

What Does it Mean for a Language Model to Preserve Privacy?

Hannah Brown¹, Katherine Lee², Fatemehsadat Mireshghallah³
 Reza Shokri¹, Florian Tramèr^{4*}

¹National University of Singapore, ²Cornell University

³University of California San Diego, ⁴Google

{hsbrown, reza}@comp.nus.edu.sg kate.lee168@gmail.com
 fatemeh@ucsd.edu tramer@google.com

Abstract

Natural language reflects our private lives and identities, making its privacy concerns as broad as those of real life. Language models lack the ability to understand the context and sensitivity of text, and tend to memorize phrases present in their training sets. An adversary can exploit this tendency to extract training data. Depending on the nature of the content and the context in which this data was collected, this could violate expectations of privacy. Thus, there is a growing interest in techniques for training language models that *preserve privacy*. In this paper, we discuss the mismatch between the narrow assumptions made by popular data protection techniques (data sanitization and differential privacy), and the broadness of natural language and of privacy as a social norm. We argue that existing protection methods cannot guarantee a generic and meaningful notion of privacy for language models. We conclude that language models should be trained on text data which was explicitly produced for public use.

1 Introduction

We use natural language to construct identities and communicate all our information in day-to-day life. Humans naturally understand when sharing a sensitive piece of information is appropriate based on context. It may be fine to share the same piece of information with one specific person or group, and a complete violation of privacy to share in another context, or at another point in time. Between humans, we trust that these implicit boundaries will be recognized and respected. As we build technologies that collect, store, and process our natural language communication, it is important that these technologies do not violate human notions of privacy or make use of data in ways beyond what is needed for the utility of the technology [71, 101].

Language models (LMs) underlie much natural language technology we regularly interact with, from autocorrect to search engines and translation systems. Over the past few years, LMs have grown in size and now utilize unprecedentedly large datasets of natural language making privacy risks in LMs a far reaching problem. Prior work has already demonstrated that such models are prone to memorizing and regurgitating large portions of their training data [12, 13, 51, 38, 91]. Worse, they are especially likely to memorize atypical data points—which are more likely to represent privacy risks for the authors or subjects of these texts.

To address these privacy concerns, there is a growing body of literature that aims to create *privacy-preserving* language models [64, 2, 56, 98, 84, 40, 79]. While humans navigate the complexities of language and privacy by identifying appropriate contexts for sharing information, LMs are not currently designed to do this [14, 72, 66, 49, 66, 50, 41]. Instead, the approach to preserving privacy in LMs has been to *attempt* complete removal of private information from training data (data sanitization), or to design algorithms that do not memorize private data, such as algorithms that satisfy differential privacy (DP) [28, 26].

Both methods make explicit and implicit assumptions about the structure of data to be protected, the nature of private information, and requirements for privacy, that do not hold for the majority of natural language data. Sanitization techniques assume that private information can

*Authors appear in alphabetical order

be formally specified, easily recognized, and efficiently removed. In contrast, the semantic privacy guarantee offered by DP is that an adversary cannot distinguish whether any individual *record* was used to train an LM, which implicitly assumes that these records are well defined and logically map to individual pieces of private information to be protected.

We argue that while these methods can provide some limited form of *data protection* for specific types of text data, they *cannot fully satisfy* the *privacy* expectations that humans endow on the text they share. Data sanitization is only able to recognize a vanishingly small portion of textual private information. In turn, differential privacy can only provide meaningful protection guarantees for information that has clearly defined borders, thereby ignoring the reality that text is inherently a means of *communication*, and that sensitive information is routinely written by or shared among groups of individuals, which blurs the borders of private information. Instead, we argue that an appropriately named “privacy-preserving” LM should guarantee that a user’s data cannot ever appear (or be inferable) outside the context they originally expected it to appear in (i.e., respect *contextual integrity* [71] in the presence of inference attacks)—an ability that cannot be achieved without a deep understanding of the context in which the private information is produced, used, and shared.

