# LADFA: A Framework of Using Large Language Models and Retrieval-Augmented Generation for Personal Data Flow Analysis in Privacy Policies

HAIYUE YUAN*, Institute of Cyber Security for Society (iCSS) & School of Computing, University of Kent, United Kingdom

NIKOLAY MATYUNIN and ALI RAZA, Honda Research Institute Europe GmbH, Germany

SHUJUN LI*, Institute of Cyber Security for Society (iCSS) & School of Computing, University of Kent, United Kingdom

Privacy policies help inform people about organisations' personal data processing practices, covering different aspects such as data collection, data storage, and sharing of personal data with third parties. Privacy policies are often difficult for people to fully comprehend due to the lengthy and complex legal language used and inconsistent practices across different sectors and organisations. To help conduct automated and large-scale analyses of privacy policies, many researchers have studied applications of machine learning and natural language processing techniques, including large language models (LLMs). While a limited number of prior studies utilised LLMs for extracting personal data flows from privacy policies, our approach builds on this line of work by combining LLMs with retrieval-augmented generation (RAG) and a customised knowledge base derived from existing studies. This paper presents the development of LADFA, an end-to-end computational framework, which can process unstructured text in a given privacy policy, extract personal data flows and construct a personal data flow graph, and conduct analysis of the data flow graph to facilitate insight discovery. The framework consists of a pre-processor, an LLM-based processor, and a data flow post-processor. We demonstrated and validated the effectiveness and accuracy of the proposed approach by conducting a case study that involved examining ten selected privacy policies from the automotive industry. Moreover, it is worth noting that LADFA is designed to be flexible and customisable, making it suitable for a range of text-based analysis tasks beyond privacy policy analysis.

CCS Concepts: • **Security and privacy → Human and societal aspects of security and privacy**; **Social network security and privacy**; • **Information systems → World Wide Web**.

Additional Key Words and Phrases: Large Language Model, LLM, Privacy Policy, Text Analysis, Data Flows, Privacy, Security, Retrieval-Augmented Generation, RAG, Framework, Automotive Industry, Connected Vehicle

---

*Corresponding co-authors.

---

Authors' Contact Information: Haiyue Yuan, h.yuan-221@kent.ac.uk, Institute of Cyber Security for Society (iCSS) & School of Computing, University of Kent, Canterbury, United Kingdom; Nikolay Matyunin, nikolay.matyunin@honda-ri.de; Ali Raza, ali.raza@honda-ri.de, Honda Research Institute Europe GmbH, Germany; Shujun Li, s.j.li@kent.ac.uk, Institute of Cyber Security for Society (iCSS) & School of Computing, University of Kent, Canterbury, United Kingdom.

---

## 1 Introduction

Privacy policies are widely used as a means of informing people about personal processing data practices of organisations. They are considered essential documents that detail information about how organisations collect, share, and manage personal data of people (i.e., data subjects in legal terms). Privacy policies are often subject to legal requirements, such as the General Data Protection Regulation (GDPR) in the EU [55] and the UK [57] and various state-level regulations in the US (e.g., the California Consumer Privacy Act (CCPA) [52] and the California Privacy Rights Act (CPRA) [51]). While these regulations mandate organisations to disclose specific information, their implementation often lacks consistency and transparency [48, 58]. Moreover, due to the lengthy and sophisticated legal language used in privacy policies, human readers often struggle to fully understand the scale and scope of data collection and sharing practices or often choose not to read them at all [11, 35, 54].

To address these challenges, some researchers have examined privacy policy consistency, clarity, and transparency to better inform and present meaningful information to human readers. Machine learning (ML) and natural language processing (NLP) techniques have been frequently explored to facilitate such analyses [1, 16, 23, 41]. More recently, with the increasing capabilities and popularity of large language models (LLMs), several studies have shown promising results of leveraging LLMs to automate the privacy policy analysis process in zero- and/or few-shot training contexts [13, 20, 36, 46, 53].

Among all information described in a privacy policy, personal data flows are of particular importance as they can tell data subjects what personal data about them will be collected by whom and with what third parties the collected personal data will be shared, under which conditions and for what purposes. Such personal data flows can be understood as information flows in the contextual integrity (CI) theory, proposed by Nissenbaum [38]. In the CI theory, privacy can be modelled as personal data flows, which are governed by contextual norms. The key parameters of such norms are *actors* (including the *subject* of information, the *sender* of information, and *recipient* of information), *attributes* (i.e., types of information), and *transmission principle* [38, 39], which refers to the condition or constraints that govern the personal data flow, restricting it to specific circumstances (e.g., a business should only give customer records to the government if there is a warrant or court order) [33].

These parameters are essential to inform data subjects about the personal data collected, with whom it is shared, and for what purposes. Although 'purposes' is not explicitly defined in the original CI theory, Nissenbaum [40] later explained that factors such as purpose do not exist in a context but help constitute it. In this sense, we consider purpose as part of *transmission principles*. For the readability and consistency, we refer to personal data flow as data flow for the rest of this paper and refer to data flows that include these attributes as comprehensive data flows. Although a few studies have explored manual [19, 65] and automatic approaches [15, 63, 64] to extract data flows from privacy policies and construct data flow graphs for visualisation and insights discovery, to the best of our knowledge, no existing study has utilised LLMs with RAG and a customised knowledge base to automatically extract comprehensive data flows from the perspective of CI theory.

To address the above-mentioned research gap, we conducted a study with following main contributions:

- **Novel framework**: we propose an end-to-end privacy policy analysis framework LADFA (short for "A Framework using **L**LMs and R**A**G for Personal **D**ata **F**low **A**nalysis in Privacy Policies"), focusing on extracting data flows, constructing and analysing data flow graphs. It consists of a pre-processor, an LLM-based processor, and a data flow post-processor.

- **Customised knowledge base**: beyond utilising LLMs, we design and construct a knowledge base for providing useful contextual information that allow the framework to leverage retrieval-augmented generation (RAG) to facilitate the LLM-based processor and the post-processor.
- **Case study with evaluation**: we conduct a case study on ten connected-vehicle mobile apps from different original equipment manufacturers (OEMs) in the automotive industry. Unlike prior work that evaluated only segments of data flows separately, we provide a collective evaluation of comprehensive data flows. Due to the lack of ground truth, manual validation was performed by three domain experts (the first three co-authors of the paper). The evaluation results showed strong agreement with the LADFA's outputs, with the average 7-Likert scores between 6 and 7 for most tasks. Gwet's AC1 and percentage agreement for identifying data types and data flows reached 0.94 and 0.82, and 0.96 and 0.86, respectively, indicating high inter-rater reliability. These results affirm LADFA's capability in processing, understanding unstructured texts, and extracting data flows from privacy policies.
- **Insights discovery**: the comparison of data flow graphs and graph network analysis demonstrated LADFA's capability in analysing and discovering privacy- and security-related insights that are often difficult to comprehend or easily overlooked by human readers of privacy policies.

The rest of this paper is organised as follows. Related work is discussed in Section 2. Section 3 provides detailed descriptions of each component of the proposed end-to-end framework for automated privacy policy analysis. Section 4 presents the details of the case study and the results, followed by a discussion of some identified limitations in Section 5. Finally, the last section concludes the paper.

## 2 Related Work

### 2.1 Privacy Policy Analysis in General

Previous research has indicated that a significant proportion of organisations or services do not adequately inform users about data-handling practices through privacy policies. A large-scale analysis conducted in 2017 highlighted that approximately 50% of popular free apps available on the Google Play app store were missing a privacy policy, despite the fact that a large portion of them (70%) are capable of processing personally identifiable information [27]. Similarly, another study [43] found that most smartphone apps do not include privacy policies: for the privacy policies that were included, only 18% of the iOS apps could be accessed, and accessibility was even lower for the Android apps (approximately 4%). Apart from such accessibility challenges, known issues related to the readability, presentation, transparency, and consistency of privacy policies often prevent more effective and informative communication to consumers [12, 14, 43]. Consequently, consumers rarely read privacy policies and perceive them as overly complex, lengthy, and difficult to comprehend, resulting in uninformed consent to data practices [31, 41, 58, 60, 65]. Such challenges not only undermine user trust but also complicate regulatory compliance for organisations [16]. More recently, Ghahremani and Nguyen [19] presented a study to systematically evaluate the transparency and comprehensiveness of privacy policies by manually annotating and identifying contextual gaps and ambiguities based on the CI framework. They argued that such manual analysis can be considered an alternative to subjective evaluations by privacy experts; however, there is a pressing need to automate the process using ML and NLP-based approaches, especially for large-scale studies and for developing an automated tool to assist human users.

## 2.2 Automated Privacy Policy Analysis

*2.2.1 ML and NLP-based Approaches.* In addition to the challenges listed above, it is essential to analyse privacy policy content to evaluate its completeness and alignment with regulatory frameworks and improve its readability, presentation, consistency, and clarity. Given the complexity and volume of text in privacy policies, researchers have explored various automated approaches that utilise ML and NLP techniques for such analyses.

Several studies have investigated and confirmed the readability issues of privacy policies from different domains. Fabian et al. [16] developed an automated toolset that utilised NLP techniques for information extraction and readability analysis. They examined 50,000 privacy policies from popular English-speaking websites and adopted the Flesh-Kincaid Grade Level (FKG) metric, a well-known metric for evaluating readability. The mean FKG score was 13.6, indicating that, on average, these documents remain difficult for human readers to comprehend. In a different study, Srinath et al. [50] presented the PrivaSeer corpus, which contains 1,005,380 privacy policies from 995,475 web domains. The readability analysis of this dataset yielded a mean FKG score of 14.87, suggesting a need for roughly two years of US college education to fully grasp a typical privacy policy. In another study on the privacy policies of mental health apps, Robillard et al. [43] calculated average scores from three established readability metrics, including Gunning Fog, FKG, and the Simple Measure of Gobbledygook (SMOG). The results reveal that these legal documents require a reading level equivalent to that of a college student, making them difficult for the average user to comprehend fully.

Such readability issues are often caused by complex legal language and the manner in which the content is presented. Numerous studies have focused on bridging this gap using various approaches, including innovative visual representations, nudging techniques, question-answering systems, and summarisation tools, to facilitate more effective and informative communication with consumers. For instance, Oltramari et al. [41] proposed PrivOnto, a semantic framework that combines crowdsourcing, ML and NLP that can represent annotated privacy policies using an ontology, allowing the development of SPARQL queries to extract information from the PrivOnto knowledge base to address user privacy-related questions and assist researchers and regulators in large-scale privacy policy analysis. Zaeem et al. [67] developed a Chrome browser extension, PrivacyCheck, which utilises trained data mining classification models using 400 privacy policies. It allows summarising an HTML-based privacy policy and presents the results as graphical icons with short descriptions and risk-level indications. Similarly, Harkous et al. [23] developed an automated framework, Polisis, which leverages a privacy-centric language model trained on 130 K privacy policies and a neural network classifier to analyse both high-level and fine-grained details of privacy practices. Polisis can automatically assign privacy icons to privacy policies, allowing human readers to learn more privacy insights for making more informed decisions. In addition, Harkous et al. [23] introduced PriBot based on Polisis. It is the first interactive question-answering system for privacy policies, offering consumers a more dynamic and engaging way to understand these documents. Furthermore, Bannihatti Kumar et al. [3] conducted a study leveraging ML algorithms to extract clean texts specifically related to opt-out options. This work further facilitates the development of a browser extension designed to help people better understand their opt-out choices. More recently, Bui et al. [8] presented an automated system, PI-Extract, which uses a neural network model to extract privacy practices. In the same paper, a follow-up user study on investigating the effects of data practice annotations to highlight the extracted privacy practices using PI-Extract was reported to help human readers better comprehend privacy policies.

In addition to the research directions mentioned above on automated privacy policy content analysis, various studies have focused on exploring the inconsistency, lack of clarity, and misalignment

of privacy policy content with the regulatory frameworks. Andow et al. [1] introduced PolicyLint, an automated tool for analysing privacy policies by identifying contradictory sharing and collection practices using advanced NLP and ontology generation techniques. Applying PolicyLint to a large dataset of privacy policies from Google Play, they found a significant occurrence of logical contradictions and narrowing definitions, highlighting the issue of misleading statements. Torre et al. [56] used NLP techniques and supervised ML approaches to develop an AI assistant to classify and check the completeness of privacy policies. A case study of applying this tool to 24 privacy policies discovered 45 out of 47 incomplete issues against the GDPR. Recently, Tang et al. [54] proposed a comprehensive GDPR taxonomy and developed a corpus of labelled privacy policies with hierarchical information. Their extensive efforts in evaluating GDPR concept classifiers aimed to enhance the accuracy and reliability of automated GDPR compliance analysis.

*2.2.2 LLM-based Approaches.* In contrast to traditional ML and NLP-based methodologies, which often require extensive efforts to label data or rely on existing annotated datasets for supervised learning, emerging LLMs have proven to be powerful and efficient for privacy policy analysis, particularly in classifying data practices and extracting valuable insights without requiring extensive manual annotations or training. Various prompt-based methods have been explored to help human readers navigate complex legal texts in privacy policies. Tang et al. [53] developed PolicyGPT, which leveraged LLMs using a "prefix prompt" to perform categorical classification and definition tasks, achieving average Macro F1 scores of 97% and 87% for classification tasks using GPT4 on two privacy policy datasets, respectively. It significantly outperformed past ML models. Similarly, Salvi et al. [46] introduced PrivacyChat, a prompt-engineering-based system using an LLM (GPT-3.5) to improve consumers' comprehension of privacy policies. More recently, Goknil et al. [20] proposed PAPEL, a framework that leverages the capabilities of LLMs through prompt engineering to analyse and convert the critical aspects of privacy policies into user-friendly summaries. Their findings demonstrated that various LLMs, including LLaMA and GPT-3.5/GPT-4, can achieve robust performance in privacy policy analysis tasks. Slightly differently, Rodriguez et al. [44] demonstrated the effectiveness and accuracy of LLM-based solutions for automating privacy policy analysis. They further proposed a set of specific configurations for ChatGPT and argued that such a solution should be considered a replacement for traditional NLP techniques for automated privacy policy processing. More recently, Chen et al. [13] developed LLM-based privacy policy concept classifiers using both prompt engineering and LoRA (low-rank adaptation) fine-tuning techniques, and the results show that such classifiers can outperform other state-of-the-art classifiers using different LLMs (including GPT-3.5, Llama 3 8B, Qwen1.5 7B). To explore LLM explainability in terms of completeness, logicality, and comprehensibility, they also applied prompt engineering to generate explanations for the classification results. Evaluation by three human annotators indicated that LLMs can provide clear and understandable justifications. Similarly, Mori et al. [36] conducted a study to evaluate LLMs' capabilities to understand privacy policies compared with real human users. The results revealed that LLMs could achieve an accuracy of 85.2% of comprehension levels, while the correct answer rate of human users was only 63%, demonstrating the potential of using LLMs to replace user studies for the evaluation of privacy policies. Apart from benefiting ordinary users, LLMs have also been applied to discover legal and reputational risks by detecting policy violations within enterprises' regulatory and operational frameworks. A recent work by [42] proposes a training-free method that consider policy violation detection as an out-of-distribution detection problem and adopts whitening techniques. Their method requires policy text and a small number of examples, yet achieves state-of-the-art performance compared with other LLM-as-a-judge and fine-tuning approaches.

Similar to our work, Xie et al. [63] conducted a study of using the Llama-3.1-70B-Instruct to support automated analysis of privacy policies, which involves two main tasks. One is the assessment to what extent a given privacy policy segment covers the required content of a specific clause. Another one is to asses personal information practices in privacy policies, specifically, focusing on categories of personal information collects and shares, the purposes of data collection and sharing, and the third-party recipients. More recently, and in close relation to our work, Yang et al. [64] developed a framework that performs a series of operations, including privacy policy text segmentation, paragraph-level classification, sentence-level classification, privacy attribute mapping, relationship extraction, and graph generation. Unlike our approach, the key component of their framework for reconstructing attributes and relationship graphs relies on the OPP-115 dataset as training and testing data, applying knowledge distillation techniques involving GPT models and BERT. Their graph generation is implemented using Neo4j, allowing data flows visualisation and interactive query features. However, the reliance solely on OPP-115 labels may not be suitable for modern privacy policies. Different from past studies, our work leverages taxonomies and findings from multiple existing studies to build a knowledge base, aiming to facilitate LLM-driven privacy policy analysis from the perspective of data flows, specifically, what personal data is collected, with whom it is shared, and for what purposes.

*2.2.3   Data Flow Analysis.* Although LLMs have not been utilised to analyse data flows, such analyses has gained significant attention because they can greatly facilitate the communication of complex concepts in a more accessible and engaging manner [58]. An early effort in this domain was PoliCheck, proposed by Andow et al. [2], who developed an entity-sensitive flow-to-policy consistency model to identify contradictions between the data flows described in mobile app privacy policies and the actual data handling practices of the apps. Similarly, Cui et al. [15] introduced PoliGraph, a knowledge graph designed to represent data flows, along with an NLP-based tool, PoliGraph-ER, for the automated extraction of data flows from privacy policy texts. Their evaluation demonstrated that PoliGraph-ER outperformed PoliCheck in identifying inconsistencies between stated and actual data flows.

Yuan et al. [65] also conducted a case study of extracting data flows from the privacy policy of Booking.com, showing significant challenges in manually reconstructing data flows to fully capture the data-sharing landscape. This emphasises the need for more advanced and automated approaches to extracting and analysing data flows in privacy policies. In our prior work [66], we employed GPT-4 as an example LLM to analyse the privacy policies of selected car brands, aiming to produce data entity relationships that capture the types of data shared, the intended purposes of data sharing, and the recipients of the shared data. This was part of a broader effort to reconstruct a vehicle-centric data ecosystem. While the primary focus of that study was not the utilisation of LLMs in privacy policy analysis, it nevertheless offered a glimpse into LLMs' potential and directly motivated the follow-up work presented in this paper.

