749f\_img.jpg\)MiniLM-L6-v2](#) modification, was used to process raw texts of privacy policies. This approach was conceived to check the GDPR compliance of complete documents and hence offered a macroscopic apprehension of the said compliances.

In the case of policies, SBERT was used to obtain embeddings of complete policies. This is the ability to process larger blocks of text to obtain an overall ‘gist’ or semantic content of the document; this is necessary for assessing the general compliance of each policy.

The same custom dataset, [SBERTPolicyEmbeddingDataset](#), was used for the analysis consisting of the GDPR compliance labels of each policy and binary vectors.

Training of the model was performed using the classifier with a combination of LSTM and attention layers that enhances the compatibility with the output of SBERT. This setup helps to further enrich the knowledge about the subtleties of compliance at the policy level based on context and subject matters.

To assess the performance of trained model on real-world it was tested on split test set. The main performance indicators that were applied in the assessment of the model include accuracy, precision, recall, and F1-score.

After validation, the model was saved for reusability or for transporting it to other environments most especially the operational environment.

This implementation also proved that SBERT is not only designed for a single sentence but also for the document and this becomes useful when looking for compliance with GDPR in the entire content of the privacy policies. It is particularly important in compliance assessing because it entails a comprehensive approach when addressing the goals of an organisation to have its policies meet all the necessary legal requirements.

### 4.2.2 BERT

#### Sentence Level

As a result of the GDPR compliance complexity at the sentence level, the BERT (Bidirectional Encoder Representations from Transformers) model with the [bert-base-uncased](#) arrangement was employed to resolve this challenge. Due to the powerful language comprehension abilities of BERT, it is very useful in analyzing complex legal text material.

To enhance the training, the initial data collected was prepared by tokenization using BERT's tokenizer. All of the OPP-115 dataset sentences were pre-processed for input to BERT and that required the following: the text was split into tokens up to a maximum length of 512 and where strings was shorter, padding was added and where strings was longer, they were truncated.

Data partitioning was done in the ratio 80 : 20 with the first 80% for training while the last 20% was used for evaluation.

Specifically, custom [WeightedTrainer](#) for the training process was used, which allows for the possibility to tune the loss function based on the class weights, which was aimed at giving balanced treatment to all classes, including the less frequent ones.

In the evaluation of the model the different metrics such as accuracy, precision, recall, F1-score were used. Every one of these metrics offers a measure of the model in classification of GDPR compliance level.

The model and tokenizer after training were saved for replicability as well as for future use of the model for testing or deployment. The training and evaluation outcomes were also backed up for record keeping and showing the efficiency of the model.

Such a wide-range setup showcases the careful usage of BERT for the textual analysis on the level of sentences, meaning that every aspect of the GDPR compliance is examined thoroughly, and the results obtained by the model are recorded in detail to be used for further reference or actual operations.

**Listing 4.2:** Code snippet of the privacy policy dataset for sentence level.

```
1 class PrivacyPolicySentenceDataset(Dataset):
2     def __init__(self, texts, labels, tokenizer, max_len):
3         self.tokenizer = tokenizer
4         self.texts = texts
5         self.labels = labels
6         self.max_len = max_len
7
8     def __len__(self):
9         return len(self.texts)
10
11    def __getitem__(self, item):
12        text = str(self.texts[item])
13        label = self.labels[item]
14        encoding = self.tokenizer.encode_plus(
15            text,
16            add_special_tokens=True,
17            max_length=self.max_len,
18            return_token_type_ids=False,
19            padding='max_length',
20            return_attention_mask=True,
21            return_tensors='pt',
22            truncation=True
23        )
24
25        return {
26            'input_ids': encoding['input_ids'].flatten(),
27            'attention_mask': encoding['attention_mask'].flatten(),
28            'labels': torch.tensor(label, dtype=torch.float)
```

29 }

---

### Entire Policy Level

As for the entire policy level, the BERT model was employed with the version of [bert-base-uncased](#), in which entire privacy policies were fed into it. This level of analysis is very useful to evaluate the overall compliance of a policy with GDPR because it is possible to identify aspects that do not conform to GDPR if they look at only one or a few sentences.

The privacy policies were processed using BERT's tokenizer with a maximum length constraint to manage the size of input texts, ensuring that they fit within the model's capacity while preserving essential information: The privacy policies were processed using BERT's tokenizer with a maximum length constraint to manage the size of input texts, ensuring that they fit within the model's capacity while preserving essential information:

The developed dataset was again partitioned into a training set and an evaluation set, with the percentage ratio of 80 : 20. This division was made in order to pass the tests for unseen data, as this could eliminate certain discrepancies in the evaluation.

The used BERT model was pre-trained to predict each policy into precise of seven compliance categories, which is linked to the principles of GDPR. Training entailed tasking model parameters for accuracy, precision, recall, and F1-score, which are peculiar and vital when it comes to high-stake compliance assessment.

The efficiency of the suggested model has been analyzed on the predetermined benchmarks, and the function to calculate the metrics with the focus on the model's performance in terms of policy classification accuracy in relation to the

level of compliance.

Once this was done, the model along with its tokenizer was stored for use more operational environments or for other assessments. This step ensures the scalability of the model and its readiness for real-time GDPR compliance measure checks.

This section serves to demonstrate the full spectrum to which BERT is applied for complete policy analysis, which characterizes the extent of the processing and analysis that goes into ensuring conformity to GDPR standards of practice.

---

**Listing 4.3:** Code snippet of the privacy policy dataset.

---

```
1 class PrivacyPolicyDataset(Dataset):
2     def __init__(self, policies_text, labels, tokenizer, max_len):
3         """
4         Args:
5             policies_text (list of str): Texts of multiple policies.
6             labels (list of list of int): Labels for each policy.
7             tokenizer: Tokenizer to be used for encoding the text.
8             max_len (int): Maximum length of the tokens.
9         """
10        self.tokenizer = tokenizer
11        self.policies_text = policies_text
12        self.labels = labels
13        self.max_len = max_len
14
15    def __len__(self):
16        return len(self.policies_text)
17
18    def __getitem__(self, idx):
19        encoding = self.tokenizer.encode_plus(
20            self.policies_text[idx],
21            add_special_tokens=True,
22            max_length=self.max_len,
```

```
23         return_token_type_ids=False ,
24         padding='max_length' ,
25         truncation=True ,
26         return_attention_mask=True ,
27         return_tensors='pt'
28     )
29
30     return {
31         'input_ids': encoding[ 'input_ids' ].flatten() ,
32         'attention_mask': encoding[ 'attention_mask' ].flatten() ,
33         'labels': torch.tensor( self.labels[ idx ] , dtype=torch.float )
34     }
```

---

### 4.2.3 GPT2

#### Sentence Level

At the sentence level, GPT2, which is a generative pre-trained transformer model, was used because of its good performance in generating the detailed context-aware embeddings. Such analysis is relevant when one has to dissect privacy policies and check each sentence for compliance with GDPR principles.

The first operation included pre processing the data in the same manner as the GPT2 tokenizer. This tokenizer prepares text for input to GPT2 by breaking it into sentences and then tokenizing each, to include padding to the desired size.

Thus, the GPT2 model was used to classify the given sentences to the GDPR-compliant or non-compliant, with the corresponding specific training scheme. This involved on changing the loss function for handling class imbalance issues by mean of weighted losses.

Hugging Face's Trainer class was used to establish a training regime for the

NER models customizing the Hugging Face Trainer to achieve a weighted class loss means that the loss function was adapted to try to better evaluate the model based on all classes irrespective of the number of data points for each class.

Accuracy, precision, recall, and F1 score were calculated after the training of the model for the complete assessment. These metrics give the number of true positives, true negatives, false positives, and false negatives to decide how accurately the given model is distinguishing between the GDPR compliant and non-complaint websites.

After validation, the trained model and tokenizer are stored so that they can be utilized or implemented for other actual GDPR compliance checks. Additionally, the training results and logs were archived for documentation and analysis purposes: Additionally, the training results and logs were archived for documentation and analysis purposes:

In this way, this complex approach to the realization of GPT2 at the level of a sentence speaks about the possibility of the given model's deep analysis of the text for adherence to the GDPR applying the advanced NLP methods to fit the privacy policies to the norms.

### **Entire Policy Level**

More specifically, GPT2, which demonstrated strong capabilities in reasoning and producing text resembling human-generated content, was used at the policies' overall level to analyze the GDPR compliance of privacy policies. This approach uses GPT2 specifically for its context and evaluation of large text passages which is suitable to this type of analysis.

The dataset formed by entire privacy policies was obtained using GPT2



tokenizer which guarantees that every policy is tokenized and encoded accurately. As for the tokenization, particular attention was paid to its efficient implementation since the option was limited due to word count.

The training was preformed on a version of GPT2 model adjusted for sequence classification for multiple labels. The specified model was designed to make the exact predictions concerning the seven GDPR principles, using the sigmoid function to probability-independent logits of the model.

To mitigate class imbalance, a weighted training was done during the training process, In the evaluation of the models' effectiveness, different measures like accuracy, precision, recall, F1-score, was used.

The efficiency of the proposed model was also verified with the help of a set of measurable criteria that would allow for the quantization of the model's performance in terms of its accuracy in GDPR compliance classification and its improvement.

After validation, the trained model and the tokenizer were saved to allow for direct use or for the actual operational GDPR compliance testing.

The elaborate setup guarantees that GPT2 is well-utilized to apply on full scraps of text documents' privacy policies, due to its efficient contextual comprehension trait to provide a holistic and accurate compliance audit of the GDPR policies.

## **4.3 Testing on Unlabeled Data**

The saved models for both the granularity levels helps apply the solutions to the ACL Coling dataset that contains unlabeled policies. This step was important in order to determine the models' prediction ability and how well they would

perform in actual situations.

### 4.3.1 SBERT

#### Sentence Level

In the last phase of the implementation, the SBERT model is used to evaluate the discrete sentences of privacy policies against the particular GDPR principles. To describe this approach, several operations are invoked, such as the loading of the model's weights, generation of predictions from the sentences, and overall evaluation of the policy.

To load the pre-trained SBERT classifier, there is a custom function to make sure that the model is in evaluation mode which is very important when making prediction without changing the learned parameters of the model.

The function `sentence\_sbert\_classify\_policy` takes a policy, breaks down sentences in the policy, extracts embeddings, and uses the classifier to proceed to GDPR compliance. The function assigns the GDPR labels for each sentence and gathers both the labels and the probabilities of them being correct.

The outcomes received from the classifier are then summed up to offer a more combined assessment of the current policy's compliance. Such measures consist in presenting the classification results referring to each of the sentences with the corresponding GDPR labels assigned and such the specific GDPR labels revealed throughout the document to supply some practical information regarding the policy's compliance with the GDPR.

It not only enables the analysis at the sentence level but also compiles those findings to provide an evaluation of the Policy's general conformity. Thus, with the help of employing the SBERT model in the suggested way, the thesis outlines

the applicability of highly developed methods in the sphere of NLP in the actual environment of GDPR compliance evaluation.

### Entire Policy Level

On the entire policy level, the trained SBERT model is used to determine compliancy of the entire privacy policies under GDPR. This entails evaluating a document with multiple sentences for banking on SBERT's capacity to analyze large texts.

A few preprocessing techniques are applied to the images the most important initial step, though, is to load the pre-trained SBERT classifier to be ready for the inference. The model is prepared to match the demands that come with evaluating entire policies, while keeping an eye on the integer's accuracy and speed.

The function for policy classification is `policy\_sbert\_classify\_policy` that performs the analysis of each policy to produce the GDPR compliance prognosis. This particular function is the one that computes the embeddings, the one that involves the classifier in order to predict compliance, and the one that applies a certain threshold and identifies the GDPR labels that are relevant.

Because the classification is done in real-time, the final output is aggregated and presented with an overview of the policy's conformity to GDPR. This entails general and specific predictions for all the policies, where a diagram with the compliance level is given as well as any particular GDPR labels that are determined.

This section shows how the use of SBERT to apply it to entire privacy policies also leads to a much more detailed and profound analysis of the GDPR compliance. Applying the new NLP methods, the thesis demonstrates a way of

using of the machine learning techniques in improving the legal compliance check system.

### 4.3.2 BERT

#### Sentence Level

For the sentence level analysis, BERT (Bidirectional Encoder Representations from Transformers) model is used to classify the GDPR compliance of all the sentences in the given privacy policies. This detailed approach enables one to distinguish between a particular policy compliance concern and a general concern, which is a very important aspect when doing policy analysis.

To harness BERT's capabilities, the pre-trained BERT was loaded along with its tokenizer for the pipeline construction. The model is trained with an objective of predicting the hidden states which can be used further for analysis or for high level classification.

A pipeline is trained with the help of Hugging Face's 'pipeline' function, with only the classification pipeline being used. This pipeline makes the process of sentence classification easier due to the fact that both the model and tokenizer can be incorporated into the pipeline as a callable.

The function of classification is called like `sentence\_bert\_classify` which runs the pipeline for each sentence. It sets a filter to know whether a particular sentence belongs to a certain GDPR tag according to the model's probability score.

Proceeding from the given presentation, having received the data, the results are summarized and presented so that it will be clear whether each proposed sentence meets or violates the requirements of standard English. This entails

predictions for all the sentences as well as a summary of the novel GDPR labels encountered within the policy.

In this part, BERT is implemented and explained at the sentence level for GDPR compliance evaluation of policies, which again proves the effectiveness of integrating policies at the corpus level and interpreting the outcomes based on the classification of individual sentences.

### Entire Policy Level

The BERT model, which is reported to have a good understanding of the language context, is employed at the policy level to give an overall score on the GDPR compliance of overall policies from top to bottom and in complete privacy policies as a whole. This approach gives a broad view of compliance that other single-sentence analysis fails to show, due to some intricacies.

The first step entails preparing the model derived from pre-trained BERT model and the tokenizer. The model is set up to return the hidden states: These could be used for further analysis or to improve on the classification bit.

The classification pipeline is defined with Hugging Face's [pipeline](#) function. This pipeline helps in the efficient and effective classification of text data through the application of the model.