Users’ private data is being constantly used to train and fine-tune various services based on language models, which can obviously violate data privacy. Instead, public sources of data (e.g., Web scrapes), seem to not pose privacy risks. Yet, public availability of language data should not be mistaken for data intended to be made public. Text may be shared by humans specifically to violate someone else’s privacy (e.g., doxing), and even public social media posts are not always intended for an audience broader than one’s acquaintances. Even if this is not the case, applications of LMs could make data usable or searchable in new, unintended ways, or make it harder for the data to be modified or erased. An understanding of context is necessary to judge whether it is appropriate to use a piece of data in training.

We further argue that individual users cannot give informed *consent* for their data to be used in a LM or not. First, researchers are still working to quantifying the privacy risks of allowing one’s data to be part of a LM training set. Second, one user’s private information is likely contained in the text of many other users. A single user would not be able to specify how all the text they have contributed is managed. We thus conclude that **data protection is not equivalent to privacy protection for natural language data** and to offer any *meaningful guarantee of privacy*, LMs should be trained on data that was explicitly intended for fully public use, both at present and into the future.

2 Background on Language Models

Language models (LMs) are essential components of state-of-the-art natural language processing pipelines, and refer to systems that are mainly trained on a large corpus of text for word sequence prediction tasks. More precisely, a language model is optimized to learn the occurrence probability of tokens¹ in any sequence, based on the co-occurrence of tokens in the training data. The ultimate objective is to find the relation between a token and its preceding or surrounding segments. To this end, language models extract various statistics and correlations from sequences of words, at the level of sentences or paragraphs.

The current trends of language modeling also shows that aggressive data collection and training enormous models are crucial for improving the performance of LMs. State of the art algorithms based on large neural networks enable effective extraction and encoding of a vast number of statistics about the training corpus, and have achieved unprecedented performance on a wide range of applications. The pervasive application of LMs and ever-larger datasets needed to train them pose serious privacy concerns.

Applications of Language Models. There is a significant interest in the research community and industry to **apply LMs in any situation where humans use natural language** such as: assisting humans in various services, or facilitating communication. For example, LMs are being used in call centers, medical applications, mobile phones and personal computers and home assistants (such as Apple Siri, Amazon Alexa, Google Assistant, Microsoft Cortana, etc), email and message auto-complete services, document translation and search, writing companions (such as SmartCompose [17], Codex and CoPilot for code completion [16]),

¹a token is an instance of a sequence of characters that are grouped together as a useful semantic unit for processing – it could be a character, a word or a sub-word.

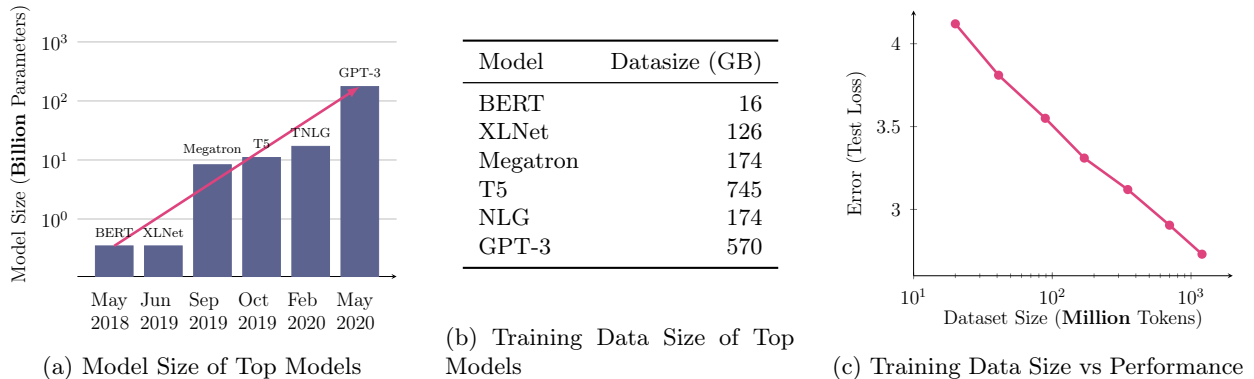


Figure 1: Recent trends in model size and training data size of language models (a and b), and the impact of training set size on model performance (c). **State-of-the-art language models require a significant amount of training data. The size of top models also increases by an order of magnitude every year.** These factors significantly increase the privacy risks of language models.