## 2.3   Identified Research Gaps

Although LLMs have been used for privacy policy analysis, as reviewed in Section 2.2.2, they have not yet been specifically applied to facilitate the extraction of comprehensive data flows and construction of data flow graphs. One possible reason is the well-documented issue of LLM hallucination [24] that can compromise accuracy and produce unreliable results. Hallucinations, particularly those originating from pre-training data, can introduce misinformation, biases, and knowledge gaps [24], making it challenging to extract reliable information about data flows in privacy policy texts. Additionally, privacy policies show significant linguistic variability, with

different organisations using different terminologies and structures to describe similar concepts, making it more complicated to conduct comparison analyses across privacy policies.

A promising approach to mitigate these challenges is retrieval-augmented generation (RAG) [28], which can effectively minimise hallucinations caused by knowledge gaps while preserving the generative capabilities of LLMs. The RAG approach augments LLMs with external knowledge sources by converting both knowledge base elements and users' queries into numerical representations called embeddings. This allows semantic similarity searches to retrieve the most relevant information from the knowledge base, which is then provided to the LLM alongside the original query. As a result, it enhances the accuracy and credibility of the generative outputs of LLMs.

One key advantage of introducing RAG is that LLMs do not need to be retrained for task-specific applications, making them more efficient and adaptable [18]. In addition, considering the diverse linguistic variability shown in different privacy policies, adopting RAG can potentially constrain and unify the generative outputs of LLMs, enabling consistent formatting for more effective and comparative analysis and visualisation.

However, to the best of our knowledge, no past studies have utilised LLMs with RAG to conduct privacy policy analysis for extracting structured data flows in a comprehensive manner. In addition, while prior research has examined privacy policies either manually or with other AI-based methods, a systematic framework for automated processing of privacy policies to extract and visualise data flows remains largely unexplored.

## 3  The Proposed Framework: LADFA



Fig. 1.  LADFA architecture

In this section, we introduce LADFA for extracting data flows from privacy policies, constructing and analysing data flow graphs using LLMs and RAG. As illustrated in Figure 1, it consists of three main components: a pre-processor, an LLM-based processor, and a data flow post-processor.[1] The pre-processor is responsible for 1) defining concepts to facilitate data flow extraction; 2) constructing a knowledge base; and 3) converting the input privacy policy into machine-readable segments. The LLM-based processor employs a hybrid approach that combines prompt chaining with LLM agents.

---

[1]The source code can be accessed from https://github.com/hyyuan/LADFA

As shown in Figure 1, each text segment is processed through multiple sequential sub-components, with each sub-component focusing on addressing a specific task. Dashed lines in the figure indicate the flow of the prompt chain, while the highlighted 'RAG' texts denote instances where the LLM agent needs dynamically access the knowledge base, retrieve relevant information, and augment the generative output. The data flow post-processor consists of a data parser, a graph generator, and an analyser, to collectively transform, structure, visualise, and interpret the extracted data flows from the LLM-based processor. Note that it is possible to leverage LLMs for the pre- and post-processors as well, but this paper focuses on applying LLMs to the analyser because it is the most complicated part of the whole pipeline so can benefit the most from the use of LLMs.

## 3.1 Pre-processor

The pre-processor consists of three sub-components: two design steps (A for preparing key questions and B for constructing relevant domain knowledge bases) and one module (TS for text segmentation). Design steps A and B can be processed once and reused thereafter, and Module TS operates dynamically.

*3.1.1 Design Step A: Preparing Key Definitions and Key Questions.* As reported in previous studies, key parameters in the CI theory that can affect users' privacy expectations and lead to different privacy comprehensions [34]. In this work, we aim to analyse privacy policies and extract comprehensive data flows through the lens of CI theory, adopting its parameters with slight modifications to better align with terminologies that have been used in existing studies [5, 6, 9, 59, 65, 66] and data protection regulations such as the EU/UK GDPR. These refined definitions are used consistently throughout the paper.

One key parameter in the CI theory is 'actors', which consists of 'subject of information', 'the sender of information', and 'recipient of information'. Here, we refer the latter twos as *data sender* and *data receiver*, respectively. The 'subject of information' in this study is the data subject defined in the GDPR, who are usually the consumer/reader of a privacy policy (but not necessarily so if the reader is taking actions on behalf of other data subjects, e.g., parents managing their children's data). Note that the data subject is not always the data sender, e.g., when an online service shares a user's data with a third-party, the data sender is the online service.

As mentioned earlier, another key parameter of CI theory is 'attributes' (i.e., types of information). In the context of this work, attributes are the types of personal information passed from *data sender* to *data receiver*. As reported in [5], personal information can be categorised hierarchically, for instance, 'personal identification information' is the top level category, where 'contact information' is a subcategory of 'personal identification information', and 'phone number' and 'email address' are subcategories of 'contact information'. For consistency and readability, we use *data type* to denote the type of personal information described in privacy policies. In such a hierarchy structure, *data type* may appear as a leaf node (e.g., name or email address), a root node (e.g., personal data), or a mid-level node (e.g., demographic data or location data), depending on the writing style of a privacy policy. In summary, these forms the *data flow* concept discussed in this paper. Specifically *data sender→data type→data receiver* indicates the how a specific type of data flows from one party to another party.

Furthermore, to align variations of *data types* appeared in privacy policies written in different styles and to support consistent comparative analyses, we introduce the concept of *data category*, aiming to further process and classify *data types* based on a customised and simplified typology. Existing studies such as [5, 59] and regulations such as the GDPR have categorised personal information, often through complex, multi-layered hierarchical structures. In this study, we propose a simplified three-level typology, intended solely to facilitate and demonstrate the effectiveness and

usefulness of the proposed framework. The root-level node of the proposed typology is 'personal data', then to differentiate nodes at remaining levels, we refer to a mid-level node as a '*data category*' and a leaf node as a '*data type*'. The term '*data category*' is aligned with the GDPR's notion of 'special categories of personal data'. More details of developing the proposed typology are provided in the paragraph "Knowledge Typology for Classifying Data Categories" later in this section. Moreover, we introduce *data consumer type*, *data processing purpose*, and *data processing method*, aiming to cover the broad scope the 'transmission principles' of the CI theory.

The concept of *data consumer type* is defined to capture the role of data consumer given a data flow. In line with the GDPR's definition, a data consumer maybe a (data) controller or a (data) processor. In this study, if the data consumer is a controller, which typically owns the privacy policy, its type is considered as first-party. If the data consumer is a processor, its type would depend on different contexts: 1) if it operates with the same organisation, it is considered as first-party; 2) if it is an external entity such as third-party service provider, the data consumer type is regarded as third-party. More detailed description can be found in the paragraph "Knowledge Typology for Identifying Data Consumer Type".

The concept of *data processing purpose* is introduced to capture the intended use or objective behind the data processing. While the concept is related to the lawful bases (conditions) defined in GDPR, we did not rely on the regulation to derive the *data processing purpose* applied in this study. Instead, we referred to existing studies [6, 9, 59] that examines privacy policies, from which we established a set of purposes tailored to this study. More details can be found in the paragraph "Knowledge Typology for Identifying Data Processing Purpose".

The *data processing method* refers to the method in which personal data is processed. Inspired by previous studies [59, 65], two main types of data processing methods are introduced in this study: 1) a method is active when the data is voluntarily entered or generated by the user while interacting with the service covered by the privacy policy; 2) a method is passive when the data is automatically collected and shared without user input, so passive from the user's perspective. More detailed discussion can be found in the paragraph "Knowledge Typology for Identifying Data Processing Method".

To this end, *data consumer type*, *data processing purpose*, and *data processing method* further enhance the context of *data flow* to form the *comprehensive data flow*. Building upon the above definitions, we break down the complex task of extracting *comprehensive data flows* into several subtasks, each of which answers a specific question. In this design step, we formulated a set of five questions.

- Q1: What are the claimed *data flows* described in the privacy policy text, in which a *data receiver* collects a *data type* from a *data sender* or *data sender* shares a *data type* with a *data receiver*?
- Q2: For each *data type* within a *data flow* identified in Q1, can it be further processed and classified as a specific *data category*?
- Q3: For each data flow identified in Q1, what is the *data consumer type*?
- Q4: For each data flow identified in Q1, what is the *data processing purpose*?
- Q5: For each data flow identified in Q1, what is the *data processing method*?

To allow LADFA to generate reliable answers to these questions, it is important to define the scope associated with these questions so that LLMs can have better contexts to understand and process the input text. To this end, it is important to establish a domain knowledge base that supports answering these questions and facilitates the construction of customised prompts.

*3.1.2 Design Step B: Constructing The Domain Knowledge Base.* To address questions set in Section 3.1.1, particularly Q2–Q5, we establish a domain knowledge base composed of four knowledge

typologies, each tailored to a specific question. The domain knowledge base is subsequently segmented and encoded into vector representations using an embedding model, then stored in a vector database to support the retrieval phase of the RAG approach. This can also help ground LLMs' behaviours to generate outputs that are more consistent and less variable.

To this end, we adopted definitions and examples from multiple existing studies, including the OPP-115 dataset [59], the personal information taxonomy study [5], a data processing purposes case study [6], a study of data-usage purposes in mobile apps [9], and findings from our own previous studies [65, 66], to build the domain knowledge base. For the rest of the paper, the domain knowledge base is referred to as KB, which consists of $KT_{data}$, $KT_{consumer}$, $KT_{purpose}$, and $KT_{method}$. They represent knowledge typologies for *data category*, *data consumer type*, *data processing purpose*, and *data processing method*, respectively. Table 1 summarises how knowledge typologies were derived using existing studies. In the remainder of this section, we describe how each of these studies contribute to the construction of the knowledge typologies in more details.

Table 1. Summary of how multiple sources were used to construct the four knowledge typologies

| Source | $KT_{data}$ | $KT_{consumer}$ | $KT_{purpose}$ | $KT_{method}$ |
|---|:---:|:---:|:---:|:---:|
| OPP-115 [59] | ✓ | ✓ | ✓ | ✓ |
| Personal information taxonomy [5] | ✓ | | | |
| data processing purposes case study [6] | | | ✓ | |
| Data-usage purposes in mobile apps [9] | | | ✓ | |
| Data flow reconstruction using a privacy policy[65] | | ✓ | ✓ | ✓ |
| Vehicle-centric data ecosystem [66] | ✓ | ✓ | | |

*Knowledge Typology for Classifying Data Categories.* In this part, we introduce how we constructed $KT_{data}$ for processing and classifying data categories to address Q2. As mentioned earlier, we introduce three-level typology, with the root-level represented by a single node, 'personal data'. Our focus is on defining the remaining levels. We use the OPP-115 dataset [59] as a baseline to construct the knowledge typology, and then further refine and expand it using the personal data taxonomies [5] and findings from the work [66].

OPP-115 dataset is a collection of 115 website privacy policies with manual annotations of 23,000 fine-grained data practices [59]. It was released in 2016 and has become a widely used resource for privacy policy research. The dataset defines ten data practice categories, each accompanied by detailed descriptions and illustrative examples. The OPP-115 dataset cover personal data in 16 main categories, including *Finance, Health, Contact, Location, Demographic, Personal Identifier, User Online Activities, User Profile, IP Addresses and Device IDs, Cookies and Tracking Elements, Computer Information, Survey Data, Generic Personal Information, Other,* and *Unspecified.*

The work conducted by Belen Saglam et al. [5] investigates how personal data evolved and was perceived across different domains. They analysed data from multiple sources, including governmental legislation/regulations, privacy policies of applications, and academic research articles, and produced a series of hierarchical personal information taxonomies (see Table 2 for a top-level overview). As shown in Table 2, some data categories identified for different domains in [5] are consistently identified from multiple (i.e., at least three) data sources, including *Demographic, Personal Identification Information, Financial Information, Health Information, Criminal Records/Court Judgements, Sex Life & Sexual Orientation,* and *Communication Data* categories, suggesting these represent common data categories independent of domain context. Among these categories defined in the personal information taxonomies, the first four categories are present in the OPP-115 dataset,

but using slightly different terminologies. For consistency and simplification, *Demographic*, *Personal Identity Identifier*, *Finance*, and *Health* are considered as mid-level nodes of KT$_{data}$.

Table 2. First order system of categorisation of personal information taxonomies from different data sources [5]

| Categories | Government | App (H) | App (F) | Academic paper (H) | Academic paper (F) |
|---|---|---|---|---|---|
| Demographic | ✓ | ✓ | ✓ | ✓ | ✓ |
| Personal Identification Information | ✓ | ✓ | ✓ | ✓ | ✓ |
| Financial Information | ✓ | ✗ | ✓ | ✗ | ✓ |
| Health Information | ✓ | ✓ | ✗ | ✓ | ✗ |
| Judicial Data | ✓ | ✗ | ✗ | ✗ | ✗ |
| Criminal Records/Court Judgements | ✗ | ✓ | ✓ | ✓ | ✓ |
| Sex Life & Sexual Orientation | ✗ | ✓ | ✓ | ✓ | ✓ |
| Technical Device Information | ✗ | ✓ | ✗ | ✓ | ✗ |
| Communication Data | ✗ | ✓ | ✓ | ✓ | ✓ |
| Property/Assets Information | ✗ | ✗ | ✓ | ✗ | ✓ |
| Security Data | ✗ | ✗ | ✗ | ✓ | ✗ |

✓: The data category is identified from the corresponding data source.
✗: The data category is not identified from the corresponding data source.
H: Health, F: Finance

In addition, *Communication Data*, as defined in [5], overlaps with several OPP-115 data categories, including *User Online Activities*, *User Profile*, *IP Addresses and Device IDs*, and *Cookies and Tracking Elements*. To avoid redundancy, *IP Addresses*, *Device IDs*, and *Cookies and Tracking Elements* are merged into one single category, *Online Identifier*. Additionally, *Device Information* was included to replace *Computer Information*, aiming to cover a broader range of modern devices (e.g., smartphones, tablets). Furthermore, because the *Security Data* category in [5] includes biometric attributes, which are absent in the OPP-115 dataset, *Biometric Information* was incorporated as another mid-level node of in KT$_{data}$.

Moreover, our previous study [66] identified governmental bodies as important parts in vehicle-centric data-sharing ecosystems, particularly concerning data collection and sharing involving legal authorities such as courts and law enforcement agencies. This closely aligns with the *Criminal Records/Court Judgements* category listed in Table 2, which is therefore considered as another mid-level node of KT$_{data}$.

Following the same practice of refining and merging results using different sources [5, 59, 66], KT$_{data}$ includes the following data categories (i.e., mid-level nodes: *Demographics*, *Contact*, *Finance*, *Health*, *Location*, *Personal Identity Identifier*, *Online Identifier*, *Device Information*, *Biometric Information*, *User Online Activities*, *User Profile*, *Criminal Records/Court Judgements*, *Generic Personal Information*, *Survey data*, *Other*, and *Unspecified*. Apart from *Other* and *Unspecified*, each data category is accompanied by a text description and a list of data types (i.e., leaf nodes) that serve as illustrative examples belonging to the corresponding data category. It is worth noting that these leaf nodes are mainly derived from [5, 59, 65] and are not exhaustive. An example branch of KT$_{data}$ is shown below:

> **Root node:** Personal data
> **Mid-level node:** Location
> Description: Geo-location information (e.g., a user's current location) regardless of granularity, which may include exact location, ZIP code, or city-level data.
> **Leaf nodes:** Location data, Global Positioning System (GPS) location data, Location history, Global System for Mobile communications (GSM) location data, Universal Mobile Telecommunications Service (UMTS) location data.

It has been well documented that, LLMs can perform better on various tasks with few-shot prompting compared to zero-shot prompting, with few-shot prompting's performance in some cases becoming comparable with fine-tuning approaches [7]. Including such descriptive text alongside an array of *data types* can further enhance LLM performance by providing relevant context and examples to support few-shot prompting. To support this as well as the implementation of RAG, the typology is encoded in JSON format[2]. The same approach is applied to the remaining knowledge typologies presented in the rest of the section.

*Knowledge Typology for Identifying Data Consumer Type.* For the knowledge typology to identify data consumer type and address Q3, we applied the same hierarchical structure as $KT_{data}$ to develop $KT_{consumer}$. The focus is on defining mid-level and leaf nodes. As introduced in Section 3.1.1, the two main data consumer types are first-party and third-party, which are considered as the mid-level nodes. To define them, we adopted the following definitions from the OPP-115 dataset:

- First-party collection and use: '*Privacy practice describing data collection or data use by the company/organisation owning the website or mobile app*', and
- Third-party collection and use: '*Privacy practice describing data sharing with third parties or data collection by third parties. A third-party is a company or organisation other than the first-party company or organisation that owns the website or mobile app.*'

To generalise these definitions beyond websites and mobile apps, we slightly change the text description and introduce non-exhaustive lists of first-party and third-party entities as leaf-nodes of $KT_{consumer}$. These leaf-nodes are informed by prior studies on examining privacy policies [65, 66], which identified various frist-party and third-party entities. As an example, the following branch demonstrates the structure of $KT_{consumer}$:

> **Root node:** Data consumer type
> **Mid-level node:** First Party
> Description: A first party is the entity, such as a website or company, that directly collects and uses personal data from individuals/customers. The company/website/application/service's actual name would be often used as indication of existence of the first party.
> **Leaf nodes:** We, Us, This website, This company, This organisation, Our website, Our company, Our organisation, Our service.

*Knowledge Typology for Identifying Data Processing Purpose.* With respect to Q4, we define $KT_{purpose}$) as a two-level typology, with the root-level node represented by a single node capturing the overarching objective of data processing. Then the focus is on developing an array of leaf nodes with corresponding descriptions that specify distinct purposes, derived from existing studies. Few existing studies have specifically focused on examining data processing purposes in privacy policies. Bhatia and Breaux [6] recruited human annotators to study five privacy policies, producing 218 data purpose annotations. Their analysis identifies six categories of data processing purposes: *Service Purpose*, *Legal Purpose*, *Communication Purpose*, *Protection Purpose*, *Merger Purpose*, and *Vague Purpose*. Bui et al. [9] developed a hierarchical taxonomy of data usage purposes by applying neural text clustering with contextualised word embeddings to group purpose clauses that have similar meanings in a large policy corpus. The hierarchical taxonomy consists of four high-level purposes and 16 low-level purposes (see Table 3). The OPP-115 dataset [59] includes detailed descriptions of 11 purposes, including *Basic Service/Feature*, *Additional Service/Feature*, *Advertising*, *Marketing*, *Analytics/Research*, *Personalisation/Customisation*, *Service Operation and Security*, *Legal*

---

[2]Knowledge typologies encoded in JSON can be accessed from https://osf.io/ab23w/overview?view_only=23e83d260dcf419899585ce868ace61b

*Requirement*, *Merger/Acquisition*, *Other*, and *Unspecified*. These purposes were manually created by three recruited law experts using a top-down approach.