The function called [policy\\\_bert\\\_classify\\\_policy](#) is specifically intended to run the classification pipe through entire policies. It rates each policy, imposes a filter to the model to decide pertinence, and indicates GDPR compliance labels.

Classification results are then summed up and reported to offer extended information to the policy's conformity. This entails putting down the list of identified GDPR labels and the corresponding confidence levels. This process not

only provides an all-inclusive assessment at the policy level but also consolidates these outcomes to provide the big picture on the policy's compliance. Thus, using BERT in this manner, the thesis illustrates how state-of-the-art NLP tools can be effectively applied to solve practical problems in the real-life framework of GDPR compliance analysis.

### **4.3.3 GPT2**

#### **Sentence Level**

At the sentence level, the GPT2 model, is used to predict the GDPR compliance of each and every sentence in the of privacy policies. That is why this approach that we focus on the analysis of each sentence separately is suitable best in combination with GPT2, as this model is best used in understanding context and generating more subtle interpretations of texts.

The first process consists of loading the GPT2 model and the tokenizer that comes with the model. The model has been made specifically to generate the hidden states which help interpret the types of produced structures and improve overall performance of the classifier.

A classification pipeline is created with a function 'pipeline' from Hugging Face. This setup will also ease the applying of the model to textual data so as to enhance the efficiency of the classification.

The function `sentence\_gpt\_classify\_policy` takes each of the sentences and feeds it through the classification pipelines. It scores each sentence for relevance using set thresholds to the model's confidence, and the final records will have GDPR compliance labels.

After classification, the outcomes of the same are summed up and displayed

so as to enable the production of a preview of each of the sentences' statutory compliance. This includes, exhibiting statics of the classification outcomes of each sentence and any fresh GDPR tags that were found in the document.

This section also exemplifies how GPT2 can be used at the lowest level in order to reflect on specific aspects of GDPR compliance in the privacy policies and show the ability to generate valuable and specific recommendations from the analysis of the classification of the sentences.

### Entire Policy Level

Using GPT2 at the policy level means evaluating the GDPR compliance of whole sections of texts, entire privacy policies, in fact. It makes use of GPT2's deep learning capacity to analyze text structures, which makes it possible for Comply to identify compliance problems in whole documents.

The first and foremost of all the operations that are required are to load the pre-trained GPT2 model and tokenizer for text preprocessing. Parameters of the model are changed to allow getting desired hidden states which would allow for better text analysis.

An end-to-end classification pipeline is set up with the help of the Hugging Face 'pipeline' widget as it wraps up an easy model application process. This setup enhances the processing of text data since small but considerable text data is processed by a small network, and large but negligible text data is processed by a large network.

The call function in this pipeline is the classification function `policy\←  
_gpt\_classify\_policy` which is expected to process each policy. This function reviews the full text of each policy and filters the predictions that cross a predefined

threshold and the corresponding GDPR compliance labels are returned.

After this, the summary of the results is produced and presented in order to give an overview of the compliance status of the policy. This ranges from listing all GDPR labels perceived with their respective confidence level thus ending up with an in-depth understanding of the policy under the GDPR regulations.

This section demonstrates how at the policy level GPT2 is employed for the purpose of comparing and evaluating GDPR compliance of whole privacy policies, thus this example perfectly proves the ability of GPT2 to provide specific and valuable recommendations founded on thorough text analysis.

## 4.4 Single Principle Checker

The Single Principle Checker was designed to address the problem of detecting cases of GDPR non-compliance at the sentence level concerning individual principles only. This implementation was specific with the “Storage Limitation” principle under the GDPR law which is very crucial.

### Model Overview

The model employed for this concerned is a variant of SBERT referred to as [all-MiniLM-L6-v2](#) which establishes proficiency in the production of customized numerical representations of the sentences hence the embeddings. The latter captures more refined semantics of the text needed for the compliance assessment.



### Data Preparation

The embeddings and labels were managed by a custom dataset class called [SBERTEmbeddingDataset](#). The dataset made sure that every sentence was morphed into a high-dimensional vector that would reflect all of the semantic information that was in the sentence, as well as a label that was either 1, if the Storage Limitation principle was followed or 0 if it was not.

### Training Process

The training process used a binary classification approach where a neural network architecture that has LSTM layers and attention was used in estimating the compliance. This setup of embeddings ensured that the model could home in on some of the distinctive characteristics in the mappings of figures, thus improving its capacity to detect hints of compliance or lack thereof.

### Optimization and Evaluation

These hyperparameters include learning rate and epochs, and in order to obtain the best configuration that gave the best performance in terms of F1 score which is perfect for binary classification problems the training was experimented under different configurations. Hence, measures such as accuracy, precision, recall, and F1-score were obtained to give a general performance of the model. The best number of neighbours was determined and the model was then trained and the saved for the model configuration.

### **Experimentation Details**

Further, more trials were made in practicing the model in the apt style, fit for rigorous testing. Various learning rates including  $5*10^{(-5)}$ ,  $2*10^{(-5)}$ ,  $1*10^{(-5)}$ ,  $1*10^{(-4)}$ , and numbers of epochs including 5, 10, 15, 20 were also experimented, and decision thresholds of 0.3, 0.5, 0.7, 0.8 were used also. The outcome showed that the function of the model differed depending on their choice – this proved the fact regarding model customization for compliance principles.

## **4.5 OpenAI Embedding Experiments**

Furthermore, GPT-3 API embeddings were also captured, and they were processed to be input to classification. This experiment again considered the fundamental generative capability of GPT-3 in order further to improve the classification accuracy.

This section outlines how OpenAI embeddings, and particularly the Ada model, can be used to classify the specified privacy policy-based sentences according to their GDPR conformity. The experiments include: loading a checkpoint file, configuring a classification pipeline, training a neural network and performance assessment.

The import of the Ada model and the processes of its tokenizer from OpenAI's API is necessary as they are used to obtain suitable embeddings that express the meanings of the privacy policy texts.

A function is made to pass through texts and obtain the corresponding embeddings from the Ada model. These embeddings serves as primary inputs to a neural network classifier.

This is to ensure that the different policies have their own embeddings and labels as they go through the training process. This makes it easier to sort data while during the training or validation of the model in the subsequent phases.

A multilayer perceptron is proposed to determine whether a given item conforms to the GDPR or not based on the embeddings created by the pre-trained language model. It includes linear layers, which are commonly denoted as weights in the neural networks, as well as activation functions.

It outlines the sequence of setting up the training and validation as well as defining the necessary components such as the computation of loss, and performance metrics like accuracy, precision, recall, and F1 score.

## 4.6 Model Evaluation

Last but not least, each model's efficiency was assessed based on the aforementioned criteria: precision, recall, F1-score, and accuracy. The results which we will review in the next chapter 5. These were compared to establish whether the ascuslevaluate training at the sentence-level as compared to the policy-level training were effective. Also to compare the performance of the developed model with and without the GPT-3 embeddings to determine the contribution of GPT-3 embeddings in the classification task.

In this chapter, the reader is given a step by step description of the implementation strategy for the NLP models that is used for the experiment that is executed in this research work. The use of the methodology ensures that the findings arrived at can be easily repeated and therefore ascertaining the reliability of applying the findings to different assessments of GDPR compliance.

## **Chapter 5**

# **Evaluation and Discussion**

This chapter delves into the performance evaluation of the machine learning models used in this thesis, namely SBERT, BERT, and GPT2, alongside an exploration of GPT-3 embeddings. The focus is on finding if they are efficient in GDPR compliance classification at both the sentence and policy levels.

### **5.1 Overview of Model Performance**

#### **Performance Metrics**

The models were evaluated based on accuracy, precision, recall, and F1-score.