and many other situations where personal and sensitive data is created and used. The following is a short list of some common language model tasks, which are the foundation of many of LM applications: part of speech (POS) tagging and parsing [25], optical character recognition (OCR) [65], automatic speech recognition [74], natural language generation [11], sentiment analysis [97], and natural language inference [9]. Applications based on these tasks process potentially private data at scale, such as user queries, sensitive documents, emails, and private conversations.

Objectives and Types of Language Modeling. Language models are trained to construct sentences that resemble natural language. They do this by learning statistical measures to capture the local role of each word among its surrounding words and its global consistency within a longer sequence of words (e.g., the whole sentence, or the paragraph).

One core feature of LMs is learning **embedding functions**: mappings from words (and phrases) to vectors in a high-dimensional space such that the closeness between two vectors reflects how close the meanings of the corresponding words (and phrases) are. Embedding functions act as a proxy to encode the semantics of words and sentences in a language and are based on the particular sentences observed in a training corpus. So, training reliable embedding models requires a significant amount of training data. The embedding functions are then used as inputs for downstream NLP tasks. Two other major state-of-the-art classes of neural language models also enable **generating and representing text**: generative LMs which focus on next-token prediction (for example, transformer-based models [94], such as GPT-3 [11]) and masked LMs with the objective of filling in blanks in a sentence (for example, BERT [21] and RoBERTa [58]).

Trends in Language Modeling. Algorithms for learning language models (notably transformer LMs) show an unprecedented performance on extremely large models with hundreds of billions of parameters trained on extremely large datasets [7, 37, 20, 100, 45]. Figure 1 illustrates this trend. What is very important to note is that **using large models, large datasets, and high amounts of compute time are all essential for achieving a high performance** [45]. Empirical results show that the error (test loss) of a transformer-based language model has a power-law relationship to its model size, dataset size, and the amount of compute used for training (see, for example, Figure 1c). Thus, an order of magnitude scale-up is needed to observe tangible improvements in model performance.

3 Privacy Risks of Language Models for their Training Data

Machine learning models learn by extracting generalizable patterns from their training dataset. Yet, it has also been posited that memorizing some parts of the training data can be necessary to optimally generalize

to long-tailed data distributions [30]. For example, nearest neighbor language models [46] which retrieve samples directly from their training dataset are shown to outperform their conventional counterparts. Data memorization can directly lead to leakage of private information from a model’s training set, where behavior of the model on samples that were present in the training set becomes distinguishable from samples that were not. Such leakage has been demonstrated in high-dimensional machine learning models [85], and recent large LMs [13]. The trend appears to get worse as both the size of LMs and their training sets increase (Figure 1). Below we discuss concrete examples of such privacy risks and their consequences.

Membership inference. Membership inference attacks reveal whether or not a given data-point was used in training a given model [85]. These attacks can be seen as privacy risk analysis tools [67], which help reveal how much the model has memorized the individual samples in its training set, and what the risk of individual users is [85, 69, 59, 81, 88]. An adversary who has no direct access to the model and its training, for example in the case of machine-learning-as-a-service, is able to identify the members of the training data by simply querying the model [85]. Membership inference attacks are alarmingly powerful against neural network models with large capacity, enabling them to identify atypical (thus sensitive) members of the training set [69]. The power of membership inference attacks have been demonstrated on natural language processing models such as NLP classifiers [83] as well as released embeddings [87, 61]. Such attacks could cause especially serious harm if they are mounted on clinical models, where membership in the training set could reveal a person’s disease condition.

