

# PRIPA: A Tool for Privacy-Preserving Analytics of Linguistic Data

Jeremie Clos<sup>1</sup>, Emma McLaughlin<sup>1</sup>, Pepita Barnard<sup>1</sup>, Elena Nichele<sup>1</sup>,  
Dawn Knight<sup>2</sup>, Derek McAuley<sup>1</sup>, Svenja Adolphs<sup>1</sup>

<sup>1</sup> University of Nottingham

{jeremie.clos, emma.mclaughlin, pepita.barnard, elena.nichele,  
derek.mccauley, svenja.adolphs}@nottingham.ac.uk

<sup>2</sup> Cardiff University

knightd5@cardiff.ac.uk

## Abstract

The days of large amorphous corpora collected with armies of Web crawlers and stored indefinitely are, or should be, coming to an end. There is a wealth of hidden linguistic information that is increasingly difficult to access, hidden in personal data that would be unethical and technically challenging to collect using traditional methods such as Web crawling and mass surveillance of online discussion spaces. Advances in privacy regulations such as GDPR and changes in the public perception of privacy bring into question the problematic ethical dimension of extracting information from unaware if not unwilling participants. Modern corpora need to adapt, be focused on testing specific hypotheses, and be respectful of the privacy of the people who generated its data. Our work focuses on using a distributed participatory approach and continuous informed consent to solve these issues, by allowing participants to voluntarily contribute their own censored personal data at a granular level. We evaluate our approach in a three-pronged manner, testing the accuracy of measurement of statistical measures of language with respect to standard corpus linguistics tools, evaluating the usability of our application with a participant involvement panel, and using the tool for a case study on health communication.

**Keywords:** privacy-preserving linguistics, corpus linguistics, software tools

## 1. Introduction

There is a wealth of hidden linguistic information which is increasingly difficult to access, hidden in personal and private data that would be unethical and technically challenging to collect using traditional methods such as Web crawling and mass surveillance of online discussion spaces. Additionally, advances in privacy regulations and changes in the zeitgeist bring into question the problematic ethical dimension of extracting such information from unaware if not unwilling participants.

Since the generation of knowledge from large amounts of empirical data is at the heart of corpus linguistics, its practitioners have long sought ways to protect the privacy of those who have generated it. However, so far the use of privacy-preserving methods has focused on post hoc processing such as automated anonymisation and de-identification. Those automated methods are severely lacking when faced with modern methods of re-identification and de-anonymisation. Non-automated methods on the other hand are not as scalable.

As a first step towards addressing this issue, we developed PRIPA<sup>1</sup>, a software tool using a distributed participatory approach and continuous informed consent by allowing participants to stay in control of their data, and only voluntarily contribute their own censored personal data on their own terms (McLaughlin et al., 2022).

We evaluate our prototype by producing a compar-

ison of word frequencies and collocate association scores between two standard state-of-the-art systems and PRIPA, showing that PRIPA is on par with those tools for some of their common features. We produce a small scale quantitative and qualitative evaluation of the tool by users of different levels of expertise, highlighting some key challenges in the production of privacy-preserving linguistic analysis tools.

This paper is structured as follows: In section 2, we discuss the overall methodology of PRIPA: general design for continuous consent, and software architecture. In section 3 we describe our evaluation methodology. We will finally conclude with key challenges and recommendations for further development in section 4.

## 2. Privacy-Preserving Corpus Linguistics

Privacy-preserving technologies allow for the processing of personal data in a way that minimises risks towards the privacy of the people who generated it (Noble et al., 2019). There are several approaches to privacy-preserving analytics, which rely on different tools to protect this privacy: trusted execution environments, homomorphic encryption, secure multi-party computation, differential privacy, and personal data stores. We opt for the personal data store approach to privacy-preserving analytics because it is the most compatible with the notion of continuous consent and granular sharing of data that is key to PRIPA, however those approaches are not mutually exclusive and further development of the tool will investigate the use of additional privacy layers such as differential privacy for statistics which cannot be computed locally.

<sup>1</sup><https://c19comms.wp.horizon.ac.uk/pripa>

Other approaches to privacy-preserving analytics use the personal data store approach. Mozilla’s Rally project (Mozilla, 2022) for example focuses on passive monitoring of data volunteers for Web-based data. One key difference with PRIPA is that Rally does not differentiate websites of interest, while PRIPA predetermines a set of websites of interest from which statistics are collected. Additionally, Rally monitors a wider set of interactions such as videos watched, time spent on each page, and all domain names of websites visited during the experiment while PRIPA focuses on specific linguistic items.