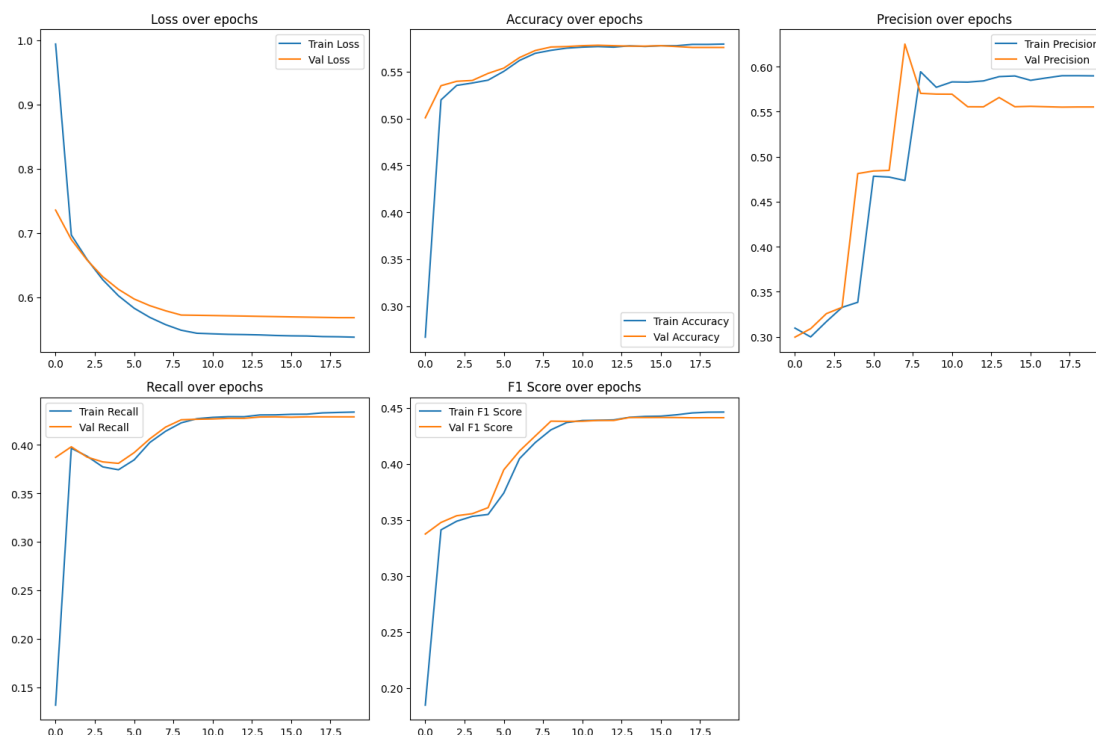


Fig. 5.1. Metrics for SBERT performance on sentence level

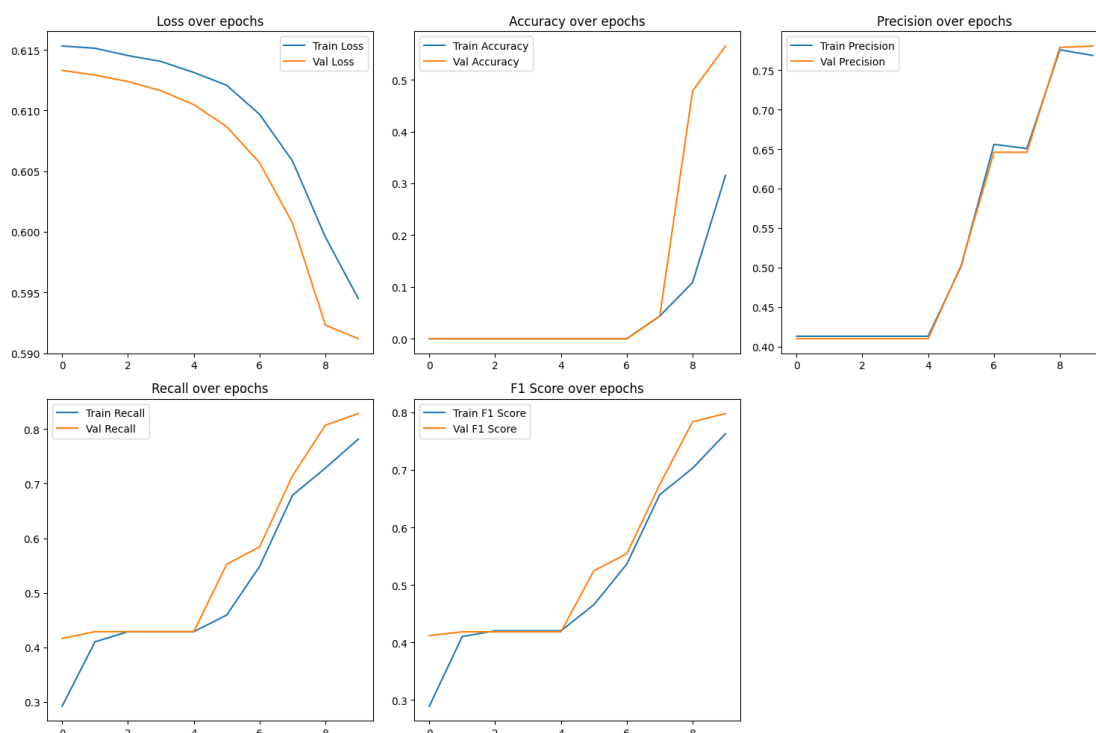


Fig. 5.2. Metrics for SBERT performance on policy level

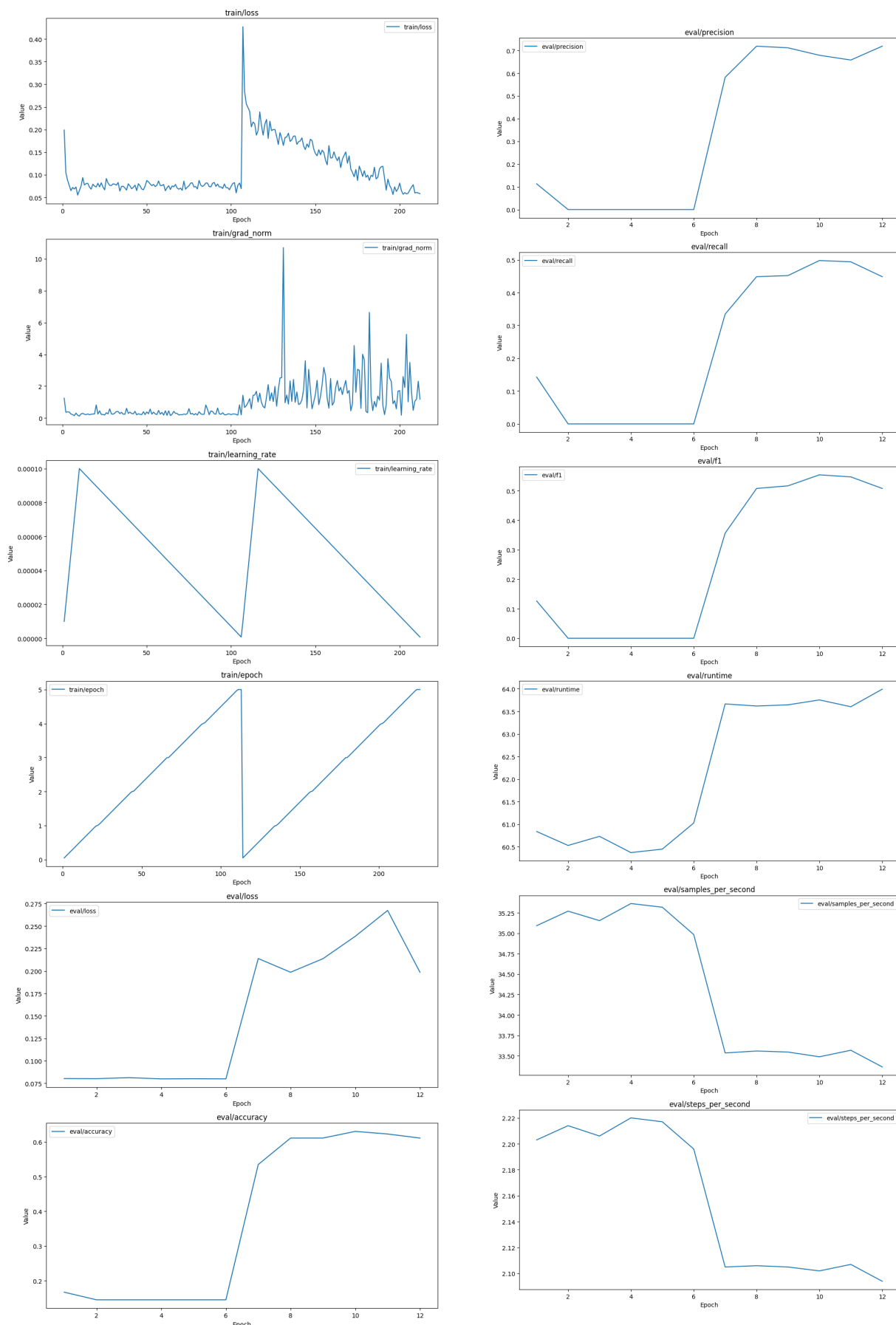


Fig. 5.3. Metrics for BERT performance on sentence level

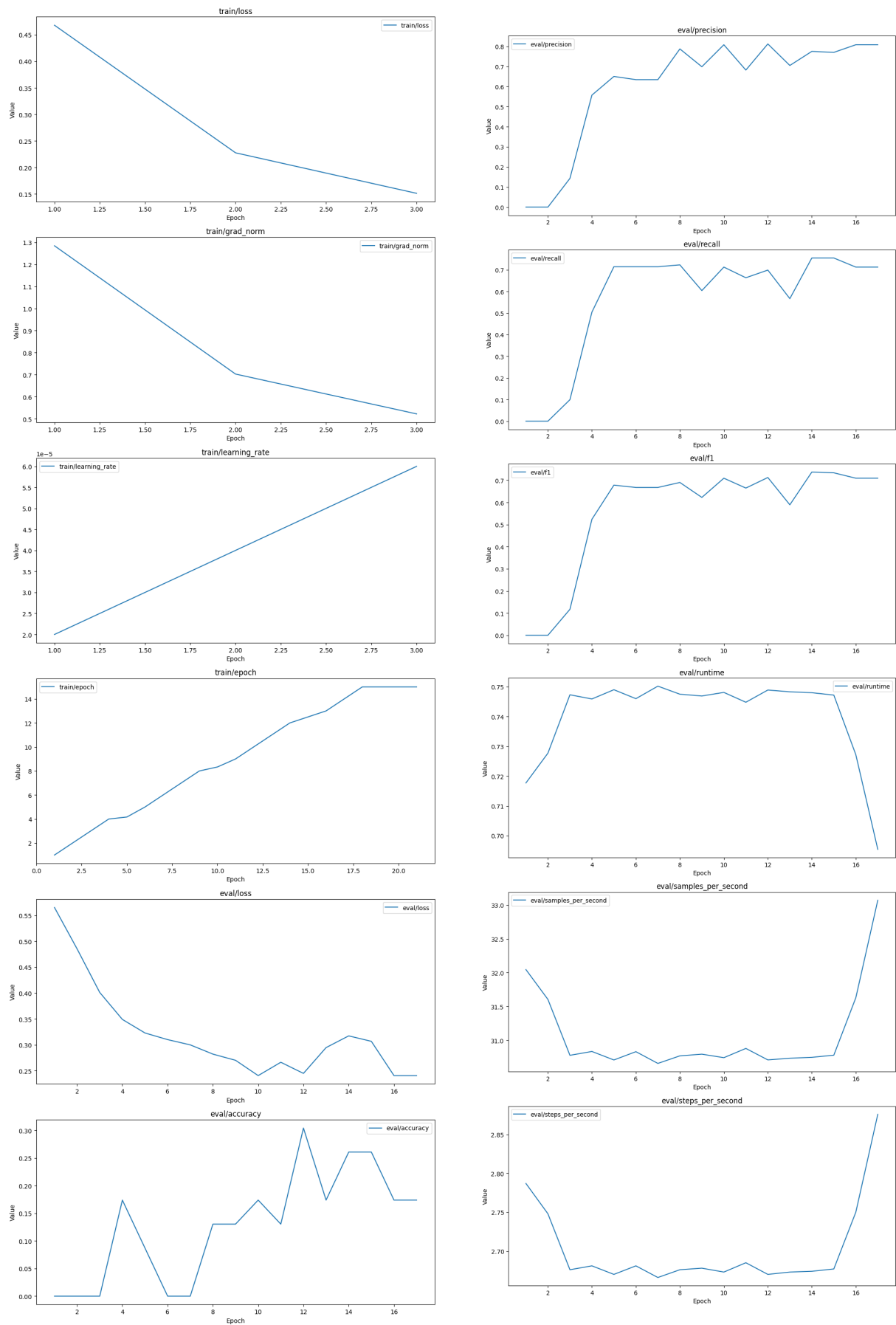


Fig. 5.4. Metrics for BERT performance on policy level

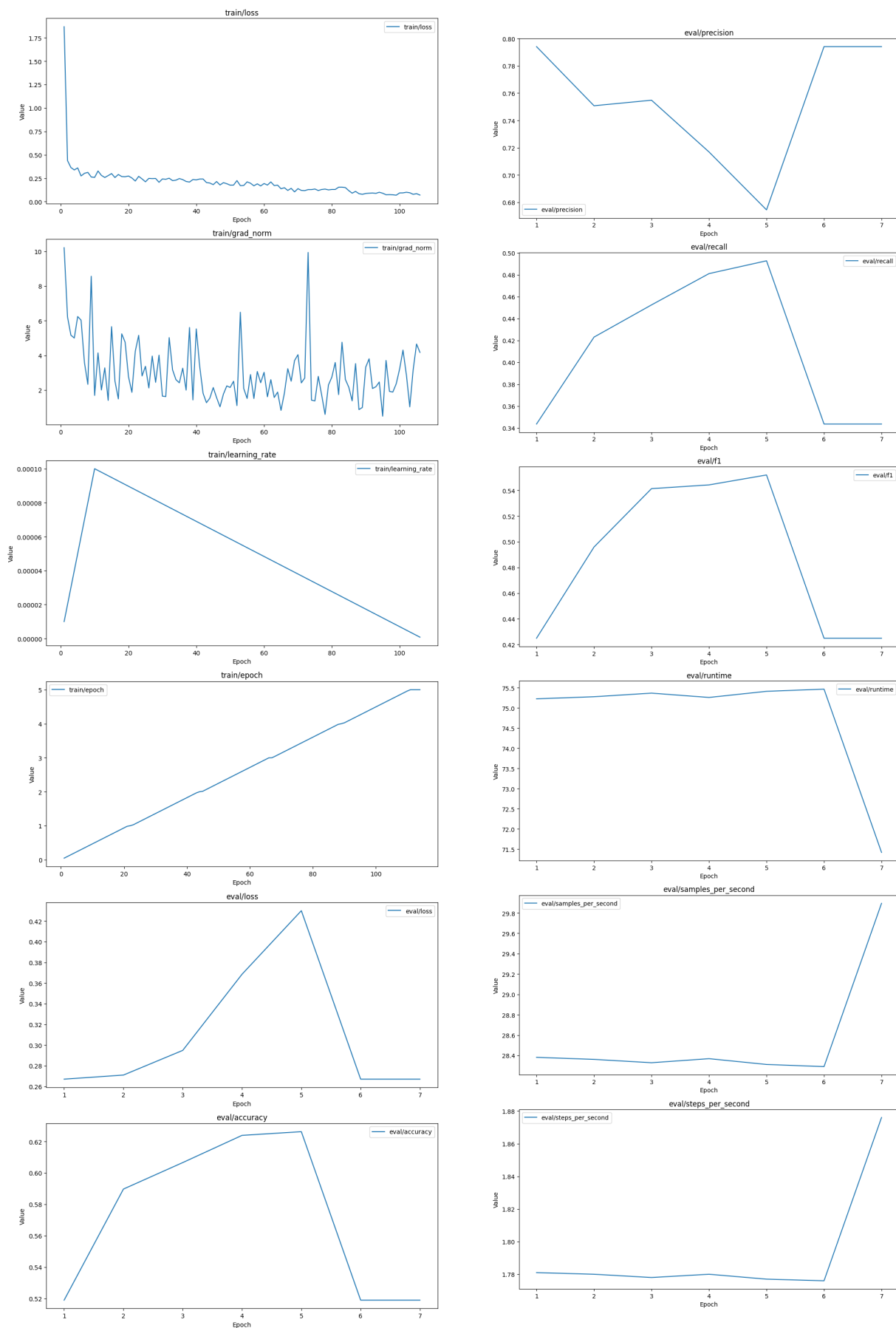


Fig. 5.5. Metrics for GPT2 performance on sentence level



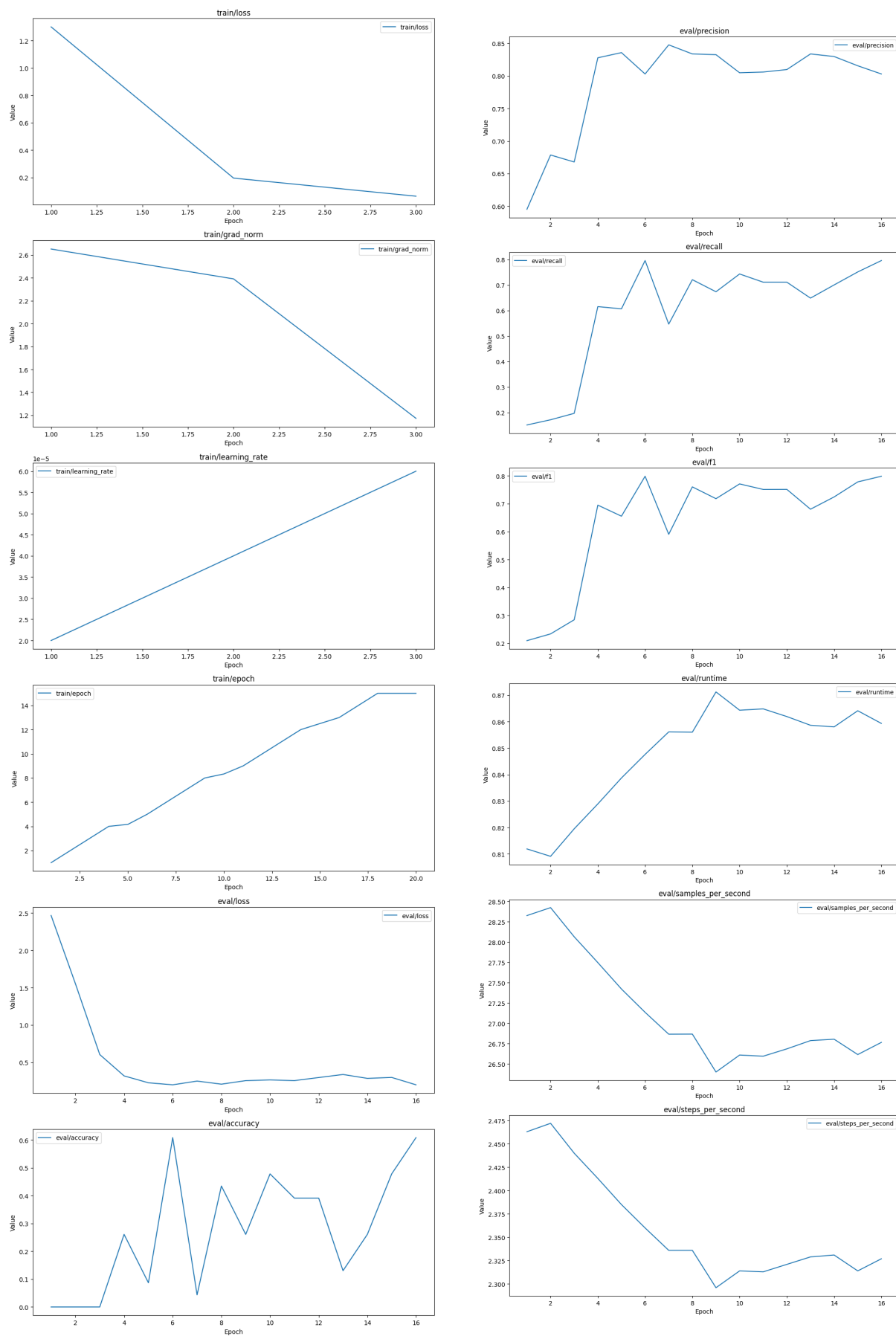


Fig. 5.6. Metrics for GPT2 performance on policy level

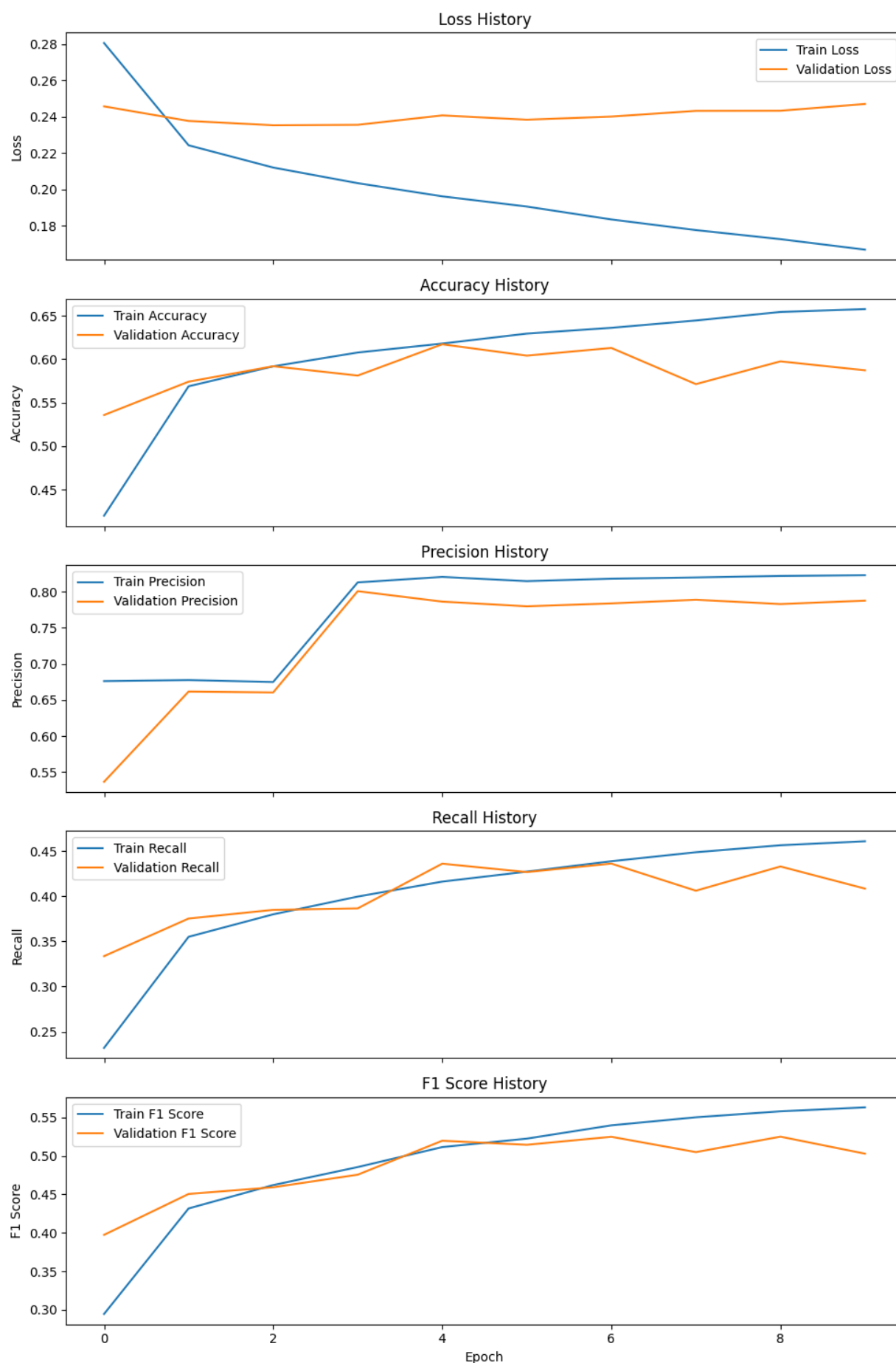


Fig. 5.7. Metrics for OpenAI embeddings performance

Comparison Across Models

Comparative analysis revealed that while GPT2 performed well in sentence-level classifications, achieving the highest recall rate, it struggled with precision, often misclassifying non-compliant instances. On the other hand, the introduction of GPT-3 embeddings significantly enhanced the accuracy and F1-scores across all models, demonstrating the power of advanced embeddings in improving classification tasks.

TABLE 5.1  
Model Evaluation Metrics

Metric	Set	Sentence Level			Policy Level			API Model
		SBERT	BERT	GPT2	SBERT	BERT	GPT2	
Accuracy	Train	0.578	-	-	0.576	-	-	0.628
	Val	0.571	0.609	0.521	0.435	0.622	0.435	0.599
Precision	Train	0.575	-	-	0.761	-	-	0.823
	Val	0.547	0.764	0.729	0.809	0.743	0.796	0.770
Recall	Train	0.430	-	-	0.857	-	-	0.434
	Val	0.428	0.857	0.429	0.757	0.464	0.760	0.431
F1-Score	Train	0.441	-	-	0.802	-	-	0.543
	Val	0.43	5 0.803	0.521	0.778	0.550	0.775	0.518

5.2 Analysis of Results

For the results, after looking at the explanation above we can see that SBERT performed the best of across all the models practically even though in terms

of metrics it was far from the best. Let's take a dive into how it works and predicts compliance for different sentences and privacy policies. We will start with checking compliance of sentences against a single principle from the GDPR article five: "Storage Limitation"

**Sample 5.2.1: Non-compliant sentence example to Storage Limitation**

we collect two kinds of information related to you:

- (a) personally identifiable information ("pii"); and
- (b) non-personally identifiable information ("non-pii").

**Sample 5.2.2: Compliant sentence example to Storage Limitation**

furthermore, some information is never completely removed from our databases due to technical (e.g., information that is stored in our "back-up" systems) and certain legal constraints.

Four sample sentences, two of which adhere to Storage Limitation and two of which don't were evaluated on the three models and two of these two sample sentences can be seen above. The models performed as seen below:

TABLE 5.2  
Comparison of metrics for single principle before class resampling

Metric	SBERT	BERT	GPT2
Accuracy	0.97	0.97	0.99
Precision	0.50	0.45	0.98
Recall	0.39	0.42	1.0
F1 Score	0.43	0.43	0.99

At first BERT and GPT2 performed even worse than the SBERT but after resampling the dataset, their performance skyrocketed.

TABLE 5.3  
Comparison of metrics for single principle after class resampling

Metric	SBERT	BERT	GPT2
Accuracy	0.88	0.99	0.99
Precision	0.83	0.99	0.98
Recall	0.96	1.0	1.0
F1 Score	0.89	0.99	0.99

Then moving on to test these model in a more practically evaluation based way, they were fed the sentences shown below. It was interesting to note that for single principle training, they all produced equivalent and correct results.

TABLE 5.4  
Compliance analysis for privacy principles using NLP models

Principle	Sentence	ChatGPT API	SBERT	BERT	GPT2
<b>Storage</b>	<u>We collect two kinds of information related</u>	Purpose	False	False	False
<b>Limitation</b>	<u>to you: (a) personally identifiable information (“PII”); and (b) non-personally identifiable information (“non-PII”).</u>	Limitation			
	<u>Furthermore, some information is never completely removed from our databases due to technical (e.g., information that is stored in our “back-up” systems) and certain legal constraints.</u>	Storage Limitation	True	True	True

Continued on next page

Table 5.4 continued from previous page

Principle	Sentence	ChatGPT API	SBERT	BERT	GPT2
	<u>This policy states that user data will be stored for a maximum of two years.</u>	Storage Limitation	True	True	True
	<u>We reveal only the last four digits of your credit card numbers when confirming an order.</u>	Data Mini- mization	False	False	False

After seeing the result, I went ahead to confirm using [ChatGPT API](#) as follows:

Listing 5.1: Code snippet of the ChatGPT request.

```

1 client = OpenAI()
2
3 completion = client.chat.completions.create(
4     model="gpt-4o",
5     messages=[
6         {"role": "user", "content": "\"\"From the sentence: *sentence placeholder↵
          * Does it comply with any of the principles from the GDPR article 5?↵
          If yes, then which is the main one?\"\""},
7         {"role": "user", "content": "\"\"Hello!\"\""}
8     ]
9 )
10
11 display_open_ai_response(completion.choices[0].message.content)

```

From the above, for the first non compliant sentence, the response was:

Hello!