Table 3. The hierarchical taxonomy of data-usage purposes in [8]

| High-level | Low-level |
|---|---|
| Production | Provide service |
| | Improve Service |
| | Personalise Service |
| | Develop Service |
| | Manage Service |
| | Manage Accounts |
| | Process Payments |
| | Security |
| Marketing | Customer Communication |
| | Marketing Analytics |
| | Promotion |
| | Provide Ad |
| | Personalise Ad |
| | General Marketing |
| Legality | General legality |
| Other | Other purposes |

By comparing all the existing sets of purposes, we found that those defined in [6] and [9] are largely reflected in the OPP-115 definitions. However, certain purposes in the OPP-115 dataset, such as *Merger/Acquisition* and *Analytics/Research*, are not explicitly captured in the hierarchical taxonomy proposed in [9]. Moreover, the definitions provided by Bhatia and Breaux [6] are relatively brief and overlap substantially with those in the other two sources. Additionally, as noted in one or our previous studies [65], various personal data are often collected and shared to enable social media integration with websites/apps. We observed that this purpose is absent from the categories in [6, 8, 59].

Considering all the observations and given the broader yet fine-grained set of data processing purposes in OPP-115, we adopted this set as the primary leaf nodes of $KT_{purpose}$ with *Social Media Integration* explicitly added as an additional leaf node in this study. In total, there are 10 leaf nodes of data processing purposes including *Basic Service or Feature*, *Additional Service or Feature*, *Advertising*, *Marketing*, *Analytics or Research*, *Personalisation or Customisation*, *Operational Integrity and Security*, *Legal requirement*, *Merger/Acquisition*, and *Unspecified*. An illustrative branch of $KT_{purpose}$ is presented below:

> **Root node:** Data processing purpose
> **Leaf node:** Advertising
> Description: To show advertisements that are either targeted to the specific user or not targeted, or other general advertising activities.

It is worth noting that the data processing purposes presented here are not intended to cover all lawful bases (conditions) under the GDPR, nor is the aim to produce a comprehensive dataset for this purpose. Nevertheless, comparing existing categorisations/taxonomies with definitions in the GDPR is an interesting topic for future research, however, it is out of the scope of this work.

*Knowledge Typology for Identifying Data Processing Method.* To help address Q5, we mainly use findings from [65] and definitions in the OPP-115 dataset as references to complete the knowledge typology. Similar to $KT_{purpose}$, $KT_{method}$ is structured as a two-level typology, with a single root node and a set of leaf nodes with the corresponding descriptions. A manual privacy policy analysis reported in [65] identified two primary methods of data processing: explicit and implicit methods. This finding aligns with the definitions in the OPP-115 dataset, which defines explicit, implicit, and unspecified data processing modes. In these definitions, explicit refers to cases where users actively provide their data while using the service/website, while implicit refers to data collection and use that occurs passively and automatically without active user involvement. Since the term 'explicit' is often associated with consent in the context of the GDPR, to ensure clarity, consistency, and avoid confusion, we decided to use the terms *Active*, *Passive*, and *Unspecified* as the three leaf nodes $KT_{method}$. To illustrate, one branch of $KT_{method}$ is shown below:

> **Root node:** Data processing method
> **Leaf node:** Active
> Description: a company or organization gathers information that a user knowingly and intentionally provides. This involves instances where users actively input data, such as filling out a web form, creating an account, making a purchase, or subscribing to a newsletter. Explicit data collection typically requires the user to be fully aware of the information being provided and often involves their direct consent.

*3.1.3 Module TS for Text Segmentation.* Several existing tools, such as ASDUS [37] and the `html2text` Python package [47], support automatic text segmentation of HTML documents. However, these tools are either specialised for extracting top-level titles only or struggle with handling table contents and (nested) bullet points. We have observed that many privacy policies frequently use tables and nested bullet points for its content, e.g., for outlining data processing practices. Existing solutions may not effectively convert such HTML content into meaningful text segments for further analysis. Considering this, a customised pipeline utilised the Beautifulsoup4 Python package[3] was developed for automatic text segmentation of an HTML page. It consists of the following main steps: I) removing elements associated with heading, footer, style, and scripts (i.e., `<head>`, `<footer>`, `<style>`, `<script>`); II) processing all non-`<table>` HTML elements following Algorithm 1 to produce an initial list of segments; III) extracting table contents with their headings, where each table is one segment; and IV) for segments that contain bullet points, identifying the associated heading or paragraph, and merging them with the associated bullet points to form one segment (See Figure 5 in Section 6 for an example).

## 3.2 LLM-based Processor

The LLM-based processor adopts a mixed approach that utilises prompt chains and LLM agents. In the context of prompting techniques for generative artificial intelligence (GenAI), an effective prompt engineering technique is to decompose tasks into several subtasks. Prompt chaining refers to sequentially prompting an LLM with a subtask and then using its response as input for the next prompt [61]. Additionally, RAG plays a crucial role in improving the model outputs by incorporating relevant external knowledge. According to the taxonomies reported in a survey [49], RAG systems in the context of GenAI can be classified as agents when considering the retrieval component as an independent tool. As illustrated in Figure 1, the LLM-based processor processes each text segment sequentially through multiple sub-components, with each one responsible for answering one specific question outlined in Section 3.1.1. The dashed lines between the sub-components in

---

[3]https://pypi.org/project/beautifulsoup4/

---

**Algorithm 1** The algorithm for extracting and processing non-`<table>` content of an HTML page

---

```
1: segments ← ∅              ▷ (Initialise an empty list)
2: for all element ∈ FindAllTags(html) do
3:     if element ∈ {p, h1, h2, h3, h4, h5, li, ol, ul} then
4:         if FindParents(element, {ol, ul}) = ∅ then
5:             segments = Process(element, segments)
6:         end if
7:     end if
8: end for
```

```
1: function Process(element, segments)
2:     if FindParent(element) ≠ table then
3:         if element is p then
4:             text ← ExtractTextWithout<a>(element)
5:             if text ≠ ∅ then
6:                 Append text to segments
7:             end if
8:         else if element in {h1, h2, h3, h4, h5} then
9:             text ← ExtractText(element)
10:            if text ≠ ∅ then
11:                Append "*" + text
12:            end if
13:        else if element is li and a not in element then
14:            if element contains no a tags then
15:                if element has nested lists then
16:                    text ← FlattenNestedLists(element)
17:                else
18:                    text ← ExtractText(element)
19:                end if
20:                if text ≠ ∅ then
21:                    Append "−" + text
22:                end if
23:            end if
24:        else if element in {ul, ol} then
25:            for all child in FindChildren(element) do
26:                segments = Process(child, segments)
27:            end for
28:        end if
29:    end if
30:    result ← Join(segments, END_OF_LINE)
31:    return result
32: end function
```

---

Figure 1 represent the use of prompt chain techniques, whereas the highlighted 'RAG' texts indicate that the corresponding sub-component function as LLM agents.

*3.2.1 LLM Screening Agent.* Assume that a single privacy policy $T$ consists of $n$ text segments, i.e., $T = \{T_i\}_{i=1}^{n}$. Not all segments are about personal data collection and data sharing practices. The main objective of this sub-component is to filter out irrelevant text segments and pass only the relevant ones to the next sub-component.

*3.2.2 LLM Data Flow Agent.* If a text segment $T_i$ passes the screening phase, it is fed into this sub-component (i.e., LLM data flow agent) to extract data flows. A customised prompt $P_i^{\text{flow}}$ is formulated, as illustrated in Eq. (1) based on multiple inputs, including the given $T_i$, question Q1 defined in Section 3.1.1, and a set of predefined answering rules $R_{\text{flow}}$. Here, $R_{\text{flow}}$ regulates how an LLM generate outputs and determine their format. Please see Figure 6 in Section 6 for a prompt example.

$$P_i^{\text{flow}} = \text{PromptGeneration}(T_i, Q1, R_{\text{flow}}) \tag{1}$$

Then LLM processes the prompt $P_i^{\text{flow}}$ and generates structured outputs $O_i^{\text{flow}}$ as shown in Eq. (2), where $O_i^{\text{flow}}$ contains $m$ individual data flows, i.e., $O_i^{\text{flow}} = \{F_{i,j}\}_{j=1}^{m}$.

$$O_i^{\text{flow}} = \text{LLM}_{\text{generate}}\left(P_i^{\text{flow}}\right) \tag{2}$$

If there are more than one data sender or data receiver for a identified data flow, multiple data flows will be generated for each unique pair of data sender and data receiver. Each data flow $F_j$ consists of a data sender $DS_j$, a data type $DT_j$, and a data receiver $DR_j$.

$$F_{i,j} = \left(DS_{i,j}, DT_{i,j}, DR_{i,j}\right), \quad \forall j \in \{1, \dots, m\} \tag{3}$$

If either the data sender or receiver is unknown, we instruct the LLM to leave the corresponding field empty. In such cases, we use the placeholder symbol $\perp$ to represent the missing value, resulting in one of two partial tuples $\left(\perp, DT_{i,j}, DR_{i,j}\right)$ or $\left(DS_{i,j}, DT_{i,j}, \perp\right)$.

*3.2.3  LLM-RAG Agents.* To further analyse data flows, we develop a set of LLM-RAG agents to identify: 1) data category, 2) data consumer type, 3) data processing purpose, and 4) data processing method. The retrieval phase of the RAG approach involves identifying the corresponding knowledge typology and then retrieve information that is relevant to a given query. Across the RAG and LLM literature, a variety of terms have been used to describe the retrieved component. Some studies refer to the retrieved information as retrieved contexts [25], others describe the retrieved text documents as additional contexts [28], while others denote the retrieved snippets as contexts [26] or contextual backgrounds [45]. For the clarity and simplicity, here we call the retrieved information knowledge contexts.

Each LLM-RAG agent performs a series of similar tasks that can be generalised as the following:

(1) As part of the RAG process, the task is to retrieve knowledge contexts that are most relevant to a given text segment and a single data flow from the set of extracted data flows, based on semantic similarity. This task can be represented in Eq. (4). Here, $KCX_{i,j}$ represents the top-$k$ retrieved knowledge contexts, where the value of $k$ is dynamically determined. The retrieval is carried out using the `VectorIndexRetriver` module of LlamaIndex [32][4] to embed the query and compares it against the stored knowledge embeddings as discussed in Section 3.1.2. By default, the similarity is computed via cosine similarity, ensuring that the retrieved knowledge contexts most semantically aligned with the query are prioritised. In practice, we use a semantic score threshold of 0.6 to retrieve knowledge contexts, therefore, only contexts with a score greater than 0.6 are selected. If a single knowledge context meets this requirement, it is solely used in the prompt. If two or more contexts satisfy the threshold, the top two are included. If no context exceeds 0.6, the one with the highest score is used. $KT_l$ represents the domain knowledge typology selected for the retrieval, where $l \in \{\text{data}, \text{consumer}, \text{method}, \text{purpose}\}$ corresponds to the knowledge typologies defined in Section 3.1.2. $F_{i,j}$ represents the extracted data flow described in Section 3.2.2, where $^*$ suggests that not all LLM-RAG agents take a completed data flow as input in the retrieving process. One example of using only part of a data flow is detailed in the paragraph "Agent for Identifying Data Category". In addition, $T_{\text{adj}}$ represents optional adjacent text segments (e.g., one example of using $T_{\text{adj}}$ is detailed in the paragraph "Agent for Identifying data processing method" of this sub-subsection).

$$KCX_{i,j} = \text{Retrieve}\left(KT_l, F_{i,j}^*, T_i, T_{\text{adj}}\right) \tag{4}$$

---

[4]LlamaIndex is a data framework specialised for building applications using LLMs

(2) Once the relevant knowledge context(s) is/are retrieved, a customised prompt $P_{i,j}$ is generated, as formulated in Eq. (5). The prompt is constructed using multiple inputs, including the original text segment $T_i$, and the retrieved contexts, $\text{KCX}_{i,j}$. Additionally, one question $Qp$ is selected, where $p \in \{2, \ldots, 5\}$ to guide the prompt formulation. Furthermore, $R$ regulates LLM's output format.

$$P_{i,j} = \text{PromptGeneration}\left(T_i, \text{KCX}_{i,j}, Qp, R\right) \tag{5}$$

(3) The generated prompt $P_{i,j}$ is then passed to the LLM for inference, producing one structured output $O_{i,j}$, as defined in Eq. (6).

$$O_{i,j} = \text{LLM}_{\text{generate}}(P_{i,j}) \tag{6}$$

To distinguish different LLM-RAG agents, we use different superscripts in the notations of Eqs. (4), (5), and (6) in the remaining part of the section.

*Agent for Identifying Data Category.* First, this LLM-RAG agent takes the data type $\text{DT}_{i,j}$ of an identified data flow $F_{i,j}$, a text segment $T_i$, and the knowledge typology of data category $\text{KT}_{\text{data}}$, to retrieve knowledge context(s) $\text{KCX}_{i,j}^{\text{data}}$. The retrieving process follows the same formulation as Eq. (4), with the notations substituted according to those introduced above.

Then, a customised prompt $P_{i,j}^{\text{data}}$ is dynamically augmented using the retrieved knowledge context(s) $\text{KCX}_{i,j}^{\text{data}}$, $R$, along with the data categorisation question Q2. This can be formalised based on Eq. (5) accordingly. See Figure 7 in Section 6 for an example of generated prompt. Finally, the LLM takes the customised prompt $P_{i,j}^{\text{data}}$ and generates one structured output for the data category (i.e., $O_{i,j}^{\text{data}}$).

*Agent for Identifying Data Consumer Type.* The objective of this agent is to determine whether the data consumer type of a given data flow $F_j$ is a first-party, third-party, or undefined. Following the same RAG approach, the context retriever takes $F_{i,j}$, $T_i$, and knowledge typology $\text{KT}_{\text{consumer}}$ to retrieve the relevant knowledge context(s) $\text{KCX}_{i,j}^{\text{consumer}}$. A customised prompt $P_j^{\text{consumer}}$ is then constructed and passed to an LLM to identify a single data consumer type (i.e., $O_{i,j}^{\text{consumer}}$).

*Agent for Identifying Data Processing Purpose.* This agent determines the data processing purpose for a given data flow $F_{i,j}$. Using the same approach, the associated domain knowledge typology $\text{KT}_{\text{purpose}}$ and a given text segment $T_i$ are used to retrieve the knowledge context(s) $\text{KCX}_{i,j}^{\text{purpose}}$, which is then used to facilitate the prompt generation. The output of this LLM-RAG agent is $O_{i,j}^{\text{purpose}}$, which identifies the data processing purpose for the given data flow.

*Agent for Identifying Data Processing Method.* As described in Section 3.1.2, the data processing method can be categorised as active, passive, or unspecified. Unlike other agents described above, determining the data processing method is more challenging, as it often depends on the surrounding textual context. Preliminary experiments indicated that incorporating adjacent text segments, i.e., the previous segment $T_{i-1}$ and the next segment $T_{i+1}$, can help improve classification accuracy. Hence, the retrieval task can be reformulated with minor changes to Equation 4 and represented as follows:

$$\text{KCX}_{i,j}^{\text{method}} = \text{Retrieve}\left(F_{i,j}, T_{i-1}, T_i, T_{i+1}, \text{KT}_{\text{method}}\right) \tag{7}$$

Similar to the above two agents, only the retrieved knowledge context with highest semantic similarity is used. The prompt and output generation processes remained the same as those in the previous sections. The final output of this agent is denoted as $O_{i,j}^{\text{method}}$.

In summary, by processing $n$ privacy policy text segments, the LLM-based processor generates an $n$ sets of outputs $O$. Each set $O_i$ includes $m$ sets of outputs, where each one of these outputs

consists of the identified data flow, the corresponding data category, data consumer type, data processing purpose, and data processing method. These can be represented formally in Eq. (8).

$$O = \left\{ O_i = \left\{ F_{i,j}, O_{i,j}^{\text{data}}, O_{i,j}^{\text{consumer}}, O_{i,j}^{\text{purpose}}, O_{i,j}^{\text{method}} \right\}_{j=1}^{m} \right\}_{i=1}^{n} \tag{8}$$

*3.2.4 Further Implementation Details.* As mentioned earlier, embedding, indexing, and retrieval were developed using LlamaIndex [32]. Specifically, the embedding uses the `BAAI/bge-small-en-v1.5` embedding model [62] to convert textual documents into vector representations that are suitable for retrieval and analysis. The execution of LLM inference is supported by Groq[5], which provides language processing units (LPUs) optimised for accelerating AI inference tasks. The Groq API is compatible with LlamaIndex and supports various open-source LLMs, facilitating seamless integration into LADFA. Different LLMs were tested to find the best ones(s). The screening agent and the data flow extraction agent utilise the `llama-3.3-70b-versatile` model[6], while the agent responsible for identifying data categories uses `llama3-70b-8192` model[7]. For computational efficiency, a lightweight model `llama-3.1-8b-instant`[8] was used for tasks involving smaller knowledge bases, i.e., $\text{KT}_{\text{consumer}}$, $\text{KT}_{\text{purpose}}$, $\text{KT}_{\text{method}}$.

Finally, prompts were construction by following the format required by the Groq Chat Completions API to ensure structured and effective interactions with the selected LLMs. In addition, for all LLMs, we set both `top_p` and `temperature` to be 0.5, aiming to reduce sampling variance and constrain token selection, aiming to produce balanced outputs that favour consistency over creativity.

