# Privacy Scan

## 0. Administrative information

| | |
|---|---|
| **Project/Study title:** | Leveraging Generative Artificial Intelligence for Personalized Customer Acquisition in Tech Startups: A Study on Marketing Strategies in the AI Age |
| **Department / Faculty:** | Faculty of Science<br>Department of Information and Computing Science |
| **Contact details of project controller(s)\*:**<br>*Role, name, job title, department, e-mail.*<br>*\*Indicate who will be the everyday contact person.* | 1. Researcher, C.E.W. (Mike) Brachten, c.e.w.brachten@students.uu.nl<br>2. 1st Supervisor, R.L. (Slinger) Jansen, slinger.jansen@uu.nl<br>3. **2nd examiner, N.A. (Nico) Brand,** n.a.brand@uu.nl |
| **Initial personal data collection start / end date:** | From September 1 up until January 31 |
| **Reviewed by:** | **Frans Huigen, MSc.,** Data Steward Privacy of the Faculty of Science, f.huigen@uu.nl, or privacy-beta@uu.nl |
| **Privacy Scan Outcome** | **Approved** |

**Comments to this Privacy Scan by Frans Huigen:** The nature of the research, its topic and the approach give reason to be cautious when it comes to the processing of personal data. Generative AI and its use in a scientific context, especially in interaction with data subjects and the processing of their personal data, are under development and the ways to go about are not set in stone (if ever). Nonetheless, the GDPR proves an instrument with which we can assess the impact of the researchers activity to the data subjects. I would like to specifically point out that subsidiarity in this case is a difficult principle to adhere to, as the majority of LLMs in circulation are proprietarily owned and operated. Open Source models, let alone open source interaction, is absent or would prove an unreasonable effort for the student-researcher to employ it. As there is little to none academic work on this topic, Mr. Brachten is unfortunately forced to make do with what is offered.

For this specific activity I take into account that this concerns a Master's thesis research. Meaning that time and resources, as well as the relationship with the student-researcher, are different than in cases with PhD or postdoc research. This assessment then becomes a balancing act between the protection of the data subjects' privacy and subsequently the mandatory compliance to the GDPR, and ultimately facilitating the proof of competence of the student-researcher in the best way possible. In my view It's most important that the student-researcher shows understanding, willingness and competence to elaborate and explain the effects of the legal framework on his endeavour. Based on our e-mail exchange, on the discussion we had on July 9, and on the privacy scan (specifically: its executive summary) now laying before you; I see that he does. I endorse his documentation on the processing activity necessary for this research.

My involvement here of course doesn't necessarily stop with this commentary; if necessary the researcher and the ERB are very welcome to reach out to me in the course of the research. I too am learning and growing in this new field of research.

**Comment by Researcher:** For this risk assessment, LLMs have been utilized to process the content of large (legal) documents.

## 1. Executive summary

In a conversation on July 9, 2024, with Frans Huigen, Privacy Officer at Utrecht University, some unclarities in this document were identified. The documentation of privacy is spread across this DPIA, the DMP, the Method section of the long proposal, and the appendices A on Privacy and Ethics in the proposal. For extra justification of the major risks identified and its measures, a short executive summary is provided here.

**Collection of personal data**
The collection of personal data is done through *public* online sources, including social media. This may be done manually but also through web scraping of such data. The information that is collected is highlighted in section 4 of this DPIA.

**Information about collection of personal data**
Recipients are informed about the collection of their data in a clear and accessible way in accordance with Articles 13 and 14. In the first email, the footer of the email contains information about the purpose of the email and the fact that their personal information has been processed, together with a link to a privacy policy and how to exercise their data subject rights or opt-out. To avoid major bias in the form of the Hawthorne effect, the information provided is kept to a minimum to be compliant with the GDPR at first. A second follow-up email is sent after two weeks informing data subjects that they were part of a study, and more information about their GDPR rights and a possibility to opt-out is also provided. This email will provide the information in a more fully transparent way.