From the sentence: "we collect two kinds of information related to you: (a) personally identifiable information ("pii"); and (b) non-personally identifiable information ("non-pii")." Does it comply with any of the principles from the GDPR article 5? If yes, then which is the main one??

Hello! Yes, the sentence you provided does comply with principles from Article 5 of the General Data Protection Regulation (GDPR). Specifically, it aligns most closely with the principle of "**Purpose Limitation**" and "**Data Minimization**."...

Meanwhile for the second sentence that was compliant, the response from ChatGPT was:

Hello!

From the sentence: "furthermore, some information is never completely removed from our databases due to technical (e.g., information that is stored in our "back-up" systems) and certain legal constraints." Does it comply with any of the principles from the GDPR article 5? If yes, then which is the main one?

...The sentence appears to relate most directly to the **Storage Limitation** principle. This principle requires that personal data be kept no longer than necessary for the purposes for which it is processed. However, the sentence explains why some data cannot be completely removed due to technical and legal constraints, which acknowledges that there are exceptions to the complete erasure of personal data...

After this little success, it was relevant to move on and try it out in a more practical setting. There was a privacy policy that was fed to the model and it had to extract the sentences that made it classify this policy as adhering to **Storage Limitation**.

### **Sample 5.2.3: Sample Privacy Policy**

**Last Updated: June 10, 2024**

#### **Introduction**

Welcome to [Company Name] (“we”, “us”, “our”). We are committed to protecting and respecting your privacy. This policy explains how we collect, use, and disclose your personal data, as well as your rights in relation to that data.

#### **Information We Collect**

We collect various types of information in connection with the services we provide, including:

- Personal Identification Information: such as name, email address, phone number, etc.

#### **Use of Information**

We use the collected information for various purposes, including:

- Providing and maintaining our services.



**Storing and Retention***a. Retention Period*

We will retain your personal data only for as long as is necessary for the purposes set out in this privacy policy. The criteria we use to determine retention periods include:

- The length of time we have an ongoing relationship with you and provide services to you.
- Whether there is a legal obligation to which we are subject (e.g., certain laws require us to keep records of your transactions for a certain period before we can delete them).
- Whether retention is advisable considering our legal position (such as in regard to applicable statutes of limitations, litigation, or regulatory investigations).

*b. Secure Deletion*

Once the retention period expires, we will securely delete your personal data. Methods of deletion may include:

- Permanently deleting electronic records.
- Securely shredding physical documents.
- The right to erasure (the “right to be forgotten”).

**4. Data Security**

We implement appropriate technical and organizational measures to protect your personal data against unauthorized access, alteration, disclosure, or destruction. These measures include encryption, access controls, and secure data storage facilities.

**5. Your Data Protection Rights**

Depending on your location and applicable data protection laws, you may have the following rights:

- The right to access your personal data.

- The right to rectify incorrect or incomplete data.
- The right to erasure (the “right to be forgotten”).
- The right to restrict the processing of your data.
- The right to data portability.
- The right to object to the processing of your data.

#### **6. Changes to This Privacy Policy**

We may update our privacy policy from time to time. We will notify you of any changes by posting the new privacy policy on this page. You are advised to review this policy periodically for any changes.

#### **Contact Us**

If you have any questions or concerns about this privacy policy, please contact us at:

- Company Name
- Address
- Email
- Phone Number

Specifically the highlighted sections show the adhering sentences as printed by the model and they comply with the definition given by the GDPR of this principle.

So far this was just training and testing on one principle. Following this, let's look at the results of the multilabel classification for all 7 principles on the sentence level. In this case all three of the models used a relatively high threshold of 0.8 and had the metrics as seen below:

TABLE 5.5  
Comparison of metrics for multi principle classification in sentence level

Metric	SBERT	BERT	GPT2
Accuracy	0.57	0.62	0.62
Precision	0.56	0.72	0.74
Recall	0.43	0.45	0.46
F1 Score	0.44	0.52	0.55

Subsequently, these models were saved and tested on entire policies. However, because they were trained on a sentence level granularity, the privacy policies were split into sentences before being analyzed by the models and finally a compliance report was generated highlighting the principles covered with sample sentences and which principles need to be worked on.