## 3.3 Data Flow Post-processor

To obtain insights from the outputs generated by the LLM-based processor, this component comprises three key sub-components: a data parser, a graph generator, and an analyser.

*3.3.1 Data Parser.* The data parser is responsible for data cleaning and disambiguation, ensuring that the raw outputs from the LLM-based processor are refined for subsequent graph generation and analysis. Despite providing detailed instructions to LLMs for structuring their output, inconsistencies and ambiguities were observed in the generated data. To address these challenges, additional processing steps are required.

*Duplicate Data Entries.* Redundant entries for the same data (e.g., data sender, data receiver, and data type) often occur due to variations in plural and singular forms or inconsistencies in abbreviations and their corresponding full representations. For instance, 'customers' and 'customer' could be extracted as two distinct data receivers, or 'vehicle identification number' and 'vin' might be recognised as separate data types. To mitigate this issue, rule-based approaches utilising the Python package inflect[9] and Python's re module[10] are applied to standardise entity representations.

*Misclassification of Data Consumer Types.* In general, LLMs may incorrectly classify the data consumer type given a data flow, where the data consumer type that should be classified as a first-party is instead treated as a third-party, or vice versa. For instance, in our case study of OEM

---

[5]https://groq.com/

[6]https://console.groq.com/docs/model/llama-3.3-70b-versatile

[7]https://console.groq.com/docs/model/llama3-70b-8192, this model, with 70B parameters, was the most up-to-date large model available when we conducted the experiment. It is deprecated now in the Groq platform.

[8]https://console.groq.com/docs/model/llama-3.1-8b-instant

[9]https://pypi.org/project/inflect/

[10]https://docs.python.org/3/library/re.html

privacy policies (Section 4), we observed cases where an OEM's mobile app was identified as a data receiver to collect and process personal data, and the corresponding data consumer type was incorrectly categorised as a third party rather than as a first party. To address this issue and support more in-depth analysis, the data parser analyses the data sender and the data receiver collectively and assigns appropriate attributes through a two-step process:

(1) A list of first-party keywords, including 'we', 'us', 'app', and 'website', were generated. Then, the NLP package spaCy[11] is used to extract the root and the possessive modifier (if present) from a given text (i.e., data receiver/sender). If the extracted root text matches any of the predefined first-party keywords or the name of the organisation owning the privacy policy, additionally, if a possessive modifier is present and also matches a keyword from the list or the name of the organisation, the attribute 'first-party' is assigned to the associated data receiver/sender. Otherwise, the attribute 'third-party' is designated.

(2) We introduce the concept of a 'user-party' as another attribute, which applies specifically to data senders. Following a similar NLP approach mentioned above, but with a different set of keywords including 'you', 'user', and 'customer', we assigned 'user-party' attributes to relevant data senders.

By implementing this approach, the data parser enriches the semantic representations of both data senders and receivers. In turn, this enable a more refined classification of data consumer type of a given data flow. As depicted in Table 4, data flows with first-party data consumer type can be represented by three distinct cases, while data flows with third-party data consumer type follow the similar pattern. For the clarity and consistency, we refer to these as first-party data flows and third-party data flows for the rest of this paper.

Table 4. Data flows cases under different data consumer types

**First-party data flow**:
(User-party) data sender → data type → (First-party) data receiver
(First-party) data sender → data type → (First-party) data receiver
(Third-party) data sender → data type → (First-party) data receiver

**Third-party data flow**:
(User-party) data sender → data type → (Third-party) data receiver
(First-party) data sender → data type → (Third-party) data receiver
(Third-party) data sender → data type → (Third-party) data receiver

**Incomplete data flows**:
data sender → data type → ?
? → data type → data receiver
? → data type → ?

Moreover, previous research [58, 65] indicated that privacy policies often fail to clearly specify how data are collected, to what extent, and with whom they are shared. To capture such insights, we instructed the LLM data flow agent, as described in Section 3.2.2, to leave the data sender or data receiver as unknown entities if the text segment does not explicitly specify whose data are processed or who is responsible for the data processing, respectively. We define such cases as incomplete data flows. Incomplete data flows consists of three distinct cases, as illustrated in Table 4,

---

[11]https://spacy.io/

where a question mark represents an unknown entity. These enhancements allow us to conduct a more granular analysis and adds interpretability to the follow-up analysis.

*3.3.2 Graph Generator.* The graph generator is responsible for converting the processed data into a graph-based model, which can be formalised as a directed graph describing how data flows between different entities (i.e., data sender, data type, and data receiver). The graph can be formally described as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{\mathcal{V}_i\}_{i=1}^{M}$ and $\mathcal{E} = \{\mathcal{E}_j\}_{j=1}^{N}$ represent a set of $M$ nodes and a set of $N$ edges, respectively. Each node $\mathcal{V}_i$ represents one entity, and different entity types with different attributes are encoded with different colours. Edges in $\mathcal{G}$ can be visualised in different colours, each representing a different data processing purpose. The implementation of the graph generator is supported by the Python packages networkX[12] and pyvis[13].

*3.3.3 Analyser.* The analyser extracts and reports various statistics to facilitate a deeper understanding of the processed privacy policy data. It mainly reports results of analyses based on 1) the network graph constructed in the Graph Generator, which represents the relationships and data flows between different entities, and 2) outputs of the LLM-based processor. This allows comparisons across different privacy policies, identification of common patterns, and insights discovery. Further details on the specific statistics and their implications are presented in the next section.

## 4   A Case Study for the Automotive Industry

Advancements in AI, the Internet of Things (IoT), and 5G/6G technologies have revolutionised the automotive industry, enabling connected vehicle services and autonomous driving through the integration of electronic control units (ECUs) and sensors that collect, process, and share vast amounts of vehicle and personal data. The extensive scale of data collection and sharing within the modern vehicle ecosystem [29, 66] raises significant privacy and security risks for consumers. A recent review published by Mozilla [10] claimed that "*Cars Are the Worst Product Category We Have Ever Reviewed for Privacy*". This review presents an in-depth data practices analysis for 25 major automotive original equipment manufacturers (OEMs) worldwide, revealing the extensive range of data collected and shared by most automotive OEMs. The work carried out was manually done using privacy policies as the main data sources, and it also highlighted the ambiguity and misleading nature of the complex language used in privacy policies. Inspired by this public review, we decided to use the automotive industry as a case study to test the effectiveness and efficiency of LADFA in automating privacy policy analysis.

   In this section, we first outline the methodology for selecting privacy policies as the dataset in Section 4.1, followed by a description of the evaluation process in Section 4.2. Finally, the experimental results are presented and discussed in Section 4.3.

### 4.1   Datasets

With the advancement of IoT and AI technologies, many modern vehicles offer connected vehicle services through mobile applications. These mobile apps offer various features to enhance driving assistance and user experience. Each mobile app should have a dedicated privacy policy to inform consumers of its data-handling practices. To systematically analyse the privacy policies governing these services, we implemented the following inclusion and exclusion criteria for selecting privacy policies:

---

[12]https://networkx.org/
[13]https://pyvis.readthedocs.io/

(1) The privacy policy must be published for a connected vehicle mobile app available in the EU or the UK. This selection criterion is motivated by the fact that legal frameworks and regulatory requirements governing data protection (e.g., the EU/UK GDPR) are closely aligned, leading to comparable structures and content and allowing systematic cross-policy analysis.

(2) The privacy policy must be specifically and solely dedicated to the associated connected vehicle mobile app. If the identified policy for a connected vehicle mobile app was identical to its OEM's website general privacy policy, it was excluded. Since the general privacy policy focuses on the data handling practices of using the OEM's website rather than the connected vehicle mobile app.

(3) The text of the privacy policy must be embedded within an HTML page rather than being provided as a linked document in formats such as PDF or Word. While LLMs can process other formats, the customised text-segmentation pipeline is optimised for HTML, which is also the dominant format for most of privacy policies.

In addition to the above criteria, automotive OEMs and their associated apps with broken or inaccessible privacy policy links were excluded from the study. Using the Mozilla study [10] as an initial reference for selecting privacy policies, we further refined our dataset by applying the inclusion and exclusion criteria. The selected privacy policies[14] are listed in Table 5, where each privacy policy is saved as a single HTML file.

Table 5. List of OEMs and their associated connected vehicle mobile apps

| OEM | Link | OEM | Link |
| --- | --- | --- | --- |
| Honda | My Honda+ | Lexus | Lexus Link+ |
| Kia | Kia Connect | Nissan | NissanConnect Services |
| Audi | Audi Privacy Policy | Vauxhall | MyVauxhall |
| Hyundai | Hyundai Bluelink Europe | Polestar | Polestar |
| Ford | FordPass | Renault | My Renault |

## 4.2 Evaluation

Due to the lack of ground truth datasets, we used manual validation to evaluate the effectiveness and correctness of LADFA. The first three co-authors of this study participated in the evaluation work as domain experts, where each independently evaluated the outputs generated by the proposed framework based on given text segments of a privacy policy. For each text segment, LADFA generates six outputs corresponding to the data type, data category, data flow, data consumer type, data processing purpose, and data processing method. In addition, the corresponding knowledge bases, as described in Section 3.1.2 were shared with all evaluators as references to facilitate the evaluation. For each output, the evaluator was instructed to give a score using a 1-7 Likert scale (i.e., 1: Strongly disagree, 2: Disagree, 3: Somewhat disagree, 4: Neither agree nor disagree, 5: Somewhat agree, 6: Agree, 7: Strongly agree) to verify its relevance, correctness, and clarity using the given knowledge bases as references.

For the evaluation experiment, we randomly sampled approximately 40% of the total outputs derived from the analysis of the My Honda+ app's privacy policy. This resulted in 150 output tuples, with each tuple containing six tasks, amounting to a total of 900 evaluation tasks per evaluator. The main reasons for choosing My Honda+ are as follows: 1) it produces the highest number of

---

[14]The dataset can be accessed from https://osf.io/zgacu/overview?view_only=ee487642d88f4ce1a14473b8402d4762

extracted data flows (see Section 4.3 for more details); and 2) the structure of the privacy policy text includes diverse formats, including tables, bullet points, and narrative paragraphs, providing a robust and varied evaluation set.

Table 6. Evaluation scores by three evaluators (mean with standard deviation), 7-Likert scale

| Evaluator | Data Type | Data Category | Data Flow | Data Consumer Type | Data Processing Purpose | Data Processing Method |
|---|---|---|---|---|---|---|
| 1 | 6.74 (0.87) | 6.11 (1.35) | 6.27 (1.21) | 6.19 (1.39) | 6.62 (1.07) | 6.78 (0.93) |
| 2 | 6.97 (0.29) | 5.62 (1.95) | 6.52 (1.25) | 6.58 (1.20) | 6.28 (1.40) | 6.58 (1.10) |
| 3 | 6.95 (0.25) | 5.29 (2.06) | 6.02 (1.55) | 6.29 (1.58) | 4.60 (1.93) | 5.92 (1.48) |

Table 6 presents a summary of the evaluation scores across the six tasks assessed by the three evaluators. Overall, the results demonstrate that evaluators generally "agree" or "strongly agree" with the LLMs' outputs, with most average scores between 6 and 7. Specifically, the 'Data Type' and 'Data Flow' scores received consistently high ratings across all evaluators. The mean values ranged from 6.27 to 6.97, with relatively low standard deviations, indicating strong and consistent confidence in LADFA's ability to identify and classify these elements correctly. However, slight variations were observed in the evaluation of the 'Data Category' and 'Data Processing Purpose'. In particular, Evaluator 3 assigned a noticeably lower average score of 4.60 (SD = 1.93) for the purpose category. Similarly, the 'Data category' task received the lower average score, particularly from Evaluators 2 and 3 (5.62 and 5.29, respectively).

The observed disparities could be caused by 1) LADFA's capability to classify data categories and data processing purposes may be less robust or more ambiguous in certain contexts, and 2) the definitions within the relevant knowledge bases used for these classifications may be inherently ambiguous, resulting in different interpretations between human evaluators and the LLMs. These highlight the need to improve LADFA's performance and the need to better define and construct knowledge bases. Nevertheless, the consistently high scores across most evaluation dimensions confirm the proposed framework's overall effectiveness in extracting structured and meaningful information from unstructured privacy policy texts.

Table 7. Gwet's AC1 and percentage agreement using 7-Likert scale and transformed 3-category scale

| Metric | Data Type | Data Category | Data Flow | Data Consumer Type | Data Processing Purpose | Data Processing Method |
|---|---|---|---|---|---|---|
| *7-Likert scale* | | | | | | |
| Gwet's AC1 | 0.85 | 0.37 | 0.51 | 0.56 | 0.21 | 0.56 |
| Percent. Agreement | 0.89 | 0.43 | 0.59 | 0.66 | 0.33 | 0.66 |
| *Transformed 3-category scale* | | | | | | |
| Gwet's AC1 | 0.94 | 0.56 | 0.82 | 0.83 | 0.48 | 0.79 |
| Percent. Agreement | 0.96 | 0.59 | 0.86 | 0.86 | 0.53 | 0.83 |

Moreover, it is essential to examine the inter-rater reliability metric to obtain more insights into the manual verification results. Since we observed that there could be high-agreement situations where traditional metrics such as Cohen's kappa coefficient and Intraclass Correlation Coefficient (ICC) might fail to reflect consensus among raters [17, 21], we used both Gwet's AC1 and the percentage agreements in this study [21, 22]. Gwet's AC1 is a statistical measure of inter-rater reliability, and the percentage agreement measures the raw consensus among evaluators. For both metrics, the value ranges from 0 to 1, where 1 indicates perfect agreement.

Because the results were evaluated using a 7-point Likert scale, it was designed to capture subtle differences in agreements and disagreements. While this granularity is valuable for in-depth assessment, it naturally reduces the likelihood of reaching exact agreements among evaluators.

Nevertheless, the agreement metrics based on the raw scores remain promising. As shown in Table 7, the overall agreement for the 'Data Type' is particularly high, with Gwet's AC1 at 0.85 and a corresponding percentage agreement of 0.889, indicating a strong consensus. However, the agreement was more moderate for other tasks. For 'Data Flow', 'Data Consumer Type', and 'Data Processing Method', Gwet's AC1 scores ranged between 0.5 and 0.56, with percentage agreement above 0.58, suggesting a reasonable level of alignment. The 'Data Category' and 'Data Processing Purpose' showed the lowest agreement, with Gwet's AC1 at 0.37 and 0.21, respectively.

To better measure agreement at a broader level, we also applied a score transformation by collapsing the 7-point Likert scale into a 3-category scale: scores of 1–3 were mapped to 1 (disagreement), 4 to 2 (neutral), and 5–7 to 3 (agreement). This scoring scheme could reduce sensitivity to minor differences in evaluators' responses and produce a more robust measure of general agreement. As shown in Table 7, we observe an improvement in both metrics across all tasks. For instance, 'Data Type' and 'Data Flow' have Gwet's AC1 scores of 0.94 and 0.82, respectively, and corresponding percentage agreements above 0.86, indicating strong agreements on the general validity of LADFA's outputs. The Gwet's AC1 and percentage agreement for 'Data Consumer Type' and 'Data Processing Method' were approximately 0.8, suggesting strong agreements between evaluators. Although the scores for 'Data Category' and 'Data Processing Purpose' improved by approximately 20%, the agreements were relatively low compared with other tasks. The results suggest that, while evaluators may differ on specific levels of agreement, there is a broader consensus on the effectiveness and accuracy of the proposed framework's capability in analysing complex privacy policy documents and extracting comprehensive data flows.

## 4.3 Analysis of Results

We aim to illustrate how interpretation and analysis can be conducted from different perspectives to generate insights of potential privacy concerns that might otherwise be hidden or neglected. For the convenience and clarity, we use OEMs' names instead of mobile app names throughout this section.

Table 8. Summary of data flow network statistics

|  | Audi | Ford | Honda | Hyundai | Kia | Lexus | Nissan | Polestar | Renault | Vauxhall |
|---|---|---|---|---|---|---|---|---|---|---|
| Edges | 84 | 198 | 510 | 79 | 93 | 61 | 201 | 151 | 8 | 176 |
| First-party nodes | 5 | 11 | 6 | 13 | 4 | 2 | 9 | 4 | 2 | 16 |
| Third-party nodes | 22 | 33 | 51 | 19 | 16 | 17 | 42 | 14 | 0 | 23 |
| User-party nodes | 2 | 3 | 3 | 3 | 2 | 1 | 3 | 3 | 1 | 1 |
| Data types nodes | 22 | 73 | 177 | 14 | 32 | 15 | 64 | 64 | 4 | 35 |

*4.3.1 Analysis of Data Flow Graphs.* We represent data flows as graphs and conduct graph-based analysis. Table 8 presents the descriptive statistics of the data flow graph networks extracted from the mobile apps' privacy policies of different OEMs. Honda has the most complex network, which contains 510 edges and more than 200 nodes. In contrast, Ford, Nissan, Polestar, and Vauxhall had moderately smaller networks, with the number of edges ranging from 170 to 270. Audi, Hyundai, Kia, and Lexus have relatively simpler network structures, whereas Renault has the smallest networks. Figure 2 presents a comparison between a complex network (Honda) and a simpler network (Renault) to illustrate the network structures derived from automated privacy policy analysis[15]. In the visualisation, pink, yellow, and green nodes represent third-party, user-party,

---

[15]All generated network graphs can be viewed from https://osf.io/zgacu/overview?view_only=ee487642d88f4ce1a14473b8402d4762

and first-party attributed nodes, respectively. For simplicity, these are referred to as third-party, user-party, and first-party nodes, respectively, in the rest of this paper. The nodes in the light blue boxes represent different types of data. Dark blue is used to annotate nodes with unknown classifications.



Fig. 2. Comparison between (a) a complex data flow network derived from My Honda+ app's privacy policy and (b) a simple data flow network derived from My Renault app's privacy policy

By taking a closer look at specific nodes as illustrated in Figure 3(b), the arrows reflect the directional nature of the network graph, indicating the direction of how data flows between entities. For instance, a third-party node 'your car's display audio' would share an array of data, including 'Android log cat', 'drive history', 'USB device information', 'mobile device ID', 'VIN', 'WiFi station'/footnoteThis is the exact wording used in the privacy policy. We understand that the actual meaning is WiFi station ID., and 'GPS information' to a third-party node 'Panasonic'. Additionally, Figure 3(a) shows a table of the raw text originally appearing in the privacy policy, demonstrating that our method can accurately extract granular information from HTML table content.