In the remaining parts of this section, we will describe the overall design principles of PRIPA and contrast them to the requirements of the General Data Protection Regulation (GDPR). We will then describe two key aspects of PRIPA for privacy-preserving corpus linguistics: the software architecture allowing data to be collected according to our key principles, and the user interface design allowing for the informed consent of users to be monitored at each key step of the data collection process.

## 2.1. Design principles of PRIPA

Being privacy-preserving by design involves the adherence to a set of principles, described in Table 1. Instead of collecting the data on the online discussion platform, we recruit participants who install a plugin into their Web browser. The PRIPA plugin then allows participants to enrol themselves into different experiments. Those experiments specify multiple things: the websites that will be watched, the words that will be observed, and the statistics that will be collected.

The principles used to develop PRIPA aim to be compatible with modern regulations in Internet privacy such as the European Union’s General Data Protection Regulation and its United Kingdom counterpart. While it is possible to use PRIPA in a malicious way, the transparency in data collection helps make this more difficult.

**Principle 1: Lawfulness, Fairness and Transparency** According to the first principle of GDPR, a service provider must specify a legal basis in order to collect data. PRIPA only collect data which are specific to an analysis which is agreed to by a participant. Additionally, PRIPA enforces the asking of consent from the user at multiple stages of the analysis, as well as allows a finer-grained control of which datapoints reach the central server. Items 1, 2, 3, 4, 6 and 7 from Table 1 correspond to this principle.

**Principle 2: Purpose limitation** The linguistic analysis is defined before the collection of the data; purpose limitation is built into PRIPA’s core.

**Principle 3: Data minimisation** According to the third principle of GDPR, a service provider must only collect data that is adequate and limited to the claimed purpose of the system. The data to be collected being

<b>P1</b>	Participants are aware of the purpose of the experiment.
<b>P2</b>	Participants are aware of the parameters (web sites, words, time scale) of the data collection.
<b>P3</b>	The features of interest (words, statistical measurements, excerpts) are described in an intelligible way for the participants.
<b>P4</b>	Participants are aware of their right to anonymity.
<b>P5</b>	Participants can consult their data before it is shared with the researchers.
<b>P6</b>	Participants can decide to exclude selected results from the data that is shared with the researchers.
<b>P7</b>	Participants can decide to withdraw completely from a study at any time.
<b>P8</b>	If participants omit to remove personally identifiable information, the researchers should remove it before long-term storage of the data.

Table 1: Key design principles of PRIPA

defined as part of the experiment, data minimisation is another core principle of PRIPA.

**Principle 4: Accuracy** By allowing participants to consult their data and choose which datapoint to communicate to the researchers, and by allowing participants to remove their data post collection, PRIPA allows the information to remain accurate. Items 4, 5, 6, 7, 8 from Table 1 correspond to this principle.

**Principle 5: Storage limitation** The fifth principle of GDPR states that the service provider must not store data for longer than needed for the claimed purpose. This is not enforced in software, but the fact that PRIPA is integrated with the Microsoft Office 365 back-end for storage of results makes it easy to set data storage policies.

**Principle 6: Integrity and confidentiality** By being integrated in the Microsoft Office 365 back-end for data storage, it is easy to enforce a higher level of security and protect whatever personal data was collected.

**Principle 7: Accountability (UK GDPR)** The United Kingdom’s version of GDPR contains a seventh principle: accountability. The principle of accountability requires the service provider to take responsibility of the way personal data is used, and have appropriate measures and records to be able to demonstrate compliance. Much like principle 6, being tied to the Office 365 ecosystem means that existing systems for limiting the use of data and logging access to those datasets can be used out of the box.

## 2.2. Data collection process

PRIPA collects 3 types of linguistic information:

**Word frequencies** Word frequencies are the raw number of occurrences for words in a specific word list, defined as part of the experiment. The word list is specific to the experiment and as such a participant that does not want to share a specific word frequency needs to withdraw from the experiment in order to preserve the integrity of the data without violating their privacy.