### GDPR Compliance Report for yola

**Summary:** Non-compliant with 2 out of 7 principles evaluated.

#### Detailed Findings:

**Lawfulness, Fairness and Transparency:** Compliant

*Example: if in addition, from time to time we may collect demographic, contact or other personal information you voluntarily provide to us, such as in connection with your participation in surveys, sweepstakes, contests, games, promotional offers, and other activities on the site.*

**Purpose Limitation:** Compliant

*Example: in general, we use personal information we collect to process your requests or transactions, to provide you with information or services you request, to inform you about other information, events, promotions, products or services we think will be of interest to you, to facilitate your use of, and our administration and operation of, the site, newsletters and for the purpose for which the information was provided.*

**Data Minimization:** Compliant

*Example: in general, we use personal information we collect to process your requests or transactions, to provide you with information or services you request, to inform you about other information, events, promotions, products or services we think will be of interest to you, to facilitate your use of, and our administration and operation of, the site, newsletters and for the purpose for which the information was provided.*

**Accuracy:** Compliant

*Example: if your personally identifiable information changes, or if you no longer desire our service, you may update your profile or delete it by clicking on the profile link after you log in and then clicking the delete account link at the bottom of the page.*

**Storage Limitation:** Non-compliant

**Integrity and Confidentiality:** Compliant

*Example: when you enter sensitive information (such as credit card number and/or social security number, national id, personal health information) on our registration or order forms, we encrypt that information using secure socket layer technology (ssl).*

**Accountability:** Non-compliant

#### Recommendations:

Review the policy to align with the following principles:

Storage Limitation, Accountability

Fig. 5.8. SBERT Compliance Report for Sentence Level

## GDPR Compliance Report for yola

**Summary:** Non-compliant with 5 out of 7 principles evaluated.

### Detailed Findings:

**Lawfulness, Fairness, and Transparency:** Compliant

*Example: our service is not intended to be used by children under the age of 13, and we do not knowingly collect personal information from children under the age of 13 except in compliance with applicable law.*

**Purpose Limitation:** Non-compliant

**Data Minimization:** Non-compliant

**Accuracy:** Non-compliant

**Storage Limitation:** Non-compliant

**Integrity and Confidentiality (Security):** Compliant

*Example: credit card details are stored in our payment gateway's highly secure payment vault that is fully compliant with the payment card industry's data security standards.*

**Accountability:** Non-compliant

### Recommendations:

Review the policy to align with the following principles:

Purpose Limitation, Data Minimization, Accuracy, Storage Limitation, Accountability

Fig. 5.9. BERT Compliance Report for Sentence Level

## GDPR Compliance Report for yola

**Summary:** Non-compliant with 3 out of 7 principles evaluated.

### Detailed Findings:

**Lawfulness, Fairness, and Transparency:** Compliant

*Example: we may collect certain information you voluntarily provide to us which may contain personal information.*

**Purpose Limitation:** Compliant

*Example: we may also share personal information with vendors, consultants and other service providers ("service providers") who are engaged by or working with us in connection with the operation of the site or the services and who need access to such information to carry out their work for us.*

**Data Minimization:** Non-compliant

**Accuracy:** Compliant

*Example: if your personally identifiable information changes, or if you no longer desire our service, you may update your profile or delete it by clicking on the profile link after you log in and then clicking the delete account link at the bottom of the page.*

**Storage Limitation:** Non-compliant

**Integrity and Confidentiality (Security):** Compliant

*Example: when you enter sensitive information (such as credit card number and/or social security number, national id, personal health information) on our registration or order forms, we encrypt that information using secure socket layer technology (ssl).*

**Accountability:** Non-compliant

### Recommendations:

Review the policy to align with the following principles:

Data Minimization, Storage Limitation, Accountability

Fig. 5.10. GPT2 Compliance Report for Sentence Level

Next for the policy level, **SBERT** only provided results when the threshold was as low as 0.5. The metrics for all models can be seen below:

TABLE 5.6  
Comparison of metrics for multi principle classification in policy level

Metric	SBERT	BERT	GPT2
Accuracy	0.60	0.43	0.43
Precision	0.76	0.80	0.79
Recall	0.85	0.75	0.75
F1 Score	0.80	0.77	0.77

In the policy level, the reports were not created because the models were not able to extract sentences and could only classify the entire policies as trained to do. The results were instead in the following format and just like with the sentence level granularity and the single policy classification, SBERT provides more detailed results.

#### Sample 5.2.4: Result Template for Policy Level

Policy Classification Results: Policy: \*inserts full policy\* Prediction: Purpose Limitation, Score: 0.9508 Prediction: Data Minimization, Score: 0.5218  
Unique GDPR Labels identified in this policy: Purpose Limitation Data Minimization

## 5.3 Challenges Encountered

The following are some of the challenges which the models experienced while in the process of the research along with the kind of approaches that were taken to handle the said challenges.

### Class Imbalance

One major issue that was experienced was the issue of class imbalance especially regarding GDPR Principle 7 (Accountability) where there were no passing cases. This affected the training process because stratification occurred in a way that favored a certain class and so strategies such as class weights had to be applied to address this issue. The weighting of the trainers mitigated the impact of one of the classes during the training of the model.

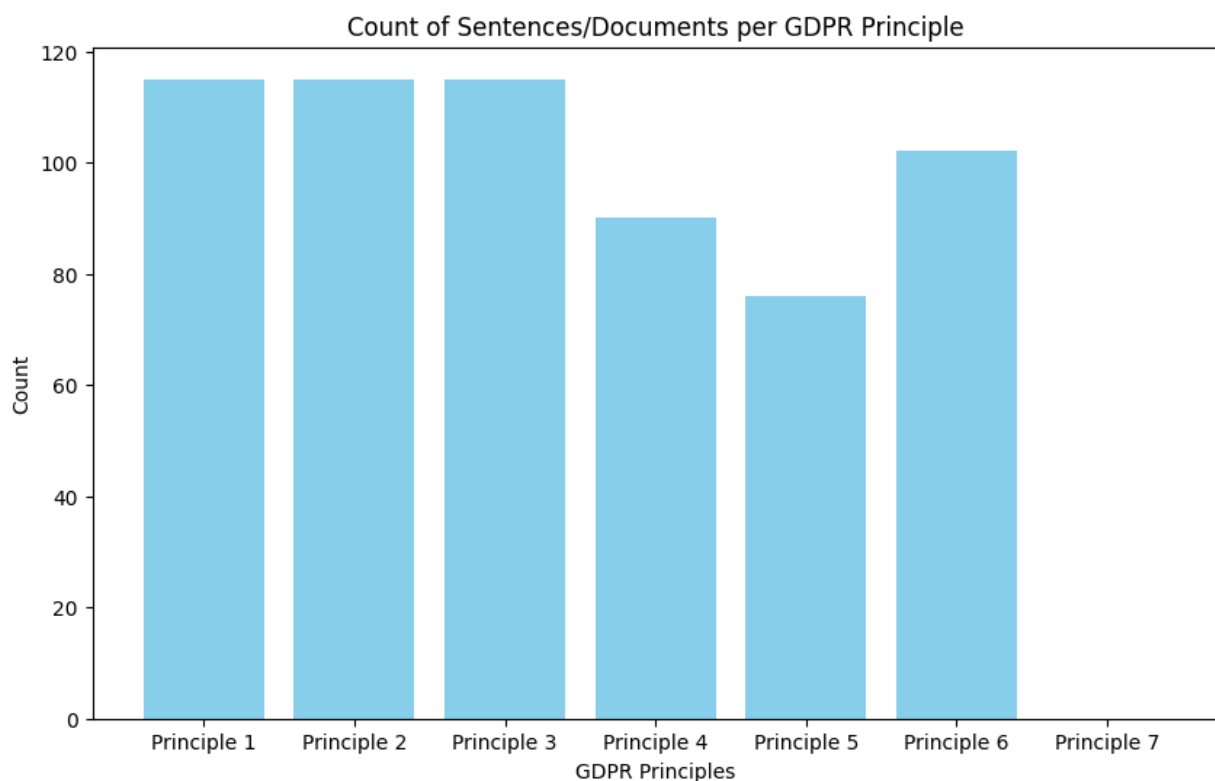


Fig. 5.11. Enter Caption

### Dataset Preparation and Preprocessing

In the aspect of preparing and preprocessing the dataset, there were multiple challenges:

- **Annotations to Labels Conversion:** The pre-existing data set had anno-

tations instead of labels in it. Furthermore, such annotations were aligned with categories of OPP-115 instead of GDPR principles. First of all, an attempt was made to assign these categories to the principles of GDPR by organizing them through clustering. However, strategy mentioned above could not prove to be fruitful.

- **Manual Mapping:** Finally, after searching for quite a while, a paper that contained the mapping between OPP- 115 categories and GDPR principles was used to solve this problem as well.
- **Sentence Formation:** Annotation to construct sentences was somehow difficult especially in counting the number of words that is supposed to constitute a sentence and handling the issue of duplicates.

### Threshold Selection and Hyperparameter Tuning

Selecting the right thresholds was also an issue. Some of the models had the maximum probability of assigning to a prediction as 0.53 does not reflect the criteria of 0.8 was wanted for the final step of prediction. Hyper parameter tuning was used to optimize the model and search out the best parameters of the model such as learning rate, threshold and training epochs.

### Model Selection

Deciding between base models and legal-special models such as `nlpaueb/legal-bert-base-uncased` was another one of the problems. Some of the time, legal versions of models did not perform as well as hypothesized or result in any significant difference, and it was decided to stick with the base models for the tasks. The

following table is to show the performance results of the legal versions of the models against the base models. It helps to compare which one of the model versions is better suited for the given task.

TABLE 5.7  
Comparison of Legal vs. Base Model for BERT

Metric	Base BERT		Legal BERT	
	Sentence Level	Policy Level	Sentence Level	Policy Level
Accuracy	0.62	0.56	0.43	0.13
Precision	0.72	0.72	0.73	0.70
Recall	0.45	0.85	0.34	0.65
F1 Score	0.52	0.77	0.44	0.67

### Max Length Limitation of Models

Another interesting problem faced was the restriction of the input of BERT and GPT2 models up to 512 tokens. This limitation meant that it was not possible to show the models full policies at once and therefore text always had to be split into sub-sections. While this approach proved beneficial in terms of the models' ability to analyze the data, it also had the side effect of limiting context and continuity between segments, which could decrease the validity and reliability of the outcomes.

### Resource and Financial Constraints

Several resource and financial constraints were encountered:

- **OpenAI API Costs:** Since the OpenAI embeddings API is not free, it incurred a lot of expenses in trying to make the deep learning model work.



This financial constraint was one of the reasons why extensive use of this method could not be applied.

- **Compute Resources:** Running models with GPU on Colab required payment for additional compute units. Moreover, due to long running times of the models getting results was a slow and costly process.

## 5.4 Discussion on GPT-3 Embeddings

GPT-3 can be regarded as one of the most efficient language models created so far, which has been designed by OpenAI. It can generate text that resembles the text written by people, provided that it has received the appropriate input. It makes use of transformer technology on 175 billion parameters to grasp and produce lengthy language transcriptions and equations **<empty citation>** In the current thesis, the GPT-3 embeddings play an important role to automate the process of GDPR compliance checks and the subsequent section will explain in detail.

### 5.4.1 Generating GPT-3 Embeddings

GPT-3 embeddings contain a lot of language knowledge and contextual features that are useful for different NLP purposes such as text classification, especially textual summarizing and query answering. For this thesis, the content of those privacy policies and other related documents was encoded with GPT-3 embeddings. These embeddings were used as the initial features for the machine Learning algorithms.

The approach of creating embeddings with GPT-3 is achieved by passing a text to the model and obtaining the weights of the hidden layers. The steps are as

follows: The steps are as follows:

- **Input Text:** The privacy policy text is tokenized and fed into the GPT-3 model.
- **Hidden Layers:** It then passes the text through several transformer layers, and produces contextual representations from each layer.
- **Embedding Extraction:** The final hidden layer's output is extracted as the embedding in order to pick the semantic context of the input text.

GPT-3 embeddings provide several advantages when it comes to checking GDPR compliance, such as semantic understanding, which allows models to capture the essence of legal language; contextual sensitivity, which is vital to accurately interpret legal documents; and exposure to a wide arrange of datasets, promoting the model's ability to analyze various types of text.

However, they also include high computational cost and the need for major processing power hence they are costly. Moreover, even though GPT-3 has become more accessible with the help of GPT-3 API, its usage is limited with certain numbers of accesses and prices, which might be unachievable for some companies. Additionally, compared to other ML algorithms, the embeddings it generates can be difficult to understand, which complicates its usage.

In the experiments conducted for this thesis, GPT-3 embeddings were compared with other models such as SBERT and BERT. The findings indicated that GPT-3 embeddings provided the least impressive performance in terms of accuracy, precision, recall, and F1-score. One future work can be on improving the GPT-3 embeddings using the proposed models by integrating embeddings with other NLP approaches like rule based methods or domain specific fine tuning.

Regardless, after training a text classifier on the embeddings gotten from the API, to evaluate the model, the same full privacy policy was used on other models was split into sentences, then embeddings generated for each sentence and the trained model had to classify based on certain thresholds. The table can be seen below:

TABLE 5.8  
Threshold and Identified Principles from GPT-3 Embeddings

Threshold	Number of Identified Principles
0.5	7
0.6	5
0.7	5
0.73	3
0.8	0
0.9	0

On the whole, GDPR compliance checks enabled by GPT-3 embeddings open new opportunities for automating the needed checks and ensuring their semantics-oriented nature together with contextual awareness. As with any tool, there are limitations to be navigated, but the possible advantages give GPT-3 embeddings a place amongst tools that can help in the analysis of legal documents and checking of compliance.

## 5.5 Implications

As such, the practical applications of this research go beyond the mere theoretical enhancement of the existing literature in terms of GDPR compliance.

This section describes how the research questions formulated were solved by the experiments and methodologies used, as well as the practical nature of the generated compliance reports.

### **5.5.1 Addressing Research Questions**

Based on the outlined research questions, the focus of the thesis was on the efficiency of NLP models, as well as, GPT-3, in automatically detecting GDPR compliance, their limitations and use cases.

#### **Effectiveness of NLP Models**

The first research question was aimed at establishing the extent to which current NLP models enable the automation of GDPR compliance check of organizational data privacy policies. Experiments conducted with GPT-3, BERT, and other models, as well as legal NLP models developed in advance, proved that NLP technologies could improve the efficiency of compliance checks many times over. Sophisticated legal texts could be analyzed to find out the relevant GDPR principles and distinguish between the compliance statuses with high accuracy, precision, recall, and F1-score in these models.

#### **Limitations of Current NLP Technology**

The second research question was intended to establish the gap in the current NLP technology in analysing and enforcing conformity to GDPR. Some of the difficulties identified include; managing legal language, variations in context of different legal documents, and a lot of feature engineering to match the annotations to GDPR principles. Nevertheless, the study identified recommendations that

followed the three future research prospects, including complex preprocessing of NLP tasks, fine-tuning for legal-specific datasets, and transfer learning to boost the performance of candidates.

### Practical Applications of NLP Tools

The third research question was concerned with the use of tools backed by NLP in order to assist the compliance officers as well as the legal professionals. According to the findings of the study, it was revealed that NLP tools could help in cutting down the time and efforts needed to do compliance checks manually. In this case, legal compliance and the extract and analysis of compliance-related information from data privacy policies by these tools help legal experts to do more focused work, and make compliance management more effective and accurate in the end.

#### 5.5.2 Practicality of the Compliance Report

One of the main practical implications of this study is the creation of compliance reports following the analysis performed by the NLP models. These reports contain a clear evaluation of the organization's compliance with GDPR principles and noncompliance and recommend procedures on how to improve the situation.