Moreover, to systematically compare networks and obtain a more comprehensive understanding, network metrics, including betweenness, closeness, and degree centrality, were computed, and the top ten nodes with the highest scores across all apps are listed in Table 9. To facilitate the analysis and visualisation, we use distinct colours to represent different categories of nodes. It is also worth noting that we chose not to merge or normalise nodes that share the same semantic meaning. For example, the nodes 'We' and 'Us' are semantically equivalent, as are 'Personal data' and 'Your data'. We retained them in their original wording to preserve fidelity to the source text and to minimise the risk of misclassification that may arise during such merging/normalisation process. To retain such raw form also allows us to directly loop up each node in the privacy policy text for cross-validation. Nevertheless, to address potential issues of semantic overlapping, we added

| Processing activity: Why we use your information? | What information is collected? | Lawful basis of processing | Where is the information collected from? | Specific retention periods |
|---|---|---|---|---|
| To provide you with your display audio functionality | Drive history, VIN, Android log cat, Mobile device ID, Wi-Fi station, GPS information, USB device information | The processing is necessary for the ongoing performance, management and facilitation of our contract with you | From your car's display audio | Panasonic will only retain this data for the period of the investigation to identify and resolve the fault. |

(a)



(b)

Fig. 3. Example of converting table content to data flows illustrated in part of a data flow network

additional attributes, including first-party, third-party, and user-party, to collection party nodes, as explained earlier in Section 3.3.1.

*Betweenness Centrality.* Betweenness centrality measures the extent to which a node lies on the shortest paths between other nodes and is often considered a "bridge" or intermediary in the network. As shown in Table 9, the top ten nodes with the highest betweenness centrality scores consistently include first-party nodes (e.g., OEMs names, 'We', 'Us'), user-party nodes (e.g., 'You', 'User') and data nodes for most apps. The only exceptions are Ford and Nissan, where Ford includes third-party service providers such as 'Vodafone' and the 'Roadside assistance provider', while Nissan has 'Our partner' as a third party node in its top ten nodes. This implies the significant role of third-party service providers in Ford's and Nissan's data flow network.

Overall, the 'Personal data'/'Personal information' node had the highest frequency, appearing in 9/10 apps' top ten lists. The 'VIN' node also appeared in more than half of all apps. Both 'personal data' and 'VIN' act as critical bridges across different data flow networks. Moreover, there are some specific sensitive data types that are worth noticing. For instance, the finance data 'Masked/partial card number' for Honda, biometric data 'Voice recording' for Kia, personal identifiable data 'Identity information' for Lexus, and health data 'Medical personal data' for Nissan. The significance of these data in relation to betweenness centrality highlights the scale of sensitive data processing and their bridging roles in associated data flow networks.

*Closeness Centrality.* Closeness centrality measures how close a node is to all other nodes in the network, highlighting its role in the information transmission of a network. By reviewing coloured stacked bars across all apps, Kia appears to have a distinct pattern, as most of its top ten nodes of closeness centrality are third-party nodes, where nodes, such as 'Kia Connected GmbH'

(associated with Germany), 'Cerence B.V.' (associated with the Netherlands), and 'Recipient outside the EU/EEA' indicate the extensive cross-border information flow.

Hyundai and Lexus share similar patterns of distribution of node types. The top ten nodes of the Hyundai network graph includes nodes, such as 'Hyundai Motor Company', 'Hyundai entity located in the Republic of Korea', 'Hyundai AutoEver Europe GmbH', 'Hyundai AutoEver Corp.', and 'Hyundai BluneLink Europe', that are affiliated with Hyundai Group. Similar to Hyundai, Lexus is part of the Toyota corporate Group, and we can see 'Toyota' associated nodes such as 'Toyota Financial Service', 'Toyota Insurance Management', and 'Toyota Financial Service', appear in the list of closeness centrality top ten nodes. These observations suggest that Hyundai and Lexus might have relatively better control over data flows within their networks if we assume that data processing practices could be more transparent within the same corporate group than those involving more third-party entities.

It is worth noting that most of the apps' data flow network graphs are associated with third-party nodes. For instance, 'Processor in third country', 'Tracking service provider', 'Web agency', 'Hosting provider', and 'IT service provider' are heavily involved in the data flow network of Audi. Similarly, 'Analytic', 'Business partner', and 'Subcontractor' are key nodes in Polestar's network. 'Dealer', 'Google', 'Our partner', 'RCI Financial Service Limited trading as Mobilize Financial Service' are essential third-party nodes in the data flow network of Nissan. This highlights the significant dependence on third-party entities to transmit information within these networks. However, the only exception is 'Renault', which does not have a single third-party node appearing in the top ten list. This may indicate either a comparatively minimal reliance on third-party entities for service provision, or that the privacy policy is written in a way that omits or obscures such details.

*Degree Centrality.* Degree centrality reveals centralised hubs/nodes of the network by computing the number of direct connections within a network. To this end, a node with the highest degree of centrality indicates its directional connection to the most nodes, signifying its central role within a network. By inspecting the top ten nodes for all apps, the top three nodes typically included a first-party node (e.g., 'We', 'Us', or the OEM's name), a user-party node (e.g., 'You', 'User', 'Customer'), and a data-type node (e.g., 'Personal data', 'Data'). Different from the observations for betweenness centrality and closeness centrality, the top ten nodes of degree centrality across all apps include all first-party, third-party, user-party, and data type nodes, suggesting a high level of interconnectivity across all node types within the data flow network graph.

It is worth noting that 'Unknown' nodes appear among the top ten nodes for both Polestar and Renault across different metrics. Manual examination of their privacy policies helped explain such nodes. For instance, for Renault's data flows 'Unknown → Cookies → Renault' and 'Unknown → IP address → Renault', the privacy policy states: "*For information relating to personal data that we automatically collect, such as IP address and cookies, a cookie policy is also available on each of Renault's websites or mobile applications.*" While LADFA correctly identified 'Cookies' and 'IP address' as categories of personal data, it did not assign a specific data sender. This is consistent with the policy text, which does not explicitly state the sender of the data. One could reasonably infer that the sender is the user of the app or device, since IP addresses and cookies originate from the user's machine. However, because the input segment of the policy text does not provide explicit attribution, we adopt a conservative stance: LADFA's classification of the sender as 'Unknown' reflects the absence of clear textual evidence. Similar patterns were observed in some of the Polestar results. We acknowledge that with carefully tailored prompt design, LLMs could potentially infer such implicit roles more reliably, although this requires balancing inference with staying true to the text. In addition, we suggest that a more refined text-segmentation strategy that incorporates more cross-paragraph context may further improve the LLM's ability. However, compared with

Table 9. Top 10 betweenness, closeness, and degree of centrality network metrics

| OEM | Metric | Top 10 nodes in Descent order | Distribution of node types |
|---|---|---|---|
| Audi | BC | Personal data, Us, Data, IP address, Information, VIN, Your data, IP address, Information about cyber security vulnerabilities/incident or hacker attack, Contact detail, URL of visited websites | |
| | CC | Log file, Audi, Chatbot, Processor in third country, Personal data, Processor, Tracking service provider, Web agency, Hosting provider, IT service provider | |
| | DC | Personal data, Internet browser, Log file, Data, IP address, User, Audi, Chatbot, Us, Your data | |
| Ford | BC | VIN, User, Data use threshold, Vodafone Global Enterprise Limited (Vodafone), Us,Ford Smart Mobility UK Limited (FSM), Information that you have purchased the subscription, Roadside assistance provider (RSA), We, Vehicle location | |
| | CC | Ford, We, Ford Smart Mobility UK Limited (FSM), Google, Vehicle data, Us, VIN, Our authorized dealer, Company or organisation, Some data | |
| | DC | You, Ford, We, User, Vehicle location, Us, Ford Smart Mobility UK Limited (FSM), Google, Speed, Vehicle information | |
| Honda | BC | Honda, You, Authentication code, Personal information, We, Us, VIN, Honda ID, Data, Masked/partial card number | |
| | CC | Honda, Data, SoundHound, Authentication code, Personal information, Honda ID, VIN, Customer ID, Us, We | |
| | DC | You, Honda, Car (via TCU - Telematic Control Unit), Personal information, Us, We, E3 Media Limited, Worldline, Customer, Manufacturer | |
| Hyundai | BC | Personal data, Your personal data, We, Hyundai Motor Company, Hyundai, Your data, Certain personal data, VIN, Security event-related data, Timestamp of the generated security event | |
| | CC | Data processor, Hyundai Motor Company, Your personal data, Hyundai entity located in the Republic of Korea, Recipient of your personal data, Cerence sub-processor, Hyundai AutoEver Europe GmbH, Hyundai AutoEver Corp., Republic of Korea, Hyundai BlueLink Europe | |
| | DC | Personal data, You, Your personal data, We, Certain personal data, Data processor, Hyundai Motor Company, Hyundai, VIN, Security event-related data | |
| Kia | BC | Personal data, Kia, Your personal data, Log-in data, Location data (GPS), Personal data relating to the contractual relationship, Communication, Commercial letter, Contact detail, Voice recording | |
| | CC | Kia Connected GmbH, Cerence B.V., Autorité de Protection de Donnée, Gegevensbeschermingsautoriteit, Government authority, Court, External advisor, Recipient outside the EU/EEA, Similar third-party that are public body, Service provider | |
| | DC | User, Kia Connected GmbH, Personal data, Kia, Log-in data, Location data (GPS), Cerence B.V., Autorité de Protection de Donnée, Gegevensbeschermingsautoriteit, Voice recording | |
| Lexus | BC | Personal data, Postcode, Address, Personal information, Name, Telephone number, Email address, VIN, Geo location, Identity information | |
| | CC | Us, Toyota, Authorised staff member, Affiliate and subsidiary company, Toyota Financial Service, Toyota Insurance Management, Toyota Insurance Manager, Authorised retailer, Authorised repairer, Personal data | |
| | DC | Personal data, You, Us, Postcode, Address, Toyota Financial Service, Personal information, Authorised staff member, Afiliate and subsidiary company, Member of our authorised retailer and authorised repairer network | |
| Nissan | BC | Personal data, Us, We, Our partner, Your Facebook website browsing data, Anonymised statistical data, Medical personal data, Data provided directly, Data relating to browsing our website, Data relating to use of our mobile application, | |
| | CC | Nissan, Any authorised dealer or repairer, Our financial partner, Us, Personal data, Dealer, We, Nissan Automotive Europe S.A.S, Nissan Motor (GB) Limited, Nissan Automotive S.A.S | |
| | DC | Personal data, You, Unknown, Nissan, Any authorised dealer or repairer, Our financing partner, Us, Our partner, Nissan Automotive Europe S.A.S, We, Dealer | |
| Polestar | BC | Polestar App, We, Personal data, Phone number, Email address, Your information, VIN, Location, Contact information, Home address | |
| | CC | Analytic, We, Customer support, Vehicl control, Personal data, Your information, Polestar app, Polestar, Business partner, Subcontractor | |
| | DC | Unkonwn, Analytic, Polestar app, We, Customer, Personal data,Customer support, Vehicl control, Phone number, Email address, | |
| Renault | BC | Personal data, Information that makes it possible to identify you, IP address, Cookie, You, Renault Group, Renault | |
| | CC | Renault Group, Renault, Personal data, Information that makes it possible to identify you, IP address, Cookie, You, Unknown | |
| | DC | You, Personal data, Renault Group, Information that makes it possible to identify you, Unknown, IP address, Renault, Cookie | |
| Vauxhall | BC | Data processor, Personal data, We, Contact detail, Our network, Your data, Aggregated information, Vehicle data, Data inferred by our activity, Data collected by the browser | |
| | CC | Us, Stellantis Europe, Partner, We, Car manufacturer, Third selected partner, Social media platform, Programmatic advertising platform, Contact detail, Our website and application | |
| | DC | Data, You, Us, Personal data, We, Contact detail, Partner, Your device, Stellantis Europe, Social media platform | |

BC: Betweenness Centrality, CC: Closeness Centrality, DC: Degree Centrality

Colour scheme: Data type nodes, first-party nodes, third-party nodes, User-party nodes, Unknown nodes

feeding shorter text segments to LLMs, we also noticed that supplying overly lengthy text is not ideal for extracting fine-grained data flows, as it can incorrectly associate unrelated data senders/data receivers with data types. Further work is therefore needed to explore more dynamic and context-aware segmentation strategies that balance contextual completeness with the precision required for detailed data-flow extraction.

Moreover, as identified in Table 8, Renault has only three and four data type nodes, respectively. This is abnormal compared to other apps. We manually examined its privacy policy. For Renault, it applies a generic approach to describe its data handling practices throughout the privacy policy, where only generic terms such as 'personal data' or 'information' are used. This raises concerns about the transparency of its privacy policies, suggesting that the privacy policies may be ambiguous or lack clarity in explicitly explaining data-handling practices.

*Summary.* Overall, graph analysis based on the centrality of betweenness, closeness, and degrees can enhance our understanding of data collection and sharing practices stated in privacy policies. Here, we summarise the insights obtained from such network analyses as follows:

- **Bridging role of data types in the data flow network**: 'Personal data' and 'VIN' (appearing in more than half of apps), and specific sensitive data (e.g., Honda's 'masked card numbers', Kia's 'voice recordings', Nissan's 'medical data') exhibit high centrality of betweenness, highlighting their roles in facilitating connections among entities as well as amplifying the potential privacy risks if compromised.
- **Third-party service providers' dependencies**: Audi, Ford, Nissan, Kia, Polestar, and Vauxhall demonstrate heavy reliance on third-party partners/service providers based on different network metrics, highlighting the complex accountability chains and potential attack surfaces.
- **Cross-border data transmission**: Hyundai, Kia, and Lexus show cross-border data flows involving international entities. This raises privacy concerns when there are mismatches in data protection standards and legislation among different regions.
- **Corporate-group centralised control**: Hyundai and Lexus display tighter intra-organisational control, with 60-75% of their top-closeness nodes being affiliated entities. Different from the heavy reliance on third-party entities observed in other apps, this could potentially enable more standardised data governance.
- **Transparency gaps**: Renault has 'Unknown' nodes in different network metrics' top ten lists. With its generic approach of describing data flows, it reflects the transparency of its associated privacy policy.

*4.3.2  Data Flow and Data Consumer Type Analysis.* As shown in Table 10, among all automotive OEMs, Honda reported the highest number of data flows (as illustrated in Figure 2 (a)), whereas Renault reported only four (as illustrated in Figure 2 (b)). Assessing privacy practices solely based on the number of identified data flows is challenging because privacy policies vary significantly in writing style and level of detail. In addition, how well a privacy policy aligns with its actual data processing practices in reality is unknown. We acknowledge the limitations of using the privacy policy as a single source for conducting the analysis alone. However, if we assume that each privacy policy maintains a consistent style and structure throughout, analysing the normalised frequencies (i.e., relative proportions) of first-party, third-party, and incomplete data flows (see Table 4 in Section 3.3.1) across different privacy policies can still provide meaningful insights. For each privacy policy, we calculated the fraction of each data flow type relative to the total number of data flows (values between 0 and 1). The results are depicted in Table 10.

Table 10. Data flow statistics across selected automotive OEMs

| | Audi | Ford | Honda | Hyundai | Kia | Lexus | Nissan | Polestar | Renault | Vauxhall |
|---|---|---|---|---|---|---|---|---|---|---|
| **Number of data flows** | 56 | 109 | 348 | 65 | 63 | 45 | 141 | 80 | 4 | 142 |
| **First-party data flow** | | | | | | | | | | |
| User-party to first-party | 0.17 | 0.19 | 0.38 | 0.20 | 0.37 | 0.22 | 0.33 | 0.025 | 0.50 | 0.31 |
| First-party to first-party | 0.04 | 0.11 | 0.01 | 0.25 | 0.02 | 0.00 | 0.01 | 0.15 | 0.00 | 0.10 |
| Third-party to first-party | 0.00 | 0.11 | 0.13 | 0.00 | 0.00 | 0.00 | 0.11 | 0.1 | 0.00 | 0.10 |
| **Total** | **0.21** | **0.41** | **0.52** | **0.45** | **0.38** | **0.22** | **0.45** | **0.28** | **0.50** | **0.54** |
| **Third-party data flow** | | | | | | | | | | |
| User-party to third-party | 0.11 | 0.36 | 0.30 | 0.14 | 0.08 | 0.76 | 0.26 | 0.15 | 0.00 | 0.41 |
| First-party to third-party | 0.36 | 0.08 | 0.05 | 0.35 | 0.54 | 0.00 | 0.24 | 0.125 | 0.00 | 0.01 |
| Third-party to third-party | 0.20 | 0.06 | 0.01 | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 | 0.00 | 0.03 |
| **Total** | **0.66** | **0.51** | **0.36** | **0.49** | **0.62** | **0.76** | **0.51** | **0.28** | **0.00** | **0.45** |
| **Incomplete data flow** | **0.13** | **0.08** | **0.12** | **0.06** | **0.00** | **0.02** | **0.04** | **0.45** | **0.50** | **0.01** |

Values under "Number of data flows" are absolute counts; all other numerical values represent normalised frequencies of each data flow type relative to the total data flows within that OEM's privacy policy.

Across all data flows, each app follows a similar pattern, where the majority of first-party data flows is from user-party nodes to first-party nodes, whereas only a small fraction is from third-party nodes to first-party nodes. Overall, more than half of the data flows in Honda, Renault, and Vauxhall had a normalised frequency greater than 0.50 of their data flows classified as first-party. In contrast, Ford, Hyundai, Kia, and Nissan had values around 0.40, while Audi, Lexus, and Polestar showed even lower proportions of first-party data flows.