**Sharing of personal data with third parties**
In order to generate AI-personalized emails, some semi-anonymous data is shared with LLM providers. A data processing addendum is signed with OpenAI, Google Gemini's privacy policies has been reviewed, and the Llama 3 LLM is an open-source locally-ran LLM. The data that is shared is pseudonymized by removal of directly identifiable information, such as name and LinkedIn-URL. A prompt to the LLM could look like:

"*Write a personalized email selling [a product]. To a 25-year-old male student of the Master Business Informatics. He lives in Amsterdam, works as a teaching assistant at Utrecht University.*"

No further information is provided to the AI, and the generated email is copied, checked manually, and sent to the recipient. At no times do third parties have access to direct identifiers such as name and email address, nor can they access email tracking information.

While pseudonymization and data processing agreements are in place to hopefully minimize the chances of a data leak, one must also take into consideration that all collected data is from public sources, and a potential attacker could already access this data by themselves. The potential effects of such a data leak are therefore slim, as the only additional information gained is that a personalized email to a person has been created with AI.

**Informed consent**
One of the major warnings of the Quick Scan was informed consent, as in the current research design informed consent is not gathered before inclusion into the study. Informing the participants beforehand about the study and gathering informed consent would introduce major bias in the form of the Hawthorne effect and also potential phishing risks. If a participants signed up for a study, and know they may receive an email in the coming months, they may be more prone to seriously consider every email and therefore also more quickly click on links or fall victim to scams.

An opt-out system is also preferred over an opt-in system, as power analysis returns a sample size of 200 for sufficient power. If an opt-in system is utilized after the email is sent, using conventional estimates of a 2,66% click-through rate in email marketing, the number of emails to be sent would increase to ±7.500 emails. This would then be more reflective of a major spam operation instead of a study, it could also introduce bias in the form of self-selection bias, and it would also be more difficult to recruit sufficient recipients. Furthermore, the data would already be collected, but would go to waste.

**Exclusion of sensitive information**
While the DPIA, DMP and long proposal mention that sensitive information is excluded from the data set, it is not made very clear how this is done. The criteria of sensitive information used is that by the European Commission as defined in the GDPR. It could be that someone posts about their sexual orientation or political affiliation publicly. This information is ignored by the researcher while building a profile manually, and with automatic collection, the data is checked manually before passing it to the LLMs. Upon discovery of sensitive data, this is immediately permanently deleted from the dataset.

**Protocol when a participant acts on their rights**
When a participant acts on their data subject rights as defined in the GDPR, their requests are honored in an appropriate way. Data collected about them will be shared upon their request, requests for deletion or data processing restriction will lead to permanent deletion of their data in all datasets, and requests for modification or limiting their processing will be handled taking into account the effect on the research. If a data subject for example wants to for example modify user engagement data or limit processing in such a way that research is affected, the data is excluded from data analysis.

2. **Project purpose(s)**
   *Project description and its purposes.*

   The research focuses on evaluating the effectiveness of AI-personalized marketing strategies for customer acquisition. The study investigates how AI-generated personalized messages impact user engagement and conversion rates compared to generic messages. The study will also compare different Large Language Models (LLMs), such as GPT-4o and Gemini 1.5 Pro, to assess their effectiveness in marketing personalization. The research aims to provide empirical insights that can help organizations optimize their marketing strategies, ultimately determining whether investing in newer AI models is beneficial.

3. **Description of data subjects -**
   **Definition**:
   *Who are your data subjects? Are they vulnerable? What in- and exclusion criteria do you use?*

   The data subjects are potential leads for the product or service that is being sold. They are not inherently vulnerable. Exclusion criteria are those younger than 18 years old, older than 65 years of age, those with (suspected) learning or communication difficulties, or those with whom the researcher has a personal relationship. This is determined based on reasonable suspicion of the researcher based on online behavior or other indicators. Age and gender are identified based on self-reporting, or by using an open source Machine Learning algorithm running locally.

   **Targeting**:
   *How are you approaching your data subjects?*

   The data subjects are approached through email twice; the first email is the actual experiment and tracks email opens, link clicks, and replies. The second email is sent about two weeks

after the initial email and explains the study, its purpose, and how the recipient can exercise their data rights.

**Number**:
*How many data subjects are you targeting? How many responses do you expect?*

The number of recipients I'm targeting is 250-300. While 200 recipients yield the statistical power I'm after, it should be accounted for that some emails will likely bounce, and that some recipients will opt-out after receiving the first or second email.