##### **Sample 5.5.1: Example of a Compliance Report**

GDPR Compliance Report for Policy XYZ

Summary: Non-compliant with 2 out of 7 principles evaluated.

Detailed Findings:

1. Data Minimisation: Compliant.
2. Integrity and Confidentiality: Non-compliant.

Example: "User data may be stored indefinitely for analytics."

Recommendations:

- Review the data retention policy to align with the 'storage limitation' principle.

The above example report demonstrates the manner in which the conclusion and recommendation from the NLP analysis are given. Each principle is considered and the specific compliance is indicated with a sentence extract from the policy. Suggestions are then offered to enable organisations to make relevant corrections for the non-compliant aspects to fit the required standards.

### **Enhancing Transparency and Accountability**

The reports that are produced from the NLP models are well known to improve compliance and this in turn will lead to increase in organizational transparency and accountability. These reports can assist organizations to get a clear picture of their compliance situation hence be in a position to deal with all the compliance issues systematically. It further serves the purpose of not being on the wrong side of the law through fines or legal prosecutions while at the same time gaining the confidence of the clients and other stakeholders.

### **Continuous Compliance Monitoring**

The other related research implication that arises from this study is ongoing compliance monitoring. Since there are changes in GDPR regulations, it is

possible for NLP models to be modified according to the current trends hence improving the compliance of an organization. The second is the dynamic capability of NLP tools that assists organisations to maintain compliance work constantly and predict possible shifts in the regulation.

### **5.5.3 Contributions to the Field of Data Privacy and Compliance**

This research in the thesis benefits the area of data privacy and compliance by providing evidence of and illustrating methods to apply NLP technologies in the augmentation of the GDPR compliance. Hence, it offers a starting point for further research and enhancement of the more complex compliance-focused NLP tools relevant to the legal and regulatory environments with the potential of delivering better solutions in the future.

In conclusion, this study aims at answering the questions that arise in the area of research related to GDPR compliance and presents practical solutions in the form of NLP tools, but also stresses the need for constant further development and improvement in the field of managing personal data protection. Based on these studies, the prospects are quite vast, for we will witness vast enhancements in compliance assurance across most organizations.

## **5.6 Future Work and Improvements**

While this thesis has focused on automating GDPR compliance checks using several NLP models like SBERT, BERT, and GPT2, and it has revealed specific directions for future research and improvement. The possible future studies and

improvements that can be made are indicated in this section.

### 5.6.1 Enhancing Model Performance

Future research could improve the performance of the NLP models and this could be achieved through several axes such as:

- **Data Augmentation:** Making the dataset more diverse and recent, because the current dataset is from before 2020. In addition, the training data set can be augmented upon to create richer datasets.
- **Fine-tuning:** Increased fine-tuning of models on larger amounts of this data in order to better interpret legal language and text.
- **Model Architecture:** Trying new architectures and putting together different models to get higher accuracy, precision, recall and the F1 score.

### 5.6.2 Real-time Compliance Monitoring

Another important area for the future work is the procedures to be used to apply real-time compliance monitoring systems effectively. Such systems could be able to constantly search for and assess new policies or modifications to existing ones to make sure of compliance to GDPR, progressively. This involves:

- **Automated Pipelines:** Developing pipelines that can be integrated with actual industry data flows. This will allowing the compliance checks to be smoother; that is checked as soon as policies are created or updated.
- **Alert Systems:** Setting up alerts to inform the organization of the areas that might not be in compliance in real-time.



### 5.6.3 Broadening the Scope of Compliance Checks

Extending the application of the compliance check itself to other regulations concerning data privacy like CCPA, HIPAA, and the other laws is another potential lane for future study. This includes:

- **Multi-regulation Frameworks:** Developing frameworks that can assess compliance with multiple regulations simultaneously.
- **Cross-jurisdictional Analysis:** Facilitating easy comparison of compliance across jurisdictions so as to make compliance checking better for multinational corporations.

### 5.6.4 User-Centric Enhancements

Future developments should also aim at enhancing the way the tool can adapt to the needs of the compliance professionals, organizations or legal users. This involves:

- **Customization:** Allowing users to have tools that allow them to customize compliance checks based on their specific business needs.
- **Integration with Existing Tools:** Integration with other current legal and compliance tools that are already being used by the organizations.

All in all, despite the work of this thesis providing a good foundation for employing NLP models to automate GDPR checks, there are several possibilities for further research and manipulation. In these aforementioned areas, it is possible to enhance the potential of new tools and improve the experiences concerning the management of data protection regulations:

# Chapter 6

## Conclusion

In this thesis, the focus was on using NLP to automate compliance checking for the GDPR. This undertaking is inspired by such factors as the growing difficulty in legal matters and the need to address data protection laws in organizations effectively. The comprehensive analysis involved leveraging state-of-the-art NLP models, including SBERT, BERT, and GPT2, across two granularity levels: on the one hand, they are at the sentence level and, on the other hand, at the entire policy level.

### 6.1 Summary of Findings

Since the aim of this study was to evaluate the NLP models to determine the degree of compliance with GDPR in privacy policies, through rigorous experimentation and evaluation, several key findings emerged:

- **Model Performance:** From the models evaluated, it can therefore be deduced that while all the models were fairly effective to a certain extent, SBERT fared best at the sentence level, thereby providing high precision

and recall of compliance matters. There were also acceptable scores for BERT and GPT2 models, also, areas of relatively strong performance were established for BERT most majored in analyzing relationships in the text. The following scores summarize the performance:

- **SBERT:**
  - \* **Sentence Level:** Accuracy: 0.57, Precision: 0.54, Recall: 0.42, F1-score: 0.43
  - \* **Policy Level:** Accuracy: 0.60, Precision: 0.76, Recall: 0.85, F1-score: 0.80
- **BERT:**
  - \* **Sentence Level:** Accuracy: 0.52, Precision: 0.72, Recall: 0.42, F1-score: 0.52
  - \* **Policy Level:** Accuracy: 0.43, Precision: 0.80, Recall: 0.75, F1-score: 0.77
- **GPT2:**
  - \* **Sentence Level:** Accuracy: 0.62, Precision: 0.74, Recall: 0.46, F1-score: 0.55
  - \* **Policy Level:** Accuracy: 0.43, Precision: 0.79, Recall: 0.75, F1-score: 0.77
- **Granularity of Analysis:** Over the provided level of the sentence analysis, there were more accurate results concerning the compliance issues, which allowed performing a more close-up examination of each clause of the agreement. Nevertheless, the whole policy-level approach provided more general patterns of compliance that can be overseen on the detailed

level. When comparing the results obtained at the two different levels of granularity, it is noted that each is particularly relevant to a certain set of compliance regulatory measures.

- **GPT-3 Embeddings:** Through embedding GPT-3, there was an improvement in the contextual awareness of the models, therefore better prediction results. Still, prime usage of computational resources was demonstrated by GPT-3 embeddings for enhancing the precise and bitwise text comprehension, together with precisions and recalls achieved in the analysis of rather nuanced non-compliance risks.

## 6.2 Discussion

In general, the application of NLP to GDPR compliance checking is full of opportunities and risks as well. Pros of NLP models especially those using GPT-3 embeddings include semantic meaning retention, contextual, and a combine of broad pre-trained knowledge. These models can add a lot of to the speed and effectiveness of compliance checks meaning that fewer rely on simply reviewing the data manually.

However, challenges remain. The use of GPT-3 embeddings is time and computationally expensive, while APIs related to GPT-3 use and cost are other studied limitations. However, predictability and interpretableness are the two major challenges of such methods, as it is crucial to understand how the model came to such decision and if it complies with the law.

The future work regarding this issue should be aimed at furthering the identification and resolution of the mentioned challenges through the consideration

of the hybrid models, rule-based systems, and methods for the enhancement of the explainability of the NLP systems. This therefore makes it necessary that the embeddings and the models need to be constantly updated to the current standards that are set in law and the use of language today.

## 6.3 Practicality

The implication of the research results is that managers in organizations who wish to improve their GDPR compliance programmes can consider implementing the above measure. By implementing and deploying NLP-based solutions and products organizations can immensely enhance their compliance management by timely and correctly completing all the regulations and the related processes. The models and methods proposed in this thesis can be integrated into the current compliance systems and thus provide relatively easy to implement, efficient solutions.

Also, the discussion reveals that there is a need to integrate legal professionals and NLP scholars to advance the knowledge produced in the two domains. This kind of partnership can result in the creation of better and contextually enhanced compliance instruments, which in turn enhance data protection measures and privacies of persons.

## 6.4 Answering the Research Questions

This thesis set out to answer several key points, and the findings gotten give some answers and insights into the chosen research questions:

- **How effective are NLP models, including GPT-3, in automating the**

**identification of GDPR compliance issues within organizational data privacy policies?** The findings show that NLP models are very efficient and accurate with primary attention to the SI model including SBERT and BERT for comprehending compliance. The already nice-looking PDF-to-text conversion is also improved with the help of GPT-3 embeddings that bring deep semantic understanding and context awareness.

- **What are the limitations of current NLP technology in interpreting and enforcing GDPR compliance, and how can these limitations be addressed?** The primary drawbacks consist of computational complexity and the necessity of utilizing less cognitively complex models such as GPT-3, along with the issues of interpretability. Overcoming these limitations can be achieved by enhancing the computational speed, making the models more available and coming up with ways through which models can be easily understood.
- **What role can NLP-powered tools play in supporting compliance officers and legal experts in maintaining GDPR compliance?** It is also demonstrated that NLP-powered tools can assist compliance officers by automating those tasks as proper identification of compliance problems in organisation, which decreases the volume of workload of the officer and increases the level of reliability of compliance checks. These tools could be helpful in compliance which in its turn would free up the legal experts' time to perform more sophisticated tasks.

## 6.5 Contributions to the Field

With respect to the domain of data privacy and compliance, this thesis aids in establishing the use of advanced NLP methods in the sphere of achieving recurrent, intricate intellectual tasks. The comparison of different models produces insights into their advantages and shortcomings, and helps to outline the trends in the further investigation in this subject. Thus, bringing in GPT-3 embeddings can be viewed as a major upgrade, which demonstrates how using the latest in AI language processing can aid in processing legal text.

Overall, this thesis shows that NLP can bring positive changes to solving the issues that relate to the GDPR regulation. Through the usage of advanced language models and embedding techniques, the regulatory environment can be better understood and dealt with to enhance the information transparency and users' friendly approaches towards the data privacy.

# Bibliography cited

- [1] J. Hirschberg and C. D. Manning, “Advances in natural language processing,” *Science*, vol. 349, no. 6245, pp. 261–266, 2015.
- [2] *AI Review for Data Processing Agreements (DPAs)* — *legalontech.com*, <https://www.legalontech.com/contracts/data-processing-agreement-dpa>, [Accessed 13-06-2024].
- [3] O. Amaral Cejas, S. Abualhaija, and L. Briand, “MI-based compliance verification of data processing agreements against gdpr,” English, Sep. 2023. [Online]. Available: <https://orbilu.uni.lu/handle/10993/55408>.
- [4] C. Meehan, K. Mrini, and K. Chaudhuri, “Sentence-level privacy for document embeddings,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, S. Muresan, P. Nakov, and A. Villavicencio, Eds., Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 3367–3380. DOI: 10.18653/v1/2022.acl-long.238. [Online]. Available: <https://aclanthology.org/2022.acl-long.238>.
- [5] M. d. C. Freitas and M. Mira da Silva, “Gdpr compliance in smes: There is much to be done,” *Journal of Information Systems Engineering & Management*, vol. 3, no. 4, p. 30, 2018.



- [6] E. Arfelt, D. Basin, and S. Debois, “Monitoring the gdpr,” en, in *Computer Security – ESORICS 2019*, K. Sako, S. Schneider, and P. Y. A. Ryan, Eds., Cham: Springer International Publishing, 2019, pp. 681–699, ISBN: 978-3-030-29959-0. DOI: [10.1007/978-3-030-29959-0\\_33](https://doi.org/10.1007/978-3-030-29959-0_33).
- [7] Wikipedia contributors, *Text mining — Wikipedia, The Free Encyclopedia*, [Online; accessed 12-June-2024], 2023. [Online]. Available: [https://en.wikipedia.org/wiki/Text\\_mining](https://en.wikipedia.org/wiki/Text_mining).
- [8] Wikipedia contributors, *Transfer learning — Wikipedia, The Free Encyclopedia*, [Online; accessed 12-June-2024], 2023. [Online]. Available: [https://en.wikipedia.org/wiki/Transfer\\_learning](https://en.wikipedia.org/wiki/Transfer_learning).
- [9] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds., Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186. DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423). [Online]. Available: <https://aclanthology.org/N19-1423>.
- [10] M. Saqr, “Is gdpr failing? a tale of the many challenges in interpretations, applications, and enforcement,” *International Journal of Health Sciences*, vol. 16, no. 5, pp. 1–2, 2022, ISSN: 1658-3639.
- [11] D. Peloquin, M. DiMaio, B. Bierer, and M. Barnes, “Disruptive and avoidable: Gdpr challenges to secondary research uses of data,” en, *European Journal of Human Genetics*, vol. 28, no. 6, pp. 697–705, Jun. 2020, ISSN: 1476-5438. DOI: [10.1038/s41431-020-0596-x](https://doi.org/10.1038/s41431-020-0596-x).