No clear pattern emerged for third-party data flows. Some apps (i.e., Ford, Honda, Lexus, Nissan, Vauxhall, Polestar) have a higher proportion of data flowing from the user-party to the third-party, while others (i.e., Audi, Hyundai, and Kia) have more flows from the first-party to the third-party. Overall, Audi, Ford, Kia, Lexus, and Nissan had more than half (normalised frequency <0.50) of their data flows classified as third-party, while the remaining apps fell between 0.20 and 0.50, except for Renault. Moreover, regarding the percentage of incomplete data flow, Polestar and Renault have a significantly higher proportion than the remaining apps, indicating potential unclear and ambiguity in their privacy policies.

To illustrate how such data can be used to gain more comprehensive insights and compare different privacy policies, we propose a generalised data flow risk score using max-normalisation of weights, aiming to provide indicative figures that can be used to compare different privacy policies in terms of potential risks and transparency. The score is computed by summing the weighted contributions of different data flow types and normalising the score to the range [0, 1]. As illustrated in Eq. (9), the numerator of the equation represents the calculation of a raw privacy score for a given instance $i$, where $w_j$ is the weight for each data flow type and $X_{j,i}$ is the corresponding value for instance $i$. The denominator of the equation calculates the maximum possible score using the maximum observed values of each data flow type, where $\max(X_j)$ is the highest value observed for the $j$-th data flow type across all instances. Finally, the normalised data flow risk score was computed by dividing the raw score by the maximum possible score:

$$S_{\mathrm{norm},i} = \frac{\sum_{j=1}^{n} w_j \cdot X_{j,i}}{\sum_{j=1}^{n} w_j \cdot \max(X_j)} \tag{9}$$

It is worth noting that the selection of different weights below is based on our own experience and judgment. They are intended for illustrative purposes, and the subsequent analysis depends on these chosen weights. For first-party data flows, cases where data originate from third-party nodes and terminate at first-party nodes are assumed to pose more risks, as the data handling practices of

third-party entities lie outside the control of both users and first parties. To reflect this, we apply weighted scores of 1 for user-party to first-party flows, 1.5 for first-party to first-party flows, and 2.25 for third-party to first-party flows. For third-party data flows, we apply weights based on the following assumption: 1 for user-party to third-party flows, 1.5 for first-party to third-party flows, and 2.25 for third-party to third-party flows. For the overall privacy scores, we assume that incomplete data flows pose a greater risk than both third-party and first-party data flows; hence, we assign a weight of 1 to overall scores of first-party data flows, 1.5 to third-party data flows, and 2.25 to incomplete data flows.

Table 11. Summary of data flow risk scores

|  | APP1 | APP2 | APP3 | APP4 | APP5 | APP6 | APP7 | APP8 | APP9 | APP10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Score for first-party data flows | 0.20 | 0.52 | 0.59 | 0.49 | 0.33 | 0.19 | 0.51 | 0.41 | 0.43 | 0.50 |
| Score for third-party data flows | 0.54 | 0.31 | 0.24 | 0.33 | 0.44 | 0.38 | 0.32 | 0.17 | 0.00 | 0.25 |
| Overall score | 0.54 | 0.49 | 0.47 | 0.48 | 0.47 | 0.51 | 0.47 | 0.61 | 0.58 | 0.49 |

It is important to acknowledge that privacy is a complex and multifaceted concept that cannot be fully captured using quantitative or qualitative rankings alone. Such measures can be misleading, or even harmful, if not considered alongside broader legal, organisational, and contextual factors. That is why we chose to anonymise apps' names in this part of the analysis, as it is mainly intended to illustrate methodological insights rather than to systematically and comprehensively rank privacy risks.

As reported in Table 11, a higher score represents higher data flow-related risks. APP2, APP3, APP4, APP7, APP9, and APP10 have relatively higher risk scores for first-party data flow (i.e., between 0.43 and 0.59) than their scores for third-party data flows. In contrast, APP1, APP5, APP6, and APP8 have higher risk scores for third-party data flows, with APP1 having the highest risk (i.e., 0.54), followed closely by APP5 and APP6 with scores of 0.44 and 0.38, respectively. It is worth noting that APP9 has a score of 0.00, indicating no third-party data flow risks, which is due to the absence of such data flows. This is shown in Figure 2, suggesting that the privacy policy is overly broad and lacks the granularity needed to better inform users about their general data handling practices.

In addition, with the contribution of incomplete data flows to the overall risk scores, APP1, APP6, APP8, and APP9 had the highest scores, approximately between 0.51 and 0.61. This aligns with their higher percentage of incomplete data flows, indicating ambiguity and a lack of transparency in privacy policies. This is also in line with the observations obtained from the graph analysis presented in Section 4.3.1, where they were singled out for potential privacy concerns related to their dependency on third-party entities and privacy policy transparency issues. The overall risk scores for the remaining apps fell just below 0.50, indicating that these apps have relatively better or clearer data handling practices stated in their privacy policies from the data flow perspective.

*4.3.3 Data Category and Data Processing Purpose Analysis.* One of the main tasks for the LLM-based processor is to further classify data types into data categories based on the definitions set in knowledge typology $KT_{data}$ to allow more in-depth analysis. As depicted in Table 12, Honda stands out with the highest number of data types across almost all data categories. Honda collects more data types, particularly in categories such as 'Online identifiers' (25), 'User online activities' (18), 'Location' (13), and 'Finance' (17), indicating its broad data processing related to online user tracking and financial transactions. However, Audi, Kia, Lexus, Polestar, Renault, and Vauxhall collected significantly less data across most data-type categories. In addition, the 'Other' data type category contains data types that cannot be categorised using the other categories listed in the table.

By manually inspecting this data category, it mostly contains vehicle-related data (e.g., 'speed', 'vehicle data', 'braking' and 'DTC (Diagnostic Trouble Code) history') and other unspecified data (e.g., 'reason for your Roadside Assistance call', 'timestamp of the generated security event').

Table 12. Data category distribution across selected automotive companies

| Data Category | Audi | Ford | Honda | Hyundai | Kia | Lexus | Nissan | Polestar | Renault | Vauxhall |
|---|---|---|---|---|---|---|---|---|---|---|
| Generic Personal Info | 1 | 5 | 7 | 4 | 3 | 3 | 9 | 5 | 1 | 2 |
| Personal ID Identifier | 1 | 3 | 4 | 0 | 2 | 1 | 2 | 1 | 0 | 0 |
| Online Identifier | 2 | 3 | 25 | 0 | 3 | 1 | 4 | 7 | 2 | 6 |
| User Online Activities | 9 | 5 | 18 | 0 | 2 | 0 | 12 | 7 | 0 | 6 |
| User Profile | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |
| Contact | 1 | 7 | 8 | 2 | 2 | 4 | 4 | 5 | 0 | 5 |
| Demographic | 0 | 3 | 13 | 2 | 0 | 1 | 4 | 5 | 0 | 2 |
| Biometric Information | 0 | 0 | 1 | 0 | 2 | 1 | 1 | 0 | 1 | 0 |
| Location | 1 | 15 | 13 | 0 | 10 | 1 | 1 | 8 | 0 | 1 |
| Finance | 0 | 2 | 17 | 0 | 0 | 0 | 12 | 5 | 0 | 0 |
| Health | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 3 | 0 | 0 |
| Criminal Records | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Device Information | 2 | 6 | 11 | 1 | 2 | 0 | 1 | 9 | 0 | 5 |
| Other | 7 | 9 | 60 | 5 | 6 | 3 | 17 | 10 | 0 | 7 |

Certain data categories, such as 'Profile', 'Biometric Information', 'Finance', 'Health', and 'Criminal Records', exhibit clear variations between apps. 'Biometric Information' is collected by multiple apps, including Honda (i.e., voice command), Kia (i.e., voice recording, voice samples), Lexus (i.e., sound or images files), Nissan (i.e., voice recordings from calls between you and the dealer), and Renault (i.e., information that makes it possible to identify you). A broad range of 'Finance' data is collected by Honda (e.g., tax code, payment card data, PAN, cardholder name, transaction information) and Nissan (e.g., purchasing history, payment method, discount granted, order history, order number). 'Health' data processing for Nissan is medical personal data. Meanwhile, Ford is the only company that explicitly claims a collection of 'Criminal Records', including crime reference numbers and vehicle theft information. Regarding 'User Profile' data, Ford collects profile pictures, Honda records information about the primary user, Nissan collects usernames and profiling details, and Vauxhall claims the processing of general profile data. Overall, the differences in data processing across apps presented in Table 12 indicate varying levels of specificity and transparency in privacy policies, with some apps explicitly detailing specific data types, while others remain more general or omit such disclosures.

Furthermore, we were interested in how different data categories are associated with different data processing purposes. To achieve this, we hand-picked data categories that were claimed to be collected and shared by two-thirds of the apps and examined the distribution of data processing purposes across different data categories collectively. Overall, as shown in Figure 4, almost all data categories are linked to multiple purposes of data processing to some extent. This highlights the multifunctional role that personal data plays in the automotive context, serving both operational needs and broader objectives.

Among all data categories, 'Operational Integrity and Security' is the most frequently associated data processing purpose across most data categories. This is partly due to its general and broad definition, making it applicable to a wide range of data flows. In addition, data processing purpose 'Legal Requirement' are consistently prominent across most of data categories. Combined with 'Operational Integrity and Security', this might mean that legal and compliance-based justifications underpin much of the identified data categories.

Fig. 4. Distribution of data processing purposes across different data categories.

Moreover, multiple data categories such as 'Contacts', 'Generic Personal Information', 'Location', 'Online Identifier', 'User Online Activities' and 'Device Information' stand out with a diverse spread, where they are frequently associated with a wide range of data processing purposes, indicating their multiple roles in the data processing. For instance, 'Contacts', its association with multiple data processing purposes suggests its dual role in both functional (i.e., basic service, communication) and business-oriented uses (i.e., marketing and promotion). 'Location', 'Device Information', 'Online Identifier', and 'User Online Activities' show the strong link to the purpose 'Analytics or Research', reflecting their additional role in possible behavioural profiling and service optimisation.

The processing of 'Personal Identity Identifier' category is commonly linked with purposes including 'Basic Service or Feature', 'Legal Requirement', 'Marketing', and 'Merger/Acquisition', in addition to 'Operational Integrity and Security', indicating the broad scale of sensitive data processing, and raising potential privacy concerns about oversharing personal sensitive data.

In general, mapping data processing purposes to data categories provides insight into the collective behaviour of data-handling practices within the automotive industry. This further confirms and raises potential concerns about necessary data minimisation, excessive data retention, and potential secondary use. Moreover, it would be interesting to take privacy policies from other domains and conduct cross-sector comparison studies; however, this is out of the scope of this paper and could be an interesting topic as part of our future work.

*4.3.4 Cross-Check with Additional Data Sources.* Although privacy policies are subject to legal requirements to disclose data-handling practices, the extent to which they provide accurate information is unclear. Hence, we need other sources of information to cross-check the results. To this end, we visited the Google Play (i.e., Android) apps' data safety sections[16] and the Apple iOS apps' privacy labels[17], where both require app developers to disclose information about mobile apps' data collection and data sharing practices.

As shown in Tables 13 and 14 in Section 6, across all selected mobile apps, disclosures via both the Google Play apps' data safety sections and Apple iOS apps' privacy labels contradict the data

---

[16]https://support.google.com/googleplay/answer/11416267

[17]https://developer.apple.com/app-store/app-privacy-details/

practices stated in their associated privacy policies. For instance, the NissanConnect Services app collects financial data as stated in its privacy policy, but it fails to declare this in its Google Play app's data safety section and Apple iOS apps' privacy labels. Moreover, even for the same mobile app, there are noticeable inconsistencies in data collection and sharing disclosures between Google Play's app safety section and Apple iOS apps' privacy labels. For instance, Audi's app does not collect any data according to its Apple iOS app's privacy labels, whereas several types of personal information are collected based on its Google Play app's safety description. This may suggest that 1) mobile developers may be negligent in accurately reporting their data-handling practices; 2) there is a lack of a clear understanding of the regulatory landscape; and 3) there is a significant discrepancy in data security and privacy between Google Play apps and Apple iOS apps. These observations highlight the difficulty in verifying the accuracy of privacy policies, while also revealing the unreliability of Google Play apps' data safety sections and Apple iOS apps' privacy labels. The significant discrepancies between different data sources raise concerns about the extent to which consumers and, perhaps, even developers can truly understand the full scope and granularity of data-handling practices. However, systematically addressing these is beyond the scope of this study, but it is worth investigating in future studies.

## 5 Limitations and Future Work

One limitation of this work is the verification process. The lack of an existing ground-truth dataset specifically designed for extracting data flows from text prevents us from conducting a large-scale evaluation. This study relied only on manual verification involving human verifiers to validate the extracted information. While this helps confirm and ensure the accuracy of our methods, it inherently limits our ability to benchmark the results with other methods. To address this limitation, *we plan to utilise LLMs teaming with human users to co-curate dedicated ground truth datasets as part of our future work*. In addition, we acknowledge that relying on three co-authors as evaluators can introduce potential biases, which could potentially affect the generalisability and robustness of the validation results. To address this limitation, *we are planning to conduct a more independent validation study with recruited human participants as part of our future work.*

It is also important to acknowledge that LLMs could be employed not only within the analyser but also for the pre-processor and data flow post-processor of the proposed framework. For instance, LLMs could be utilised to assist in refining input format, segmenting text, or enhancing insights discovery. However, in this study we restricted the use of LLMs to the analyser component, as the analyser represents the most complex and demanding stage of the pipeline, where the benefits of integrating LLMs are expected to be most promising. While this design choice allowed us to focus on demonstrating the feasibility and effectiveness of LADFA, it also constitutes a limitation. Future research could explore and evaluate the use of LLMs to support to other components of the framework.

Moreover, we consider that the construction of the knowledge bases used in the RAG implementation is another limitation of this work, since the accuracy of these knowledge bases can affect the results generated by the LLMs. Our findings demonstrate that vague or unclear definitions can lead to ambiguous and incorrect outputs from LLMs. For instance, concepts such as active and passive/automatic data processing are inherently difficult to define with great clarity. Even for human readers, the interpretations of the same text by different people would be different. To this end, we would recommend *that knowledge bases should be co-created with domain experts to enhance clarity and improve the reliability of LLMs' outputs while also helping to reduce the risk of hallucinations.*

As introduced in Section 3.1.3, LLMs process text segments sequentially to complete a set of tasks. This approach may overlook cross-paragraph (i.e., cross-segment) information, potentially

leading to unreliable or incomplete outputs from LLMs. However, we found that for a granular task such as extracting data flows, feeding the entire privacy policy as the input would lead to 1) the loss of granularity compared to our current approach of sequentially feeding smaller text segments and 2) incorrect association of unrelated data senders/data receivers with data types. Nevertheless, we would like to emphasise that we also made some efforts to mitigate this limitation by using paragraph-based text segmentation, adding headings to bullet point-based segments, and adding a table header row to each table row for additional contextual information. In a nutshell, we would *recommend that further research is needed to 1) determine the optimal text segmentation approach with considerations of cross-paragraph contexts to improve LLMs' performance on such fine-grained tasks and 2) explore other approaches beyond RAG such as contextual-retrieval*[18] *and knowledge augmented generation (KAG)* [30].

Furthermore, our current network graphs of data flows reveal strong connections between corporate groups and third parties. It would be useful to allow the propose framework to integrate with existing tools (e.g., NetVizCorpy [4], a tool that reconstructs business-to-business relationship graphs using Wikidata) to enrich and expand the scope of our network analysis by providing a more comprehensive view of inter-organisational data flows extracted from other data sources, as well as helping validating the outputs generated by LADFA.

Last but not least, while this work focuses on analysing privacy policies for the automotive industry, the aim is to demonstrate its capability and effectiveness of conducting complicated tasks using LLMs with RAG. We would also like to *highlight that LADFA is not only capable of analysing privacy policies but also generalisable and adaptable for examining any text-based documents for other domains.* Additionally, many of the parts, such as the text segmentation tool, domain knowledge bases, and LLM-based agents in the proposed framework, *can be updated and customised for different tasks to provide better flexibility and generalisability.* In our future work, we will explore applications of LADFA to different types of documents and different tasks.

## 6  Conclusion

This paper reports our work on developing and evaluating an end-to-end framework, LADFA, for automating privacy policy analysis using LLMs and RAG. LADFA consists of a pre-processor, an LLM-based processor, and a data flow post-processor. It can 1) preprocess privacy policy texts by converting and segmenting HTML-based text and constructing local knowledge bases; 2) utilise several LLM agents to analyse privacy policy text segments and extract comprehensive data flows; and 3) post-process LLMs' outputs and utilise data flow graphs to discover data protection and privacy-related insights. To demonstrate the usefulness and effectiveness of the proposed framework, we conducted a case study involving the analysis of privacy policies from ten selected connected vehicle mobile apps. The results show that LADFA can effectively and accurately understand privacy policies and extract key information to generate comprehensive data flows, as well as reveal privacy insights that may require considerable time for consumers to read and comprehend.