**Collocates** Collocates are pairs of words of interest (defined in a word list as part of the experiment) along with their strength of association, given a pre-specified window of words. The list of word pairs is specific to the experiment, and, like word frequencies, a participant that does not want to share a specific word pair needs to withdraw from the experiment.

**Concordance lines** Concordance lines are lines of text showing the context for a particular word, along with the source of that line. The size of the context is specified in the experiment, and the participant can review the list of concordance lines and exclude the ones they do not want to share.

## 2.3. Architecture and design

PRIPA is built in a client-server architecture, where the server hosts experiments which are defined in a specific format using JSON syntax<sup>2</sup>. Figure 1 describes the format. The query allows for six big types of parameters: (1) the title of the study, (2) meta-instructions which apply for the entire experiment and contain details about the way text is meant to be processed (e.g., punctuation, casing, etc.), (3) an allow list which specifies which websites need to be observed

### 2.3.1. Client-side data collection

The client of the application sits in a plug-in for Chromium-based Web browsers (e.g., Google Chrome, Microsoft Edge). We make use of the JavaScript regular expression engine in order to process word lists which are downloaded from the experiment server. Once the user selects an experiment they would like to take part in and accept the disclaimers regarding the way their data will be processed and how they can access/modify/remove it, the PRIPA extension downloads an experiment specification file and watches for the opening or closing of specific websites (depending on the specification of the experiment). When such action (open/close) is triggered, PRIPA attempts to extract the core of the webpage by ignoring banner ads and other informational noise, and runs the analysis based on the word lists provided in the experiment file. The data is stored in the Web browser itself, never leaving the participant's device until they have decided to share their data with the researcher.

---

<sup>2</sup>a lightweight data-interchange format documented at <https://www.json.org>

**Monitoring on tab open/close** Being able to collect data on either the opening or the closing of a tab/window is an important distinction for linguistic analysis. Since some websites dynamically load data based on user input (e.g., Twitter feed, Facebook messages), collecting data at opening would not be effective. Collecting data at close allows for more flexibility in the data collection process by asking participants, for example, to scroll through a month of Twitter feed before closing the tab to start the analysis.

### 2.3.2. Server-side aggregation

The statistical measures collected by PRIPA can be aggregated after the fact. Word frequency can be aggregated with a simple sum, and collocate strength is measured using pointwise mutual information (Bouma, 2009) which can be aggregated using simple frequency measures and information about document length. Considering that the Pointwise Mutual Information of two words  $w_1$  and  $w_2$  in a document  $d$  can be computed as  $PMI(w_1, w_2, d) = \log\left(\frac{P_d(w_1, w_2)}{P_d(w_1) \cdot P_d(w_2)}\right)$  and that  $P_d(w) = \frac{\text{freq}(w)}{|d|}$  where  $|d|$  is the length of document  $d$ , we only need to communicate individual and joint word frequencies as well as length of the web pages in order to aggregate that measure over all participants.

### 2.3.3. Consent monitoring

In order for PRIPA to adhere to the principles laid out in the beginning of the project, consent of the participants needs to be monitored at regular intervals when user data is manipulated. This is done at the following stages:

**During the enrolment stage** The first stage of consent is whether the participant wants to enrol in the experiment.

**During the activation stage** The second stage of consent is whether the participant accepts the collection of data from their device. Participants are asked to explicitly enable the data collection, which will start the monitoring of a specific and explicit set of websites. By explicitly enabling this monitoring, participants are informed that they can disable it at any moment.

**During the review stage** When reviewing concordance lines, participants can choose to exclude specific data points they do not want to share by simply disabling a checkbox, as shown in Figure 2. A number representing the percentage of data censored by the participant is communicated in the results, so that the researcher can make an informed decision about whether to consider this data point.

**At the submission stage** As shown in Figure 3, when submitting results to the researchers participants are asked to consent to the process of sending their data, and can instead opt to stop the experiment and delete their data.