- [12] A.-J. Aberkane, G. Poels, and S. V. Broucke, “Exploring automated gdpr-compliance in requirements engineering: A systematic mapping study,” *IEEE Access*, vol. 9, pp. 66 542–66 559, 2021, 5 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 2169-3536. DOI: [10.1109/ACCESS.2021.3076921](https://doi.org/10.1109/ACCESS.2021.3076921).
- [13] O. A. Cejas, M. I. Azeem, S. Abualhaija, and L. C. Briand, “Nlp-based automated compliance checking of data processing agreements against gdpr,” *IEEE Transactions on Software Engineering*, vol. 49, no. 9, pp. 4282–4303, Sep. 2023, 6 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1939-3520. DOI: [10.1109/TSE.2023.3288901](https://doi.org/10.1109/TSE.2023.3288901).
- [14] Ž. Spalević and K. Vićentijević, “Gdpr and challenges of personal data protection,” en, *The European Journal of Applied Economics*, vol. 19, no. 1, pp. 55–65, 2022, ISSN: 2406-2588, 2406-3215. DOI: [10.5937/EJAE19-36596](https://doi.org/10.5937/EJAE19-36596).
- [15] S. Sirur, J. R. Nurse, and H. Webb, “Are we there yet? understanding the challenges faced in complying with the general data protection regulation (gdpr),” in *Proceedings of the 2nd International Workshop on Multimedia Privacy and Security*, ser. MPS ’18, New York, NY, USA: Association for Computing Machinery, Jan. 2018, pp. 88–95, ISBN: 978-1-4503-5988-7. DOI: [10.1145/3267357.3267368](https://doi.org/10.1145/3267357.3267368). [Online]. Available: <https://doi.org/10.1145/3267357.3267368>.
- [16] R. E. Hamdani, M. Mustapha, D. R. Amariles, A. Troussel, S. Meeùs, and K. Krasnashchok, “A combined rule-based and machine learning approach for automated gdpr compliance checking,” en, in *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, São Paulo

- Brazil: ACM, Jun. 2021, pp. 40–49, ISBN: 978-1-4503-8526-8. DOI: 10.1145/3462757.3466081. [Online]. Available: <https://dl.acm.org/doi/10.1145/3462757.3466081>.
- [17] A. Nazarenko, F. Lévy, and A. Wyner, “A pragmatic approach to semantic annotation for search of legal texts – an experiment on gdpr,” E. Schweighofer, Ed., 3 citations (Semantic Scholar/DOI) [2023-12-08] Book Title: *Frontiers in Artificial Intelligence and Applications* DOI: 10.3233/FAIA210313, IOS Press, Dec. 2021, ISBN: 978-1-64368-252-5. DOI: 10.3233/FAIA210313. [Online]. Available: <https://ebooks.iospress.nl/doi/10.3233/FAIA210313>.
- [18] P. A. Bonatti, S. Kirrane, I. M. Petrova, and L. Sauro, “Machine understandable policies and gdpr compliance checking,” en, *KI - Künstliche Intelligenz*, vol. 34, no. 3, pp. 303–315, Sep. 2020, 27 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1610-1987. DOI: 10.1007/s13218-020-00677-4.
- [19] M. Galle, A. Christofi, and H. Elsahar, “The case for a gdpr-specific annotated dataset of privacy policies,” en,
- [20] O. Amaral, S. Abualhaija, M. Sabetzadeh, and L. Briand, “A model-based conceptualization of requirements for compliance checking of data processing against gdpr,” in *2021 IEEE 29th International Requirements Engineering Conference Workshops (REW)*, 2 citations (Semantic Scholar/DOI) [2023-12-08], Sep. 2021, pp. 16–20. DOI: 10.1109/REW53955.2021.00009. [Online]. Available: <https://ieeexplore.ieee.org/document/9582337>.

- [21] Z. S. Li, C. M. Werner, N. A. Ernst, and D. Damian, “Gdpr compliance in the context of continuous integration,” *ArXiv*, Feb. 2020. [Online]. Available: <https://www.semanticscholar.org/paper/71e16573d39360b98306b3bfa5482c10b4e7>
- [22] K. Mori, T. Nagai, Y. Takata, and M. Kamizono, “Analysis of privacy compliance by classifying multiple policies on the web,” en, in *2022 IEEE 46th Annual Computers, Software, and Applications Conference (COMPSAC)*, Los Alamitos, CA, USA: IEEE, Jun. 2022, pp. 1734–1741, ISBN: 978-1-66548-810-5. DOI: 10.1109/COMPSAC54236.2022.00276. [Online]. Available: <https://ieeexplore.ieee.org/document/9842614/>.
- [23] A. Qamar, T. Javed, and M. Beg, *Detecting Compliance of Privacy Policies with Data Protection Laws*. Feb. 2021.
- [24] S. Sousa and R. Kern, “How to keep text private? a systematic review of deep learning methods for privacy-preserving natural language processing,” en, *Artificial Intelligence Review*, vol. 56, no. 2, pp. 1427–1492, Feb. 2023, 12 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1573-7462. DOI: 10.1007/s10462-022-10204-6.
- [25] M. Srinath, S. Wilson, and C. L. Giles, “Privacy at scale: Introducing the privaseer corpus of web privacy policies,” no. arXiv:2004.11131, Apr. 2020, 28 citations (Semantic Scholar/arXiv) [2023-12-08] arXiv:2004.11131 [cs]. DOI: 10.48550/arXiv.2004.11131. [Online]. Available: <http://arxiv.org/abs/2004.11131>.
- [26] P. Silva, C. Gonçalves, C. Godinho, N. Antunes, and M. Curado, “Using natural language processing to detect privacy violations in online contracts,” in *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, ser. SAC ’20, 12 citations (Semantic Scholar/DOI) [2023-12-

- 08], New York, NY, USA: Association for Computing Machinery, Mar. 2020, pp. 1305–1307, ISBN: 978-1-4503-6866-7. DOI: [10.1145/3341105.3375774](https://doi.org/10.1145/3341105.3375774). [Online]. Available: <https://doi.org/10.1145/3341105.3375774>.
- [27] H. Harkous, K. Fawaz, R. Lebre, F. Schaub, K. G. Shin, and K. Aberer, “Polisis: Automated analysis and presentation of privacy policies using deep learning,” en, 2018, pp. 531–548, ISBN: 978-1-939133-04-5. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity18/presentation/harkous>.
- [28] V. Leone and L. Di Caro, “The role of vocabulary mediation to discover and represent relevant information in privacy policies,” en, in *Frontiers in Artificial Intelligence and Applications*, S. Villata, J. Harašta, and P. Křemen, Eds. IOS Press, Dec. 2020, ISBN: 978-1-64368-150-4. DOI: [10.3233/FAIA200851](https://doi.org/10.3233/FAIA200851). [Online]. Available: <http://ebooks.iospress.nl/doi/10.3233/FAIA200851>.
- [29] N. M. Müller, D. Kowatsch, P. Debus, D. Mirdita, and K. Böttinger, “On gdpr compliance of companies’ privacy policies,” en, in *Text, Speech, and Dialogue*, K. Ekšte, Ed., ser. Lecture Notes in Computer Science, 12 citations (Semantic Scholar/DOI) [2023-12-08], Cham: Springer International Publishing, 2019, pp. 151–159, ISBN: 978-3-030-27947-9. DOI: [10.1007/978-3-030-27947-9\\_13](https://doi.org/10.1007/978-3-030-27947-9_13).
- [30] H. T. Alattas, F. M. Almassary, N. R. AlMahasheer, *et al.*, “Extract compliance-related evidence using machine learning,” in *2022 14th International Conference on Computational Intelligence and Communication Networks (CICN)*, 0 citations (Semantic Scholar/DOI) [2023-12-08], Dec. 2022,

- pp. 537–542. DOI: 10.1109/CICN56167.2022.10008324. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10008324>.
- [31] D. Torre, S. Abualhaija, M. Sabetzadeh, *et al.*, “An ai-assisted approach for checking the completeness of privacy policies against gdpr,” in *2020 IEEE 28th International Requirements Engineering Conference (RE)*, 33 citations (Semantic Scholar/DOI) [2023-12-08], Aug. 2020, pp. 136–146. DOI: 10.1109/RE48521.2020.00025. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9218152>.
- [32] S. Arora, H. Hosseini, C. Utz, *et al.*, “A tale of two regulatory regimes: Creation and analysis of a bilingual privacy policy corpus,” in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, Marseille, France: European Language Resources Association, Jun. 2022, pp. 5460–5472. [Online]. Available: <https://aclanthology.org/2022.lrec-1.585>.
- [33] M. Bokaie Hosseini, P. K C, I. Reyes, and S. Egelman, “Identifying and classifying third-party entities in natural language privacy policies,” in *Proceedings of the Second Workshop on Privacy in NLP*, O. Feyisetan, S. Ghanavati, S. Malmasi, and P. Thaine, Eds., 7 citations (Semantic Scholar/DOI) [2023-12-08], Online: Association for Computational Linguistics, Nov. 2020, pp. 18–27. DOI: 10.18653/v1/2020.privatenlp-1.3. [Online]. Available: <https://aclanthology.org/2020.privatenlp-1.3>.
- [34] H. Harkous, K. Fawaz, R. Lebrete, F. Schaub, K. G. Shin, and K. Aberer, “Polisis: Automated analysis and presentation of privacy policies using deep learning,” en,
- [35] Y. Ling, K. Wang, G. Bai, H. Wang, and J. S. Dong, “Are they toeing the line? diagnosing privacy compliance violations among browser exten-

- sions,” in *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*, ser. ASE '22, 7 citations (Semantic Scholar/DOI) [2023-12-08], New York, NY, USA: Association for Computing Machinery, Jan. 2023, pp. 1–12, ISBN: 978-1-4503-9475-8. DOI: 10.1145/3551349.3560436. [Online]. Available: <https://dl.acm.org/doi/10.1145/3551349.3560436>.
- [36] D. Sánchez, A. Viejo, and M. Batet, “Automatic assessment of privacy policies under the gdpr,” en, *Applied Sciences*, vol. 11, no. 4, p. 1762, Feb. 2021, 10 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 2076-3417. DOI: 10.3390/app11041762.
- [37] P. Silva, C. Gonçalves, C. Godinho, N. Antunes, and M. Curado, “Using nlp and machine learning to detect data privacy violations,” in *IEEE INFOCOM 2020 - IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS)*, 14 citations (Semantic Scholar/DOI) [2023-12-08], Jul. 2020, pp. 972–977. DOI: 10.1109/INFOCOMWKSHPS50562.2020.9162683. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/9162683>.
- [38] O. Amaral, S. Abualhaija, D. Torre, M. Sabetzadeh, and L. C. Briand, “Ai-enabled automation for completeness checking of privacy policies,” *IEEE Transactions on Software Engineering*, vol. 48, no. 11, pp. 4647–4674, Nov. 2022, 12 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1939-3520. DOI: 10.1109/TSE.2021.3124332.
- [39] T. A. Rahat, M. Long, and Y. Tian, “Is your policy compliant? a deep learning-based empirical study of privacy policies’ compliance with gdpr,” in *Proceedings of the 21st Workshop on Privacy in the Electronic Society*,



- ser. WPES'22, 1 citations (Semantic Scholar/DOI) [2023-12-08] 1 citations (Crossref) [2023-12-04], New York, NY, USA: Association for Computing Machinery, Nov. 2022, pp. 89–102, ISBN: 978-1-4503-9873-2. DOI: 10.1145/3559613.3563195. [Online]. Available: <https://dl.acm.org/doi/10.1145/3559613.3563195>.
- [40] C. Bartolini, A. Giurciu, G. Lenzini, and L. Robaldo, “A framework to reason about the legal compliance of security standards,” English, Nov. 2016. [Online]. Available: <https://orbilu.uni.lu/handle/10993/28786>.
- [41] P. A. Bonatti, S. Kirrane, I. M. Petrova, and L. Sauro, “Machine Understandable Policies and GDPR Compliance Checking,” en, *KI - Künstliche Intelligenz*, vol. 34, no. 3, pp. 303–315, Sep. 2020, 27 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1610-1987. DOI: 10.1007/s13218-020-00677-4. [Online]. Available: <https://doi.org/10.1007/s13218-020-00677-4> (visited on 11/20/2023).
- [42] D. Bui, J.-M. Choi, and J. Shin, “Automated extraction and presentation of data practices in privacy policies,” *Proceedings on Privacy Enhancing Technologies*, vol. 2021, pp. 88–110, Apr. 2021, 31 citations (Semantic Scholar/DOI) [2023-12-08]. DOI: 10.2478/popets-2021-0019.
- [43] A. Sleimi, N. Sannier, M. Sabetzadeh, L. Briand, and J. Dann, “Automated extraction of semantic legal metadata using natural language processing,” in *2018 IEEE 26th International Requirements Engineering Conference (RE)*, 48 citations (Semantic Scholar/DOI) [2023-12-08], Aug. 2018, pp. 124–135. DOI: 10.1109/RE.2018.00022. [Online]. Available: <https://ieeexplore.ieee.org/document/8491129?denied=>.