## References

[1] Benjamin Andow, Samin Yaseer Mahmud, Wenyu Wang, Justin Whitaker, William Enck, Bradley Reaves, Kapil Singh, and Tao Xie. 2019. PolicyLint: Investigating Internal Privacy Policy Contradictions on Google Play. In *Proceedings of the 28th USENIX Security Symposium*. USENIX Association, 585–602. https://www.usenix.org/conference/usenixsecurity19/presentation/andow

[2] Benjamin Andow, Samin Yaseer Mahmud, Justin Whitaker, William Enck, Bradley Reaves, Kapil Singh, and Serge Egelman. 2020. Actions Speak Louder than Words: Entity-Sensitive Privacy Policy and Data Flow Analysis with

---

[18]https://www.anthropic.com/news/contextual-retrieval

PoliCheck. In *Proceedings of the 29th USENIX Security Symposium*. USENIX Association, 985–1002. https://www.usenix.org/conference/usenixsecurity20/presentation/andow

[3] Vinayshekhar Bannihatti Kumar, Roger Iyengar, Namita Nisal, Yuanyuan Feng, Hana Habib, Peter Story, Sushain Cherivirala, Margaret Hagan, Lorrie Cranor, Shomir Wilson, Florian Schaub, and Norman Sadeh. 2020. Finding a Choice in a Haystack: Automatic Extraction of Opt-Out Statements from Privacy Policy Text. In *Proceedings of The Web Conference 2020*. ACM, 1943–1954. doi:10.1145/3366423.3380262

[4] Zsofia Baruwa, Haiyue Yuan, Shujun Li, and Zhen Zhu. 2025. Constructing and Analysing Global Corporate Networks With Wikidata: The Case of Electric Vehicle Industry. *Global Networks* 25, 4, Article e70029 (2025), 17 pages. doi:10.1111/glob.70029

[5] Rahime Belen Saglam, Jason R. C. Nurse, and Duncan Hodges. 2022. Personal information: Perceptions, types and evolution. *Journal of Information Security and Applications* 66, Article 103163 (2022), 31 pages. doi:10.1016/j.jisa.2022.103163

[6] Jaspreet Bhatia and Travis D. Breaux. 2017. A Data Purpose Case Study of Privacy Policies. In *Proceedings of the 2017 IEEE 25th International Requirements Engineering Conference*. IEEE, 394–399. doi:10.1109/RE.2017.56

[7] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Article 159, 25 pages. https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf

[8] Duc Bui, Kang G Shin, Jong-Min Choi, and Junbum Shin. 2021. Automated Extraction and Presentation of Data Practices in Privacy Policies. *Proceedings on Privacy Enhancing Technologies* 2021, 2 (2021), 88–110. doi:10.2478/popets-2021-0019

[9] Duc Bui, Yuan Yao, Kang G. Shin, Jong-Min Choi, and Junbum Shin. 2021. Consistency Analysis of Data-Usage Purposes in Mobile Apps. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2824—-2843. doi:10.1145/3460120.3484536

[10] Jen Caltrider, Misha Rykov, and Zoë MacDonald. 2023. *It's Official: Cars Are the Worst Product Category We Have Ever Reviewed for Privacy*. https://www.mozillafoundation.org/en/privacynotincluded/articles/its-official-cars-are-the-worst-product-category-we-have-ever-reviewed-for-privacy/ Accessed: 2025-07-27.

[11] Fred H. Cate. 2010. The Limits of Notice and Choice. *IEEE Security & Privacy* 8, 2 (2010), 59–62. doi:10.1109/MSP.2010.84

[12] Baiqi Chen, Tingmin Wu, Yanjun Zhang, Mohan Baruwal Chhetri, and Guangdong Bai. 2023. Investigating Users' Understanding of Privacy Policies of Virtual Personal Assistant Applications. In *Proceedings of the 2023 ACM Asia Conference on Computer and Communications Security*. ACM, 65–79. doi:10.1145/3579856.3590335

[13] Yuxin Chen, Peng Tang, Weidong Qiu, and Shujun Li. 2025. Using LLMs for Automated Privacy Policy Analysis: Prompt Engineering, Fine-Tuning and Explainability. 12 pages. doi:10.48550/arXiv.2503.16516 arXiv:2503.16516 [cs.CL]

[14] Caitlin D. Cottrill and Piyushimita 'Vonu' Thakuriah. 2013. Privacy in context: an evaluation of policy-based approaches to location privacy protection. *International Journal of Law and Information Technology* 22, 2 (2013), 178–207. doi:10.1093/ijlit/eat014

[15] Hao Cui, Rahmadi Trimananda, Athina Markopoulou, and Scott Jordan. 2023. PoliGraph: Automated Privacy Policy Analysis using Knowledge Graphs. In *Proceedings of the 32nd USENIX Security Symposium*. USENIX Association, 1037–1054. https://www.usenix.org/conference/usenixsecurity23/presentation/cui

[16] Benjamin Fabian, Tatiana Ermakova, and Tino Lentz. 2017. Large-scale readability analysis of privacy policies. In *Proceedings of the 2017 International Conference on Web Intelligence*. ACM, 18–25. doi:10.1145/3106426.3106427

[17] Alvan R. Feinstein and Domenic V. Cicchetti. 1990. High agreement but low Kappa: I. The problems of two paradoxes. *Journal of Clinical Epidemiology* 43, 6 (1990), 543–549. doi:10.1016/0895-4356(90)90158-L

[18] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024. Retrieval-Augmented Generation for Large Language Models: A Survey. doi:10.48550/arXiv.2312.10997 arXiv:2312.10997 [cs.CL]

[19] Shahram Ghahremani and Uyen Trang Nguyen. 2024. Comprehensive evaluation of privacy policies using the contextual integrity framework. *Security and Privacy* 7, 4, Article e380 (2024), 26 pages. doi:10.1002/spy2.380

[20] Arda Goknil, Femke B. Gelderblom, Simeon Tverdal, Shukun Tokas, and Hui Song. 2024. Privacy policy analysis through prompt engineering for LLMs. doi:10.48550/arXiv.2409.14879 arXiv:2409.14879 [cs.CL]

[21] Kilem Li Gwet. 2008. Computing inter-rater reliability and its variance in the presence of high agreement. *Brit. J. Math. Statist. Psych.* 61, 1 (2008), 29–48. doi:10.1348/000711006X126600

[22] Kevin A. Hallgren. 2012. Computing Inter-Rater Reliability for Observational Data: An Overview and Tutorial. *Tutorials in Quantitative Methods for Psychology* 8, 1 (2012), 23–34. doi:10.20982/tqmp.08.1.p023

[23] Hamza Harkous, Kassem Fawaz, Rémi Lebret, Florian Schaub, Kang G. Shin, and Karl Aberer. 2018. Polisis: Automated Analysis and Presentation of Privacy Policies Using Deep Learning. In *Proceedings of 27th USENIX Security Symposium*. USENIX Association, 531–548. https://www.usenix.org/conference/usenixsecurity18/presentation/harkous

[24] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2024. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information System* 43, 2, Article 42 (2024), 50 pages. doi:10.1145/3703155

[25] Bowen Jin, Jinsung Yoon, Jiawei Han, and Sercan O. Arik. 2024. Long-Context LLMs Meet RAG: Overcoming Challenges for Long Inputs in RAG. doi:10.48550/arXiv.2410.05983 arXiv:2410.05983 [cs.CL]

[26] Hailey Joren, Jianyi Zhang, Chun-Sung Ferng, Da-Cheng Juan, Ankur Taly, and Cyrus Rashtchian. 2025. Sufficient Context: A New Lens on Retrieval Augmented Generation Systems. doi:10.48550/arXiv.2411.06037 arXiv:2411.06037 [cs.CL]

[27] Simon Leigh, Jing Ouyang, and Chris Mimnagh. 2017. Effective? Engaging? Secure? Applying the ORCHA-24 framework to evaluate apps for chronic insomnia disorder. *BMJ Ment Health* 20, 4, Article e20 (2017), 7 pages. doi:10.1136/eb-2017-102751

[28] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 9459–9474. https://proceedings.neurips.cc/paper_files/paper/2020/file/6b49 3230205f780e1bc26945df7481e5-Paper.pdf

[29] Yunxuan Li, Pascal Hirmer, and Christoph Stach. 2023. CV-Priv: Towards a Context Model for Privacy Policy Creation for Connected Vehicles. In *Proceedings of the 2023 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events*. IEEE, 583–588. doi:10.1109/PerComWorkshops56833.2023.10150231

[30] Lei Liang, Mengshu Sun, Zhengke Gui, Zhongshu Zhu, Zhouyu Jiang, Ling Zhong, Yuan Qu, Peilong Zhao, Zhongpu Bo, Jin Yang, Huaidong Xiong, Lin Yuan, Jun Xu, Zaoyang Wang, Zhiqiang Zhang, Wen Zhang, Huajun Chen, Wenguang Chen, and Jun Zhou. 2024. KAG: Boosting LLMs in Professional Domains via Knowledge Augmented Generation. doi:10.48550/arXiv.2409.13731

[31] Thomas Linden, Hamza Harkous, and Kassem Fawaz. 2020. The Privacy Policy Landscape After the GDPR. *Proceedings on Privacy Enhancing Technologies* 2020 (2020), 47–64. doi:10.2478/popets-2020-0004

[32] Jerry Liu. 2022. *LlamaIndex*. doi:10.5281/zenodo.1234

[33] Nathan Malkin. 2023. Contextual Integrity, Explained: A More Usable Privacy Definition. *IEEE Security & Privacy* 21, 1 (2023), 58–65. doi:10.1109/MSEC.2022.3201585

[34] Kirsten Martin and Helen Nissenbaum. 2016. Measuring Privacy: An Empirical Test Using Context to Expose Confounding Variables. *Columbia Science and Technology Law Review* 18, 1 (2016), 176–218. doi:10.7916/stlr.v18i1.4015

[35] Aleecia M. McDonald and Lorrie Faith Cranor. 2008. The Cost of Reading Privacy Policies. *I/S: A Journal of Law and Policy for the Information Society* 4, 3 (2008), 543–568. http://hdl.handle.net/1811/72839

[36] Keika Mori, Daiki Ito, Takumi Fukunaga, Takuya Watanabe, Yuta Takata, Masaki Kamizono, and Tatsuya Mori. 2025. Evaluating LLMs Towards Automated Assessment of Privacy Policy Understandability. In *Proceedings of the 2025 Symposium on Usable Security and Privacy*. 19 pages. doi:10.14722/usec.2025.23009

[37] Abhijith Athreya Mysore Gopinath, Shomir Wilson, and Norman Sadeh. 2018. Supervised and Unsupervised Methods for Robust Separation of Section Titles and Prose Text in Web Documents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. ACL, 850–855. doi:10.18653/v1/D18-1099

[38] Helen Nissenbaum. 2009. *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford University Press, Redwood City. doi:10.1515/9780804772891

[39] Helen Nissenbaum. 2011. A Contextual Approach to Privacy Online. *Daedalus* 140, 4 (2011), 32–48. doi:10.1162/DAED _a_00113

[40] Helen Nissenbaum. 2019. Contextual Integrity Up and Down the Data Food Chain. *Theoretical Inquiries in Law* 20, 1 (2019), 221–256. doi:10.1515/til-2019-0008

[41] Alessandro Oltramari, Dhivya Piraviperumal, Florian Schaub, Shomir Wilson, Sushain Cherivirala, Thomas B. Norton, N. Cameron Russell, Peter Story, Joel Reidenberg, and Norman Sadeh. 2018. PrivOnto: A semantic framework for the analysis of privacy policies. *Semantic Web* 9, 2 (2018), 185–203. doi:10.3233/SW-170283

[42] Oren Rachmil, Roy Betser, Itay Gershon, Omer Hofman, Nitay Yakoby, Yuval Meron, Idan Yankelev, Asaf Shabtai, Yuval Elovici, and Roman Vainshtein. 2025. Training-Free Policy Violation Detection via Activation-Space Whitening in LLMs. doi:10.48550/arXiv.2512.03994 arXiv:2512.03994 [cs.LG]

[43] Julie M. Robillard, Tanya L. Feng, Arlo B. Sporn, Jen-Ai Lai, Cody Lo, Monica Ta, and Roland Nadler. 2019. Availability, readability, and content of privacy policies and terms of agreements of mental health apps. *Internet Interventions* 17, Article 100243 (2019), 8 pages. doi:10.1016/j.invent.2019.100243

[44] David Rodriguez, Ian Yang, Jose M. Del Alamo, and Norman Sadeh. 2024. Large language models: a new approach for privacy policy analysis at scale. *Computing* 106 (2024), 3879–3903. doi:10.1007/s00607-024-01331-9

[45] Pranab Sahoo, Ayush Kumar Singh, Sriparna Saha, Vinija Jain, Samrat Mondal, and Aman Chadha. 2025. A Systematic Survey of Prompt Engineering in Large Language Models: Techniques and Applications. doi:10.48550/arXiv.2402.07927 arXiv:2402.07927 [cs.AI]

[46] Rohan Charudatt Salvi, Catherine Blake, and Masooda Bahir. 2024. PrivacyChat: Utilizing Large Language Model for Fine-Grained Information Extraction over Privacy Policies. In *Wisdom, Well-Being, Win-Win: 19th International Conference, iConference 2024, Changchun, China, April 15–26, 2024, Proceedings, Part I*. Springer, 223–231. doi:10.1007/978-3-031-57850-2_17

[47] Alireza Savand and Aaron Swartz. 2025. *html2text: Convert HTML to Markdown-formatted text.* https://github.com/Alir3z4/html2text/

[48] Florian Schaub, Travis D. Breaux, and Norman Sadeh. 2016. Crowdsourcing privacy policy analysis: Potential, challenges and best practices. *it - Information Technology* 58, 5 (2016), 229–236. doi:10.1515/itit-2016-0009

[49] Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The Prompt Report: A Systematic Survey of Prompting Techniques. doi:10.48550/arXiv.2406.06608 arXiv:2406.06608 [cs.CL]

[50] Mukund Srinath, Shomir Wilson, and C. Lee Giles. 2021. Privacy at Scale: Introducing the PrivaSeer Corpus of Web Privacy Policies. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. ACL, 6829–6839. doi:10.18653/v1/2021.acl-long.532

[51] State of California, USA. 2004. California Online Privacy Protection Act (CalOPPA). US state law. https://leginfo.legislature.ca.gov/faces/codes_displayText.xhtml?division=8.&chapter=22.&lawCode=BPC

[52] State of California, USA. 2020. California Consumer Privacy Act (CCPA). US state law. https://leginfo.legislature.ca.gov/faces/codes_displayText.xhtml?division=3.&part=4.&lawCode=CIV&title=1.81.5

[53] Chenhao Tang, Zhengliang Liu, Chong Ma, Zihao Wu, Yiwei Li, Wei Liu, Dajiang Zhu, Quanzheng Li, Xiang Li, Tianming Liu, and Lei Fan. 2023. PolicyGPT: Automated Analysis of Privacy Policies with Large Language Models. doi:10.48550/arXiv.2309.10238 arXiv:2309.10238 [cs.CL]

[54] Peng Tang, Xin Li, Yuxin Chen, Weidong Qiu, Haochen Mei, Allison Holmes, Fenghua Li, and Shujun Li. 2024. A Comprehensive Study on GDPR-Oriented Analysis of Privacy Policies: Taxonomy, Corpus and GDPR Concept Classifiers. doi:10.48550/arXiv.2410.04754 arXiv:2410.04754 [cs.CR]

[55] The European Parliament and The Council of The European Union. 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance). EU law, Official Journal of the European Union, L 119. https://eur-lex.europa.eu/eli/reg/2016/679/oj

[56] Damiano Torre, Sallam Abualhaija, Mehrdad Sabetzadeh, Lionel Briand, Katrien Baetens, Peter Goes, and Sylvie Forastier. 2020. An AI-assisted Approach for Checking the Completeness of Privacy Policies Against GDPR. In *Proceedings of the 2020 IEEE 28th International Requirements Engineering Conference*. IEEE, 136–146. doi:10.1109/RE48521.2020.00025

[57] UK Parliament. 2018. Data Protection Act 2018. UK law. https://www.legislation.gov.uk/ukpga/2018/12/contents/enacted

[58] Karl van der Schyff, Suzanne Prior, and Karen Renaud. 2024. Privacy policy analysis: A scoping review and research agenda. *Computers & Security* 146, Article 104065 (2024), 14 pages. doi:10.1016/j.cose.2024.104065

[59] Shomir Wilson, Florian Schaub, Aswarth Abhilash Dara, Frederick Liu, Sushain Cherivirala, Pedro Giovanni Leon, Mads Schaarup Andersen, Sebastian Zimmeck, Kanthashree Mysore Sathyendra, N. Cameron Russell, Thomas B. Norton, Eduard Hovy, Joel Reidenberg, and Norman Sadeh. 2016. The Creation and Analysis of a Website Privacy Policy Corpus. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. ACL, 1330–1340. doi:10.18653/v1/P16-1126

[60] Shomir Wilson, Florian Schaub, Frederick Liu, Kanthashree Mysore Sathyendra, Daniel Smullen, Sebastian Zimmeck, Rohan Ramanath, Peter Story, Fei Liu, Norman Sadeh, and Noah A. Smith. 2018. Analyzing Privacy Policies at Scale: From Crowdsourcing to Automated Annotations. *ACM Transactions on the Web* 13, 1, Article 1 (2018), 29 pages. doi:10.1145/3230665

[61] Tongshuang Wu, Michael Terry, and Carrie Jun Cai. 2022. AI Chains: Transparent and Controllable Human-AI Interaction by Chaining Large Language Model Prompts. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*. ACM, Article 385, 22 pages. doi:10.1145/3491102.3517582

[62] Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. 2024. C-Pack: Packed Resources For General Chinese Embeddings. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 641–649. doi:10.1145/3626772.3657878

[63] Qinge Xie, Karthik Ramakrishnan, and Frank Li. 2025. Evaluating Privacy Policies under Modern Privacy Laws at Scale: An LLM-Based Automated Approach. In *Proceedings of the 34th USENIX Conference on Security Symposium*. USENIX Association, Article 298. https://www.usenix.org/conference/usenixsecurity25/presentation/xie

[64] Mian Yang, Vijayalakshmi Atluri, Shamik Sural, and Ashish Kundu. 2025. Automated Privacy Policy Analysis Using Large Language Models. In *Data and Applications Security and Privacy XXXIX: 39th IFIP WG 11.3 Annual Conference on Data and Applications Security and Privacy, DBSec 2025, Gjøvik, Norway, June 23-24, 2025, Proceedings*. Springer, 23–43. doi:10.1007/978-3-031-96590-6_2

[65] Haiyue Yuan, Matthew Boakes, Xiao Ma, Dongmei Cao, and Shujun Li. 2023. Visualising Personal Data Flows: Insights from a Case Study of Booking.com. In *Intelligent Information Systems: CAiSE Forum 2023, Zaragoza, Spain, June 12–16, 2023, Proceedings*. Springer, 52–60. doi:10.1007/978-3-031-34674-3_7

[66] Haiyue Yuan, Ali Raza, Nikolay Matyunin, Jibesh Patra, and Shujun Li. 2024. A Graph-Based Model for Vehicle-Centric Data Sharing Ecosystem. In *Proceedings of the 2024 IEEE 27th International Conference on Intelligent Transportation Systems*. IEEE, 3587–3594. doi:10.1109/ITSC58415.2024.10919888

[67] Razieh Nokhbeh Zaeem, Rachel L. German, and K. Suzanne Barber. 2018. PrivacyCheck: Automatic Summarization of Privacy Policies Using Data Mining. *ACM Transactions on Internet Technology* 18, 4, Article 53 (2018), 18 pages. doi:10.1145/3127519

# Appendix

**Using your personal data with transparencye**

In the course of its activities, Renault Group collects, uses and stores some of your personal data, i.e. information that makes it possible to identify you. Renault Group intends to ensure the greatest transparency in the processing it performs on the personal data you provide to it or on the personal data it collects through the various contacts you may have with it.