**Nature of the relationship between data subjects and controllers**:
*What is the nature of the relationship between data subjects and the project controllers? Is there a (possible) imbalance of power?*

The relationship between data subjects and the project controllers can be defined as recipient/sender, customer/marketeer, and subject/researcher. There exists no inherent imbalance of power in these relationships, although it could be argued that a possible imbalance of power is present in the form of being unaware that they are part of a study. This is addressed by including opt-out options in all emails and providing more information in the follow-up email.


4. **Description of the categories of personal data**
   *List the type of data and the purpose for collecting each type of personal data.*

   Identifiable information
   - **Name:** Collected to personalize the email, for example in the greeting. Collected from social media.
   - **Email address:** Required to send the email. Collected from social media, public online sources, or by guessing based on company email structure (i.e. firstname.lastname@company.com).
   - **URLs:** Indexed to document the source of the data, and if possible, for automatic data collection by scraping.

   Demographic information
   - **Gender:** Used for AI email personalization purposes (e.g. pronouns). Collected from publicly available information, researcher judgement, or a local open source ML model.
   - **Age:** Used for AI email personalization purposes (e.g. formality). Collected from publicly available information, researcher judgement, or a local open source ML model.
   - **Location:** Used for AI email personalization purposes. Collected from publicly available information.

   Professional information
   - **Job title:** Used for AI email personalization purposes. Collected from publicly available information.
   - **Organization:** Used for AI email personalization purposes. Collected from publicly available information.

   Personal interests
   - **Hobbies, interests, and likes:** Used for AI email personalization purposes. Collected from publicly available information or researcher judgement (e.g. playing sports in profile picture). Here, manual human scrutiny is involved to prevent collection of sensitive personal data.

Personal descriptions
- **Introduction, about me section, or social media bio:** Used for AI email personalization purposes. Collected from publicly available information and subject to human scrutiny for deletion of sensitive or directly identifiable information.

5. **Description of the processing of personal data -**

**Data source**:
*Directly from individuals, secondary data?*

The data is collected from public sources on the internet, i.e., their social media and organizations' website. The personal data is therefore likely directly published and created by the individuals.

**Data storage and processing**:
*What tools/resources will you use to collect, store and process data?*

The data is collected manually or through a (Python) web scraper. The data is then stored in a csv file and manually checked using Microsoft Excel. Then, the pseudonymized data is passed to a (proprietary or open source) LLM through APIs. Finally, the email is sent though self-hosted servers, and tracked through a self-hosted tracking pixel and tracked links. The data is then collected from the server after a couple of weeks and analyzed with SPSS.

**Data access**:
*Who will have access to the personal data in every stage of your research project, and why?*

Only the researcher will have access to directly identifiable personal data. For grading and archiving purposes, fully anonymized email tracking data may be archived and shared within the university.

**Data retention**:
*How long will data be stored and when will you delete/anonymize the data?*

The personal data the AI is trained on will be stored until the emails are generated, after which they are pseudonymized for scientific archiving purposes.

**Collection times:**
*When and how often will you collect data?*

The data will be collected when the AI models and all tools are ready; likely around October-November 2024. The data collection will take about a month.

**Collection/processing location**:
*Where (geographically) is data collection/processing taking place?*

Data collection is taking place in the Netherlands. Processing the pseudonymized data will take place in the European Union, due to sharing it with LLM providers (signed data processing agreements are in place where applicable).

**Data minimisation measures**:
*How are you minimizing the amount of personal data you are processing? E.g., pseudonymization, encryption, removing data, avoiding collecting (sensitive) personal data, etc.*

The personal data used for generating emails is pseudonymized after generating the emails (id, name, email <-> personal data). Furthermore, two weeks after the email is sent, a follow-up email is sent, and the final data can be pseudonymized (id <-> name, email). Before the data is fed to the LLM, it is pseudonymized, and manually and automatically (using locally ran, open source LLMs) reviewed for senstitive data.

6. **Description of information provided to data subjects**
   *Describe how and what information is provided to data subjects about their personal data processing.*

   The recipients will be informed about what information is collected and for which purpose in the footer of the first email. This email is tracked, and a recipient has the option to opt-out. Furthermore, a follow-up email is sent two weeks after the initial email with information about the study and the option to review or delete their data.