- [44] A. Nazarenko, F. Lévy, and A. Wyner, “A Pragmatic Approach to Semantic Annotation for Search of Legal Texts – An Experiment on GDPR,” E. Schweighofer, Ed., 3 citations (Semantic Scholar/DOI) [2023-12-08] Book Title: *Frontiers in Artificial Intelligence and Applications*, IOS Press, Dec. 2021, ISBN: 978-1-64368-252-5 978-1-64368-253-2. DOI: [10.3233/FAIA210313](https://ebooks.iospress.nl/doi/10.3233/FAIA210313). [Online]. Available: <https://ebooks.iospress.nl/doi/10.3233/FAIA210313> (visited on 11/20/2023).
- [45] J. M. Del Alamo, D. S. Guaman, B. García, and A. Diez, “A systematic mapping study on automated analysis of privacy policies,” en, *Computing*, vol. 104, no. 9, pp. 2053–2076, Sep. 2022, 10 citations (Semantic Scholar/DOI) [2023-12-08], ISSN: 1436-5057. DOI: [10.1007/s00607-022-01076-3](https://doi.org/10.1007/s00607-022-01076-3).
- [46] R. E. Hamdani, M. Mustapha, D. R. Amariles, A. Troussel, S. Meeùs, and K. Krasnashchok, “A combined rule-based and machine learning approach for automated gdpr compliance checking,” in *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, ser. ICAIL ’21, 14 citations (Semantic Scholar/DOI) [2023-12-08], New York, NY, USA: Association for Computing Machinery, Jul. 2021, pp. 40–49, ISBN: 978-1-4503-8526-8. DOI: [10.1145/3462757.3466081](https://doi.org/10.1145/3462757.3466081). [Online]. Available: <https://doi.org/10.1145/3462757.3466081>.
- [47] E. Poplavska, T. Norton, S. Wilson, and N. Sadeh, “From prescription to description: Mapping the gdpr to a privacy policy corpus annotation scheme,” in Dec. 2020, ISBN: 978-1-64368-150-4. DOI: [10.3233/FAIA200874](https://doi.org/10.3233/FAIA200874).
- [48] N. Mousavi Nejad, P. Jabat, R. Nedelchev, S. Scerri, and D. Graux, “Establishing a strong baseline for privacy policy classification,” en, in *ICT Systems Security and Privacy Protection*, M. Hölbl, K. Rannenberg, and

- T. Welzer, Eds., ser. IFIP Advances in Information and Communication Technology, 21 citations (Semantic Scholar/DOI) [2023-12-08], Cham: Springer International Publishing, 2020, pp. 370–383, ISBN: 978-3-030-58201-2. DOI: [10.1007/978-3-030-58201-2\\_25](https://doi.org/10.1007/978-3-030-58201-2_25).
- [49] S. Liu, B. Zhao, R. Guo, G. Meng, F. Zhang, and M. Zhang, “Have you been properly notified? automatic compliance analysis of privacy policy text with gdpr article 13,” in *Proceedings of the Web Conference 2021*, ser. WWW ’21, 24 citations (Semantic Scholar/DOI) [2023-12-08], New York, NY, USA: Association for Computing Machinery, Jun. 2021, pp. 2154–2164, ISBN: 978-1-4503-8312-7. DOI: [10.1145/3442381.3450022](https://doi.org/10.1145/3442381.3450022). [Online]. Available: <https://doi.org/10.1145/3442381.3450022>.
- [50] J. Giner-Miguel, A. Gómez, and J. Cabot, “Datadoc analyzer: A tool for analyzing the documentation of scientific datasets,” en, in *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 0 citations (Semantic Scholar/DOI) [2023-12-08], Birmingham United Kingdom: ACM, Oct. 2023, pp. 5046–5050, ISBN: 9798400701245. DOI: [10.1145/3583780.3614737](https://doi.org/10.1145/3583780.3614737). [Online]. Available: <https://dl.acm.org/doi/10.1145/3583780.3614737>.
- [51] *Usable Privacy Policy Project*—[usableprivacy.org](https://usableprivacy.org), <https://www.usableprivacy.org/data>, [Accessed 13-06-2024].
- [52] A. Ravichander, A. W. Black, S. Wilson, T. Norton, and N. Sadeh, “Question answering for privacy policies: Combining computational and legal perspectives,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, China:

- Association for Computational Linguistics, Nov. 2019, pp. 4949–4959. DOI: 10.18653/v1/D19-1500. [Online]. Available: <https://www.aclweb.org/anthology/D19-1500>.
- [53] R. Ramanath, F. Liu, N. Sadeh, and N. Smith, “Unsupervised alignment of privacy policies using hidden markov models,” in *Proceedings of ACL*, Association for Computational Linguistics, Jun. 2014.
- [54] E. Poplavska, T. B. Norton, S. Wilson, and N. M. Sadeh, “From prescription to description: Mapping the gdpr to a privacy policy corpus annotation scheme,” in *International Conference on Legal Knowledge and Information Systems*, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:229377855>.
- [55] I. L. Nicolaidou and C. Georgiades, “The gdpr: New horizons,” en, in *EU Internet Law: Regulation and Enforcement*, T.-E. Synodinou, P. Jogleux, C. Markou, and T. Prastitou, Eds. Cham: Springer International Publishing, 2017, pp. 3–18, ISBN: 978-3-319-64955-9. DOI: 10.1007/978-3-319-64955-9\_1. [Online]. Available: [https://doi.org/10.1007/978-3-319-64955-9\\_1](https://doi.org/10.1007/978-3-319-64955-9_1).

# Appendix A

## Full Policy Report

The best model shows that the policy below adheres to the following GDPR principles:

1. Lawfulness, Fairness and Transparency
2. Purpose Limitation
3. Data Minimization
4. Accuracy
5. Integrity and Confidentiality

### Sample A.0.1: Result Yola Privacy Policy

yola, inc. privacy policy last updated january 1, 2013 yola is a licensee of the trustee privacy program. trustee is an independent, organization whose mission is to build user's trust and confidence in the internet by promoting the use of fair information practices. this privacy statement covers the web site www.yola.com. because this web site wants to demonstrate its commitment to your privacy, it has agreed to disclose its information practices and have

its privacy practices reviewed for compliance by trustee. if you have questions or concerns regarding this statement, you should first contact [support@yola.com](mailto:support@yola.com). if you do not receive acknowledgement of your inquiry or your inquiry has not been satisfactorily addressed, you should contact trustee. trustee will then serve as a liaison with us to resolve your concerns. this privacy policy ("policy") explains how personal information is collected, used, and disclosed by yola, inc., also doing business as synthasite, inc., and synthasite (pty) ltd. (collectively "yola") with respect to your use of the web site located at [www.yola.com](http://www.yola.com) (the "site") so you can make an informed decision about using the site and the web site creation and hosting services offered by yola. please note that this policy does not apply to, and we are not responsible in any manner for, any information you provide to third parties in connection with web sites hosted by yola ("user sites") or how the third parties who operate such user sites may use the personal information you provide to them. such information is subject to the privacy practices of such user site, and we encourage you to become familiar with their privacy practices before disclosing information directly to them. registration forms contain both mandatory fields (as indicated) and non-mandatory fields; if you choose not to complete the mandatory fields then yola may not be able to respond to your request or activate your option choices. by using the yola toolset and/or the yola web site, you consent to any transfer of personal information, collected by yola, outside your country for the purposes of storing the information where yola and/or its agents maintain their facilities. we reserve the right to change this policy at any time. if we make any material changes to the way we use your personally identifiable information, we will notify you here, post a new policy on our site and update the "last updated" date set forth above, or by email. therefore, we encourage you to check our policy whenever you use the service to see if it has been updated since your last visit. your use of the site after the revised policy has been posted will constitute your consent to such revised policy. if you are concerned about how your personal information is used, you should subscribe to our blog rss feed at <http://www.yola.com/blog> or e-mail us at [support@yola.com](mailto:support@yola.com). this will not prevent yola from using your personal information for the administration or operation of the site or other nonmarketing purposes described in this policy. our postal address is: 201 mission st, suite 2250, san francisco, ca 94105 we can be reached via e-mail at [support@yola.com](mailto:support@yola.com). we collect two basic types of information with our service,

namely personal information and aggregate information. as used herein, the term “personal information” means information that specifically identifies an individual (such as a name and e-mail address), and demographic and other information about an individual when directly linked to personally identifiable information. a user name may be personal information if it includes personally identifying information such as a first and last name. our definition of personal information does not include “aggregate” information. aggregate information is data we collect about a group or category of services or users from which individual user identities have been removed. in other words, information on how you use our service may be collected and combined with information about how others use the service, but no personal information will be included in the resulting data. aggregate data helps us understand trends in our users’ needs so that we can better consider new features or otherwise tailor our service. this policy in no way restricts or limits our collection and use of aggregate information, and we may share aggregate information about our users with third parties for various purposes, including to help us better understand and improve our service, and for advertising and marketing purposes. if your personally identifiable information changes, or if you no longer desire our service, you may update your profile or delete it by clicking on the profile link after you log in and then clicking the delete account link at the bottom of the page. alternatively, you can email our customer support at support@yola.com. google apps™ users: deleting the yola app from your google apps™ control panel does not delete your yola account. to delete your yola account, please follow the instructions provided above. personal information may be collected in a number of ways when you visit our site. we may collect certain information you voluntarily provide to us which may contain personal information. for example, we may collect your name, address, email address, user name and other contact and demographic information when you register and set up an account or contact us by e-mail or other means for any reason. we may also collect payment information if applicable. if in addition, from time to time we may collect demographic, contact or other personal information you voluntarily provide to us, such as in connection with your participation in surveys, sweepstakes, contests, games, promotional offers, and other activities on the site. for each registered user of the site (each a “registered user”) to the yola web site, our web server automatically recognizes any header information that is shared with our web server. yola collects the ip address and header

information, and may collect the following: the e-mail address of registered users of our web site and the yola toolset; the e-mail addresses of registered users who post messages to our bulletin board, communicate with us via e-mail, and/or make postings to our chat areas; aggregate information on what pages registered users access or visit; registered user-specific information on what pages registered users access or visit; and personal information may also be collected if you leave comments or other communications on the site, such as on the yola bulletin board and/or chat areas, and/or send email directly to yola. by doing so your name and the content of your communication may be made public and can be read, collected, or used by other users of these forums, and could be used to send you unsolicited messages. we are not responsible for the personally identifiable information you choose to submit in these forums. as mentioned above, we automatically receive certain types of information whenever you interact with us. for example, when you use the web site, our systems may automatically collect your ip address and the type of operating system or browser you use. we may also collect information pertaining to your account activity, and standard access information, such as the time and date of your accessing the service and your usage of the service. we use such information for purposes such as compiling aggregated statistics about service usage. we may also use cookie technology to collect information. among other things, the use of cookies enables us to: store registered user's preferences; record session information, such as items that visitors and registered users add to their shopping cart; record registered user-specific information on what pages registered users access or visit; alert registered users to new areas that we think might be of interest to them when they return to our site; record past activity at a site in order to provide better service when visitors and registered users return to our site; and deliver advertising targeted to your interests and ensure that visitors and registered users are not repeatedly served the same messages or offers. third party vendors, including google, show your ads on sites on the internet. third party vendors, including google, use cookies to serve ads based on a user's prior visits to your website. users may opt out of google's use of cookies by visiting the google advertising opt-out page. a cookie is a small amount of data, often including an anonymous unique identifier, which is sent to your browser from a web site's computers and stored on your computer's hard drive. most browsers automatically accept cookies as the default setting. you can modify the setting to reject cookies or to prompt

you before accepting a cookie from the sites you visit by editing browser options. however, if a browser is set not to accept cookies or if a user rejects a cookie, some portions of the site and the service may not function properly. for example, you may not be able to sign in and may not be able to access certain site features or services. the use of cookies by our partners and affiliates is not covered by our privacy statement. we do/do not have access or control over these cookies. our partners and affiliates use session id cookies on their websites once you leave yola to give you access to yola discounts and promotions. if you click on a link to a third party site, you will leave the yola site you are visiting and go to the site you selected. because we cannot control the activities of third parties, we cannot accept responsibility for any use of your personally identifiable information by such third parties, and we cannot guarantee that they will adhere to the same privacy practices as yola. if you visit a third party website that is linked to from a yola site, you should read that website's privacy statement before providing any personally identifiable information. when you use sharethis to bookmark and share our products and site, you will be taken to a third party site and note that their privacy policy governs the collection and use of the information collected on those sites. if you choose to use the get satisfaction's referral service on <http://forum.yola.com> to share a support request about our site, get satisfaction will ask you for your friend's name and email address. they will automatically send your friend a one-time email inviting him or her to visit the <http://forum.yola.com>. the get satisfaction referral program is not associated with yola and yola does not store or use that information. our service is not intended to be used by children under the age of 13, and we do not knowingly collect personal information from children under the age of 13 except in compliance with applicable law. in general, we use personal information we collect to process your requests or transactions, to provide you with information or services you request, to inform you about other information, events, promotions, products or services we think will be of interest to you, to facilitate your use of, and our administration and operation of, the site, newsletters and for the purpose for which the information was provided. we may also use the personal information we collect to improve the content of our web page and to customize the advertising and the content and/or layout of our page for our visitors and individual registered users. we also use this information to notify registered users about updates to our web site, and to contact registered users for marketing



purposes. we will not share the personal information we collect from you through the site with any other company except as provided in this policy, as it may be revised from time to time. we reserve the right to disclose or transfer such information to a third party (a) if you request us to do so or we otherwise have your consent, (b) if we believe in good faith that we are lawfully authorized or required to do so or that doing so is reasonably necessary or appropriate to comply with the law or with legal process or authorities, respond to any claims, or to protect the rights, property or safety of yola, our users, our employees or the public, including without limitation to protect yola or our users or the public from fraudulent, abusive, inappropriate or unlawful use of our site, (c) if we believe in good faith that an emergency involving immediate danger of death or serious physical injury to any person, or other irreparable injury, including economic interests of yola or a third party, requires disclosure of the information, and (d) in connection with an acquisition, merger, or sale of all or a substantial portion of our business. we may also share personal information with vendors, consultants and other service providers (“service providers”) who are engaged by or working with us in connection with the operation of the site or the services and who need access to such information to carry out their work for us. in some cases, the service provider may be directly collecting the information from you on our behalf. we may also share personal information with other third parties when you request that we do so or otherwise give us your express or implied consent. we are not responsible for the actions of service providers or other third parties, nor are we responsible for any additional information you provide directly to these service providers or other third parties, and we encourage you to become familiar with their privacy practices before disclosing information directly to them. note that nothing herein restricts the sharing of aggregate information, which may be shared with third parties without your consent. we use reasonable efforts to prevent unauthorized release of or access to your personal data. however, we cannot guarantee that your information will not be disclosed or accessed by accidental circumstances or by the unauthorized acts of others. if yola learns of a security systems breach, then we may attempt to notify you electronically so that you can take appropriate protective steps. yola may also post a notice on the site if a security breach occurs. depending on where you live, you may have a legal right to receive notice of a security breach in writing. to receive a free written notice of a security breach, you should notify us

at support@yola.com. the security of your personal information is important to us. when you enter sensitive information (such as credit card number and/or social security number, national id, personal health information) on our registration or order forms, we encrypt that information using secure socket layer technology (ssl). credit card details are stored in our payment gateway's highly secure payment vault that is fully compliant with the payment card industry's data security standards. we will communicate with you via email. yola intends to send regular email updates and occasional promotional offers that may be of interest to you. you may opt-out of these mailings at any time. if you do not want to receive marketing and promotional e-mails from us in the future, please let us know by sending us an e-mail to support@yola.com or follow the unsubscribe instructions that are included in each email communication. this will not prevent us from sending you operational, account related or other non-marketing emails. if you wish to no longer receive these emails please deactivate your account by clicking on the profile link after you log in and then clicking the delete account link at the bottom of the page. alternatively, you can email our customer support at support@yola.com. if you feel that this site is not following its stated information policy, you may contact us at support@yola.com.