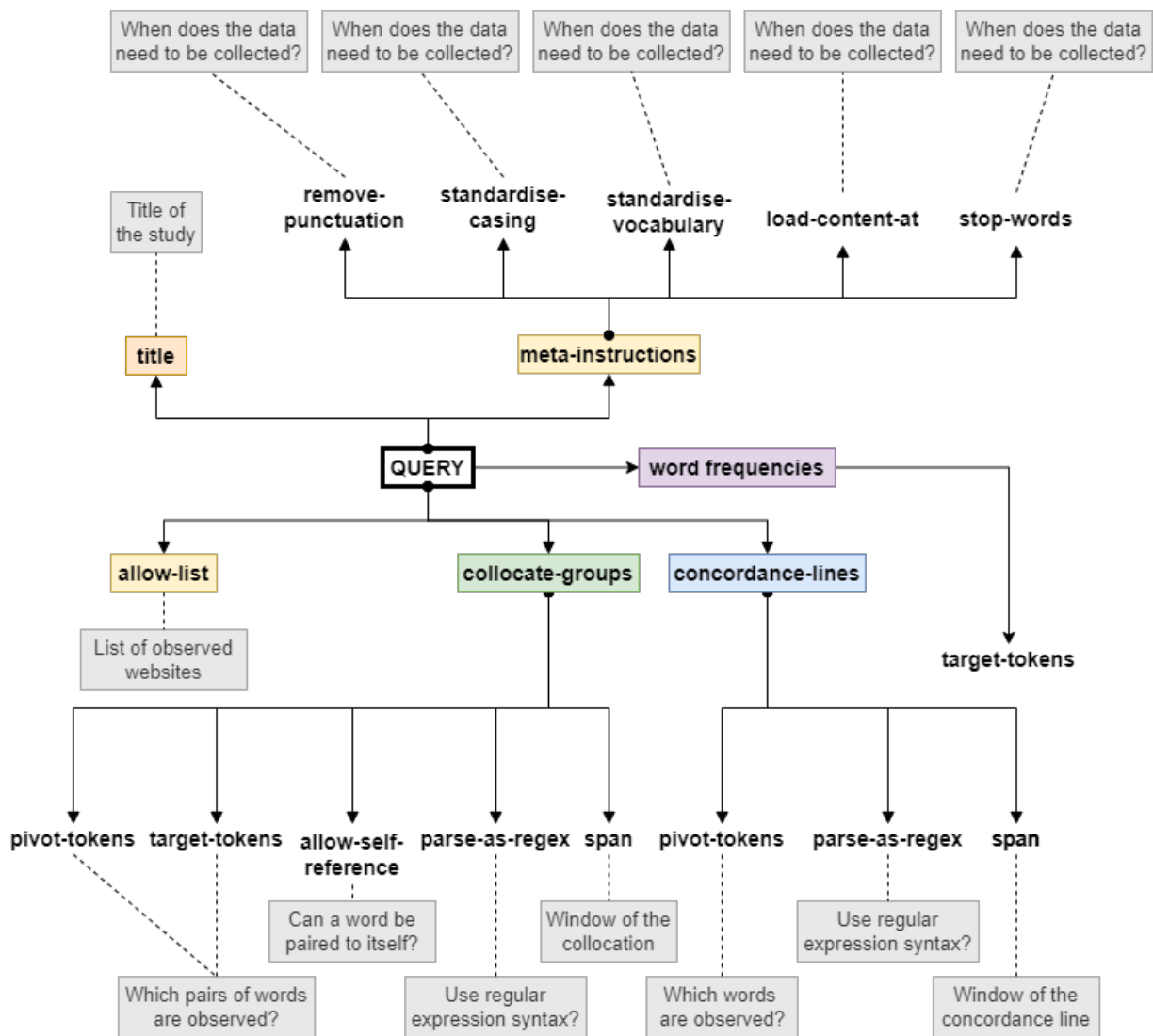


Figure 1: A graphical representation of the experiment file format

See/hide concordance lines ▼

#	Left side	KWIC	Right side	Count	Source	Exclude?
1	play raises serious issues affecting women in London's south Asian	community	but the confused, sitcomy tone is frustratingly like Goodness Gracious	1	www.theguardian.com	<input type="checkbox"/>

Figure 2: Interface allowing participants to remove individual datapoints

### 3. Evaluation and results

We evaluated our system in a three-pronged approach:

**Accuracy of word counting** As pointed out by Anthony (2013), corpus linguistics applications often differ in their measurements due to having different standards in the way they process text. For example,

some software would break "We'll" into two word tokens, while some would keep it as a singular word token. Small variations, repeated over large corpora, can lead to vastly different linguistic measurements and affect interpretation. As such, we calibrated our measurement so that it is close to standard tools such as