Converting non-table web text to a text segment

\*\*\* Using your personal data with transparency
In the course of its activities, Renault Group collects, uses and stores some of your personal data, i.e. information that makes it possible to identify you. Renault Group intends to ensure the greatest transparency in the processing it performs on the personal data you provide to it or on the personal data it collects through the various contacts you may have with it.

| Processing activity: Why we use your information? | What information is collected? | Lawful basis of processing | Where is the information collected from? | Specific retention periods |
|---|---|---|---|---|
| My Honda Plus mobile app | | | | |
| To set up and manage your Honda account, including sending you service notifications. | · Country<br>· Language<br>· Email address<br>· Address (including postcode)<br>· Mobile number<br>· Password<br>· Honda ID, which is a unique identifier generated for each Honda account holder during the registration process<br>· Your device ID which will be linked to your Honda ID where you opt to login using the touch ID, voice command or facial recognition features of your phone. We do not process your biometric data which remains on your device.<br>· Your preferences which you set in the App | To take steps at your request to enter a contract with you, for the ongoing performance, management and facilitation of such contract.<br><br>A failure to provide this information will unfortunately mean you will not be able to open a Honda account with us.<br><br>To the extent that the processing goes beyond what is necessary for the contract, the processing is necessary for our legitimate interest to provide you with a good customer experience and to provide you with features to keep your account secure. | From you (via the App) | |

Converting table web text to a text segment

_table_ | Processing activity: Why we use your information? | What information is collected? | Lawful basis of processing | Where is the information collected from? | Specific retention periods |
| My Honda Plus mobile app | To set up and manage your Honda account, including sending you service notifications. | Country, Language, Email address, Address (including postcode), Mobile number, Password, Honda ID, which is a unique identifier generated for each Honda account holder during the registration process, Your device ID which will be linked to your Honda ID where you opt to login using the touch ID, voice command or facial recognition features of your phone. We do not process your biometric data which remains on your device., Your preferences which you set in the App | To take steps at your request to enter a contract with you, for the ongoing performance, management and facilitation of such contract., A failure to provide this information will unfortunately mean you will not be able to open a Honda account with us., To the extent that the processing goes beyond what is necessary for the contract, the processing is necessary for our legitimate interest to provide you with a good customer experience and to provide you with features to keep your account secure. | From you (via the App) |

Fig. 5. Example of converting non-table and table text to text segments

**Prompt**

content:

[1] You are an expert to analyse the **TEXT SEGMENT** to extract data flows.

[2] If **TEXT SEGMENT** starts with _table_: Treat | as separators; Treat first line as the table heading; Treat second line as the table content.

[3] Read and understand the **TEXT SEGMENT** and then strictly follow the below rules to produce your responses:

    (a) If the **TEXT SEGMENT** at least talk about one party collects data or personal information from another party, or a party shares data or personal information to another party, OUTPUT the extracted data flows in multiple JSON objects.

    (b) The JSON objects must use the format:

```
{"Output":[{
    "data_sender": "",
    "data_type": [],
    "data_receiver": []
}]}
```

    (c) Respond only with valid JSON.

    (d) Each data flow represents one party (i.e., data_receiver) collects personal data (i.e., data_type) from another party (i.e., data_sender), or a party (i.e., data_sender) shares personal data (i.e., data_type) to another party (i.e., data_receiver).

    (e) For data_types: extract all atomic personal data (i.e., data_type) following these rules: (1) when dealing with sentences that have combined data_types, split them into individual data_types for a clearer representation; (2) each data_type MUST appear in **TEXT SEGMENT**, and do not change the cases; (3) DO NOT INCLUDE any other text in the answer such as input or query text or your deduction or your explanation; (4) remove Pronouns in the identified strings; (5) DO NOT INCLUDE specific addresses, postcodes, email addresses, companies, organisations, or geographical information; (6) if you can not identify a data_type, leave it empty.

    (f) For data_sender/data_receiver: (1) when dealing with sentences that have combined data_receivers or data_senders, split them into individual data_receiver or data_sender for a clearer representation; (2) each data_sender or data_receiver string MUST appear in **TEXT SEGMENT**, do not change the cases; (3) if no data_sender is explicitly stated in the **TEXT SEGMENT**, leave data_sender empty; (4) if no data_receiver is explicitly stated in the **TEXT SEGMENT**, leave data_receiver empty.

    (g) OUTPUT only "None" for other unrelated scenarios

content:

    **TEXT SEGMENT**: One Privacy Policy Text Segment

Fig. 6. Prompt template for extracting data flows following the Groq Chat Completions API format

**Prompt**

role: system content:

[1] You are an expert in categorising the **INPUT DATA TYPE** provided by the user to extract key information about data collection/sharing.

[2] Use the **TEXT SEGMENT** given by the user, along with the knowledge and understanding of data category and description provided in the **CONTEXT** list to perform the categorisation.

[3] Please strictly follow the below rules when answering questions:

(a) Output must follow the JSON format shown below.

```
{"Output": [{
    "DataCategory": "data_category",
    "DataType": "data_type",
    "InputText": "input_text",
}]}
```

(b) Respond only with valid JSON.

(c) Do not include any other text (e.g., input or query text).

(d) DataCategory is the identified data category defined in **CONTEXT**. Do not create new categories.

(e) DataType is the **INPUT DATA TYPE**.

(f) InputText is the **TEXT SEGMENT**.

(g) OUTPUT only "None", if no category can be found.

role: user content:

[1] **INPUT DATA TYPE**: Example Data Type

[2] **TEXT SEGMENT**: Privacy Policy Text Segment

[3] **CONTEXT**:

Data category: Data category 1
Data description: Description of data category 1
Example data types: data1, data2, data3, ...

[4] **CONTEXT**:

Data category: Data category 2
Data description: Description of data category 2
Example data types: data1, data2, data3, ...

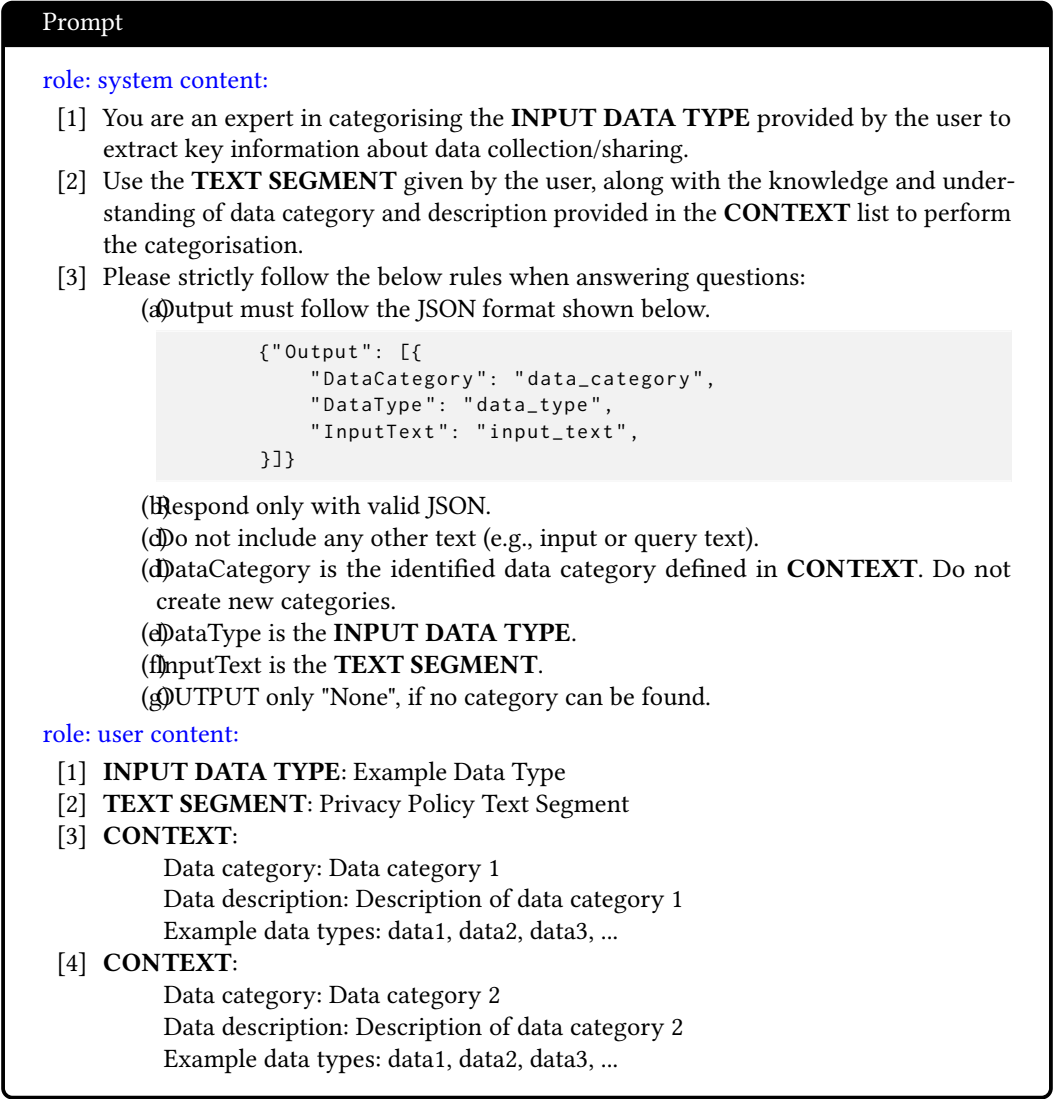Fig. 7. Prompt template for identifying data categories following the Groq Chat Completions API format

Table 13. Data collection and sharing practices stated in Google app safety section

| Data Category | Audi | | Ford | | Honda | | Hyundai | | Kia | | Lexus | | Nissan | | Polestar | | Renault | | Vauxhall | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | S | C | S | C | S | C | S | C | S | C | S | C | S | C | S | C | S | C | S | C |
| **Location** | | | | | | | | | | | | | | | | | | | | |
| Approximate location | | | | | | | | | | | | | | | | | Y | | Y | |
| Precise location | | Y | | | | | | | | | | | | | | | Y | | Y | |
| **Personal Info** | | | | | | | | | | | | | | | | | | | | |
| Name | | Y | | Y | Y | Y | | Y | | Y | Y | Y | | | | Y | Y | Y | Y | Y |
| Email address | | Y | | Y | Y | Y | | Y | | Y | Y | Y | | Y | Y | Y | Y | Y | Y | Y |
| User IDs | Y | Y | | Y | Y | Y | | | | | | Y | | | Y | Y | Y | Y | Y | |
| Address | | Y | | Y | Y | Y | | | | | Y | Y | | | | | Y | Y | | Y |
| Phone number | | Y | | Y | | | | Y | Y | | Y | Y | | | Y | Y | Y | Y | Y | Y |
| Race and ethnicity | | | | | | | | | | | | | | | | | | | | |
| Political or religious beliefs | | | | | | | | | | | | | | | | | | | | |
| Sexual orientation | | | | | | | | | | | | | | | | | | | | |
| Other info | | | | | | Y | | | | | | | | | | | Y | Y | Y | Y |
| **Financial Info** | | | | | | | | | | | | | | | | | | | | |
| User payment info | | | | | | | Y | | | | | | | | | | | | | |
| Purchase history | | | | | | | | | | | | | | | | | | | | |
| Credit score | | | | | | | | | | | | | | | | | | | | |
| Other financial info | | | | | | | | | | | | | | | | | | | | |
| **Health & Fitness** | | | | | | | | | | | | | | | | | | | | |
| Health info | | | | | | | | | | | | | | | | | | | | |
| Fitness info | | | | | | | | | | | | | | | | | | | | |
| **Messages** | | | | | | | | | | | | | | | | | | | | |
| Emails | | | | | | | | | | | | | | | | | | | | |
| SMS or MMS | | | | | | | | | | | | | | | | | | | | |
| Other in-app messages | | | | | | | | | | | | | | | | | | | | |
| **Photos and Videos** | | | | | | | | | | | | | | | | | | | | |
| Photos | | | | | | | | | Y | | Y | | | | | Y | | | | |
| Videos | | | | | | | | | | | | | | | | | | | | |
| **Audio files** | | | | | | | | | | | | | | | | | | | | |
| Voice or sound recordings | | | | | | | | | | | | | | | | | | | | |
| Music files | | | | | | | | | | | | | | | | | | | | |
| Other audio files | | | | | | | | | | | | | | | | | | | | |
| **Files and docs** | | | | | | | | | | | | | | | | | | | | |
| Files and docs | | | | | | | | Y | Y | | | | | | | | | | Y | |
| **Calendar** | | | | | | | | | | | | | | | | | | | | |
| Calendar events | | Y | | | | | | | | | | | | | | | | | | |
| **Contacts** | | | | | | | | | | | | | | | | | | | | |
| Contacts | | | | | | | | Y | Y | | | | | | | | | | | |
| **App activity** | | | | | | | | | | | | | | | | | | | | |
| App interactions | | Y | | Y | Y | Y | | Y | | | Y | | | | | | Y | | Y | Y |
| In-app search history | | | | | | | | | | | | | | | | | | | | |
| Installed apps | | | | | | | | | | | | | | | | | | | | |
| Other user-generated content | | | | | Y | Y | | | | | Y | | | | | Y | | | | |
| Other actions | | | | | | | | | | | | | | | | | | | Y | Y |
| **Web browsing** | | | | | | | | | | | | | | | | | | | | |
| Web browsing history | | | | | | | | | | | | | | | | | | | | |
| **App info and performance** | | | | | | | | | | | | | | | | | | | | |
| Crash logs | | | | Y | Y | Y | | Y | Y | | Y | | | | | | Y | | Y | |
| Diagnostics | | | | Y | Y | Y | | Y | Y | | | Y | | | | | Y | | Y | |
| Other app performance data | | | | | | | | | | | | Y | | | | | Y | | | |
| **Device or other IDs** | | | | | | | | | | | | | | | | | | | | |
| Device or other IDs | | Y | | Y | Y | Y | | Y | Y | | Y | Y | | | | Y | Y | Y | Y | Y |

S: Data shared
C: Data collected

Table 14. Data collection and sharing practices represented using Apple privacy labels

| Data Category | Audi | | | Ford | | | Honda | | | Hyundai | | | Kia | | | Lexus | | | Nissan | | | Polestar | | | Renault | | | Vauxhall | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N | T | L | N |
| **Location** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Precise location | | | | | | | | | | | | | | Y | | | Y | | | | Y | | Y | | | | Y | | | Y |
| Coarse location | | | | | | | | | | | | Y | | | | | | | | | Y | | Y | | | | Y | | | |
| **Contact info** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Name | | | | | Y | | | Y | | | Y | | | Y | | | | | | | Y | | Y | | | Y | | | Y | |
| Email address | | | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | |
| Phone number | | | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | | | Y | | | | | | Y | |
| Physical address | | | | | Y | | | Y | | | | | | | | | Y | | | | Y | | Y | | | | | | Y | |
| Other user contact info | | | | | | | | | | | | | | | | | Y | | | Y | | | | | | | | | | |
| **Health and fitness** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Health | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Fitness | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Financial info** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Payment info | | | | | | | | | | | | | | | | | | | | | | | Y | | | | | | | |
| Credit info | | | | | | | | | | | | | | | | | | | | | | | Y | | | | | | | |
| Other financial info | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Sensitive info** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Contacts** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **User content** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Emails or text messages | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Photos or videos | | | | | | | | | | | | | | | | | Y | | | | | | | | | | | | | Y |
| Audio data | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Gameplay content | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Customer support | | | | | Y | | | | | | | | | | | | | | | | | | Y | | | | | Y | Y | |
| Other user content | | | | | | | | | | | | | | | | | | | | | | | Y | | | | | | | |
| **Browsing history** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Search history** | | | | | | | | | | | | | | Y | | | | | | | | | | | | | | | | Y |
| **Identifiers** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| User ID | | | | | Y | | | Y | | | | | | Y | | | Y | | | | | | Y | | Y | Y | | | Y | |
| Device ID | | | | | Y | | | Y | | | | | | Y | | | | | | | Y | Y | Y | | | | | | Y | |
| **Purchase history** | | | | | | | | | | | | | | | | | | | | | | | Y | | | | | | | |
| **Usage data** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Product interaction | | | | | Y | | | Y | | | | | | | | | | Y | Y | | Y | Y | | | | | Y | | | Y |
| Advertising data | | | | | | | | | | | | | | | | | | | | | | | | Y | | | | | | Y |
| Other usage data | | | | | | | | | | | | | | | | | | | | Y | | | | | | | | | | |
| **Diagnostics** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Crash data | | | | | Y | | | Y | | | | | | | | | | Y | | Y | | | | Y | | Y | | | | Y |
| Performance data | | | | | Y | | | Y | | | | | | | | | Y | | | Y | | | | Y | | | | | | Y |
| Other diagnostic data | | | | | | | | Y | | | | | | | | | | | | | | | | | | | | | | |
| **Surroundings** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Environment scanning | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Body** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Hands | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Head | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| **Other data types** | | | | | | | | Y | | | | | | Y | | | | | | | | | | | | | | | | |

T: "Data used to track you"

L: "Data linked to you"

N: "Data not linked to you"