7. **Description of how data subjects can exercise their data subject rights**
   *How do you respond to data subjects exercising their data subject rights?*

   The recipients have the right to access or delete their data, and also the option to restrict processing upon demand. Contact details are provided, and the online tool (which is to be developed) will contain options for subjects to exercise their data subject rights. If requests to access the collected data emerge, this will be provided by email. In case requests for deletion, modification, or limiting processing of the data arise, the data is deleted from the study and the servers.

8. **Description of the lawful basis for processing**
   *Which lawful basis are you relying for (each of) your data processing?*

   The lawful basis of data processing is Article 6e of the GDPR, for a task contributing to public interest. This is in the form of research fulfilling a knowledge gap in current research regarding the effectiveness of LLMs in personalized emails.

9. **Description of measures to ensure compliance by processors and/or joint controllers**
   *E.g, data processing agreements (DPAs), joint controllers agreements, etc.*

   **OpenAI**
   A Data Processing Agreement (DPA) has been signed with OpenAI for GPT-4/GPT-3.5 API access. This DPA ensures compliance with applicable data protection laws such as GDPR and U.S. Privacy Laws. OpenAI, as a data processor, is obligated to process data only per the written instructions of the data controller (customer) and in a manner that meets or exceeds required privacy protections. OpenAI's obligations include ensuring confidentiality, engaging subprocessors with comparable protections, and providing assistance in responding to data subjects' rights requests and data protection impact assessments.

10. **Description of planned transfers of personal data to other countries outside the EU**
    *Will such transfers take place? If yes, how will you ensure appropriate transfer mechanisms?*

    **OpenAI**
    Personal data transfers outside the EU are planned under the agreement with OpenAI. OpenAI Ireland Ltd will process Customer Data originating in the EEA or Switzerland. Transfers to other OpenAI affiliates in jurisdictions without adequate data protection will be based on intra-group agreements incorporating Standard Contractual Clauses (SCCs) or other approved mechanisms. For UK-based customers, data transfers will be governed by SCCs as amended by the UK Addendum. These measures ensure compliance with GDPR and provide appropriate safeguards for data transfers.

11. **Obtaining, consulting and dealing with data subjects' views of the processing**

*As you are planning to process data subject's personal data, you should consult them before you start your project, to find out what they think about your plans.*

Peer students, my supervisors, a privacy officer, a data manager, a member of the ethics board, and many others in my personal life have been consulted for advice and their opinion on the processing.

## 12. Preliminary risk assessment
*Are there potential risks of physical, material, or non-material damages to data subjects derived from the processing of their data, or from a data breach? What safeguards/measures have you adopted to minimise these risks?*

The risks that have been identified are re-identification of data subjects, data leak in user engagement data or profiles before pseudonymization,

## 12a. Safeguards and measures

- **Data Minimization and Anonymization:** Pseudonymization to relative anonymity or full anonymity before processing to ensure that personal identifiers are not easily linked back to individuals. This reduces the risk of re-identification if the data is accessed improperly.
- **Manual and Automated Reviews:** The data is thoroughly reviewed manually and by using locally-run open-source language models at every stage (collection, before passing to LLM, before sending, before archiving) to detect and eliminate any sensitive information before processing.
- **Encryption:** The drive containing the stored data is encrypted using FileVault. This ensures data is unreadable without proper decryption keys.
- **UU-Hosted servers:** Where possible, a UU-hosted server (emailstudy.science.uu.nl) is used.
- **Data Processing Agreement (DPA):** The signed DPA with OpenAI requires OpenAI to implement appropriate technical and organizational measures to protect personal data.
- **Opt-out Mechanisms:** Easy-to-use opt-out mechanisms are provided, allowing data subjects to withdraw from the study or request data deletion. Initia land follow-up communications ensure that participants are aware of these options.
- **Data Subject Rights:** Data subjects are informed of their rights under GDPR, including the right to access, rectify, erase, and restrict the processing of their personal data.
- **Secure Transfer Mechanisms:** When transferring data, secure methods such as encrypted connections (TLS/SSL) are used.