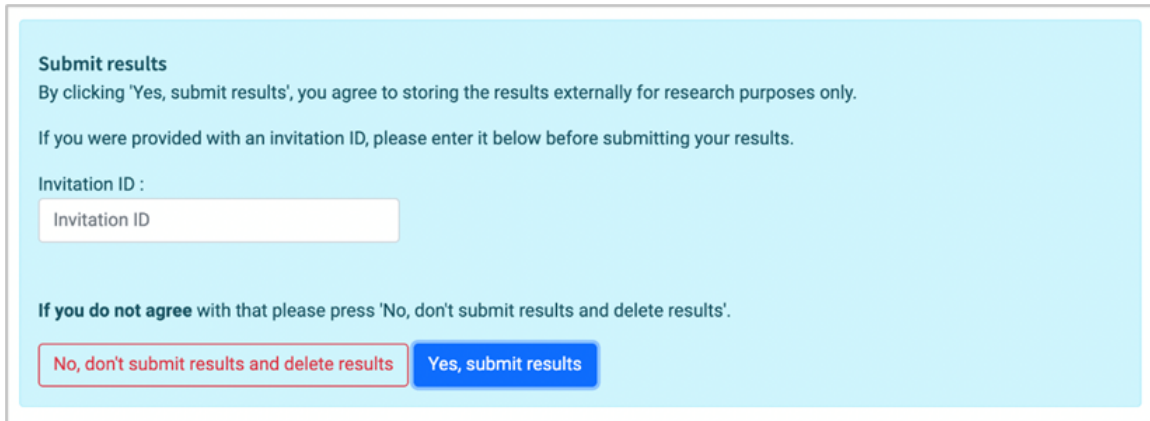


Figure 3: Interface allowing participants to confirm submission of their results, or delete them from the browser

	PRIPA	AntConc	LancsBox
may	33	34	33
might	16	16	15
must	15	15	15
should	29	29	29
would	39	39	39
could	30	30	30
can	93	98	91
will	126	126	125
shall	0	0	0
ought to	0	0	0
<b>total</b>	<b>381</b>	<b>353</b>	<b>377</b>

Table 2: Comparative analysis of PRIPA, AntConc and LancsBox on term frequency of modal verbs on a selected corpus (coloured cells indicate identical counts).

AntConc (Anthony, 2005) and LancsBox (Brezina et al., 2018). We designed a set of test web pages with minimal noise and hosted them on a university website, analysing them both offline with AntConc and LancsBox and online through PRIPA.

In Table 2 we show a comparison of frequencies of single words when running a study on modal verbs on a pre-selected corpus. We can observe that counts mostly match. A visual inspection determined that readings which were not matching were due to tokenisation differences when handling punctuation and apostrophes.

In Figure 4 we show a comparative histogram of the differences between measurements of collocation strength between PRIPA, AntConc, and LancsBox on an experiment measuring collocation strength between modal verbs and pronouns. We can see from this graph that out of our samples, most measurements fell within  $[0, 0.2[$  of LancsBox and  $[0, 0.3[$  of AntConc. A visual inspection showed that the readings that did not match were due to tokenisation differences, like with standard term frequencies.

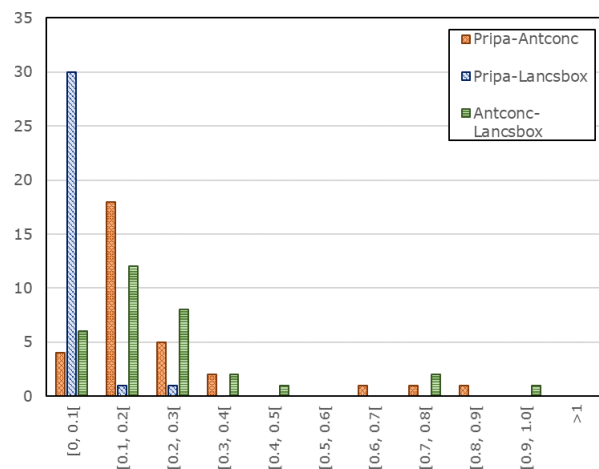


Figure 4: Histogram of differences between PRIPA, AntConc and LancsBox in calculating strength of association between collocates on a sample corpus. Difference between LancsBox and AntConc also provided for baseline.

**Usability of the software** Since participants are rarely researchers themselves, it is important that the software produced is adapted for laypeople and general non-experts. To test this, we ran a usability questionnaire with a small participant involvement panel of 6 people. The quantitative results of the study are summarised in Table 3.

We can see from the data that most participants felt confident in using PRIPA, but had a difficult time understanding the goal of the application. This raises the issue of the importance of a clear user interface and shows that PRIPA can be improved with respect to its first key design principle: participants are aware of the purpose of the experiment. Additionally, we note from the quantitative data reported in Table 3 as well as from qualitative data collected during the same survey that participants were concerned about the privacy of their data. This is partly explained due to the permission model of Chrome-based extensions, which require ask-

Question	Median
Q1 I think that I would like to use this extension frequently.	3
Q2 I found it difficult to understand what the extension does.	4
Q3 I found it easy to set up and run the project in the extension.	4.5
Q4 I think that I would need the support of a technical person to be able to use this extension	1.5
Q5 I found the analyses and results were clearly explained in the extension	2.5
Q6 I felt very confident using this extension	4
Q7 I would imagine that most people would learn to use this extension very quickly	3.5
Q8 I am concerned about the privacy and security of my personal data (i.e., who may be able to access my personal information and how it is protected) when using the extension	2.5

Table 3: Usability questionnaire given to 6 participants - Median value of the Likert data (1 = strongly disagree, 5 = strongly agree).

ing the participants access to their entire browsing experience and them trusting that we will filter only the websites and the data that is stated in the experiment details. Recent updates in the Chrome permission models allow for fine-grained website permissions at runtime and therefore that problem will soon be patched out of PRIPA.

**In-depth study of health communication** In order to evaluate our tool in the field, we ran a study of health communication from the British government during the COVID-19 pandemic. We defined a list of websites of interest based on an empirical study of the most visited news websites in the UK, on which to carry out a pilot study to examine modality markers surrounding key terms from health messages (e.g., "mask", "vaccine", "lockdown", and more). Results from our study shows that PRIPA allows us to access language data from the perspective of the people consuming it. However, it also highlighted a weakness of PRIPA in that when dealing with communication-oriented web applications such as Twitter direct messages or Facebook Messenger, it cannot differentiate between language being produced by the participant and language being consumed. Such information would be useful from a linguistic perspective and will therefore be added in future versions of PRIPA.

#### 4. Conclusion

In this paper we present PRIPA, an early prototype of a new family of corpus linguistics tools that allow for collecting personal data in a privacy-preserving way. PRIPA is an early prototype and therefore a work in progress, but its development raised a number of questions and helped us uncover a set of research directions and good practices for a more trustworthy privacy-preserving type of linguistic analysis.

#### 5. Acknowledgements

This research is funded by the Arts and Humanities Research Council (AHRC), grant reference AH/V015125/1 supported by the Horizon Digital Economy Research Institute. The authors would like

to thank Tino Tom for his work on the initial prototype browser plugin, on which PRIPA is based.

#### 6. References

- Anthony, L. (2005). Antcon: design and development of a freeware corpus analysis toolkit for the technical writing classroom. In *IPCC 2005. Proceedings. International Professional Communication Conference, 2005.*, pages 729–737. IEEE.
- Anthony, L. (2013). A critical look at software tools in corpus linguistics. *Linguistic Research*, 30(2):141–161.
- Bouma, G. (2009). Normalized (pointwise) mutual information in collocation extraction. *Proceedings of GSCL*, 30:31–40.
- Brezina, V., Timperley, M., and McEnery, A. (2018). #lancsbox v. 4. x.
- McCloughlin, E., Nichele, E., Adolphs, S., Barnard, P., Clos, J., Knight, D., McAuley, D., Aydt, M., Tom, T., and Lang, A. (2022). Privacy preserving corpus linguistics: Investigating the trajectories of public health messaging online. Technical report.
- Mozilla. (2022). Mozilla Rally. <https://rally.mozilla.org>. permanent URL hosted at <https://perma.cc/U9TU-N2XV>.
- Noble, A., Cohen, G., Crowcroft, J., Gascón, A., Oswald, M., and Sasse, A. (2019). Protecting Privacy in Pracice: the current use, development and limits of Privacy Enhancing Technologies in data analysis. Technical report, The Royal Society.