Training data extraction. Training data extraction refers to the risk of partially extracting training samples by interacting with a trained language model [80, 12, 99, 13]. An adversary can use membership inference attacks as an oracle to generate sentence samples that have a high chance to be in the training set. This attack is demonstrated on the GPT-2 (Generative Pre-trained Transformer) language model family, which consists of three generative models, with different sizes [13]. The attack can successfully recover a person’s full name, address, and phone number from the largest GPT-2 variant (Table 1). The empirical results show that the larger the model is, the more training samples it memorizes: demonstrating once again *the curse of high-dimensionality for data privacy*. Mounting the same type of data extraction attack on BERT-based models trained on de-identified clinical notes shows that more than 4% of generated sentences with a patient’s name also contain one of their true medical conditions [52].

Algorithms behind inference attacks only improve over time. Thus, current attack results under-estimate the privacy risks of large machine learning algorithms, notably language models. Given the privacy risks of LMs, there is an increasing attempt towards designing *privacy-preserving* language models, which can learn the overall distribution and structure of human language, yet do not memorize sensitive information. This can help preserving some notions of privacy, and preventing the out-of-context exposure of training data to unauthorized users.

Existing techniques for building privacy-preserving language models fall into two broad classes: (1) *data sanitization* techniques that find pieces of private information in text and remove these before any further processing, and (2) *differentially private* training algorithms that mitigate the risks of memorization. Section 5 dives deeper into these approaches, and argues that neither is adequate for creating language models that properly preserve users’ privacy.

4 What does preserving privacy in language modeling require?

To claim a language model is privacy preserving, it must only reveal private information (aka “secrets”) in the right contexts and to the right people. While this goal is easy to state, the definition is comprised of three parts, each of which is challenging to determine: (1) in what contexts a secret can be shared without violating privacy (2) what information is contained in the secret, and (3) which people know the secret (the “in-group”).

Far too often, the standard for data protection extends only to not revealing information that harms an individual. Inference attacks, such as those described in Section 3, show the possibility of information leakage in language models. It is not enough to claim privacy is preserved because attacks are not able to

Formatted	Owners	In-group	In-group sharing	Examples
●	1	1	-	Personal password file, secret key
●	1	>1	●	SSN, password, credit card sent to others
●	1	∞	⦿	A developer posts their name, address, and phone number as contact information on Github. Their personal information is “public” on the Web, but in a well defined context.
●	>100	>100	●	A company credit card is shared with employees.
○	1	1	-	Personal search history
○	1	2	●	Bob suffers a mental health crisis and texts a support hotline. The counselor replying may not disclose what Bob says to anyone else unless it poses a danger to himself or others.
○	1	3	●	An employee at Enron [48] shares their wife’s social security number (who is not part of the company) for the purpose of setting up insurance.
○	1-2	>1	○	Alice texts her friends Bob and Charlie about her divorce. Bob further texts Charlie about the matter (c.f. Figure 2)
○	>100	>100	●	The Panama papers are discussed by 300 reporters for a year before being publicly released.

Table 1: Examples of private information, and the contexts in which they might be shared. A piece of private information is “owned” by one or more users (e.g., a credit card that belongs to one user vs. a company credit card that is shared by many). Private data can be shared within a group (the in-group) of variable size. Members of the in-group may be allowed or prohibited from further sharing or discussing the information with other members of the group. Private data can be “formatted” such as a social security number (SSN), or a credit card number, or be referenced in arbitrary prose.

extract information from a model. These attacks improve over time, so while a model that current attacks can extract only a small amount of data from is at low risk for *privacy violation*, this is insufficient for the claim that the model *fully preserves privacy*.

In this section, we illustrate the wide variety of forms private information may take, and how only by understanding context and following privacy norms, can we construct language models that fully preserve privacy. Finally, we discuss how humans approach decisions about when to reveal private information and draw parallels to language models and common privacy defenses. To motivate our main arguments, we provide illustrative examples of different types of personal information shared via natural language in Table 1. These examples cover four axes of variation:

- Some secrets (typically) follow a specific format (e.g., a credit card number), while others are embedded in prose.
- Secrets relate to (or are owned by) a single individual or multiple.
- Secrets are shared with a group (the “in-group”) of one or more individuals.
- Individuals in the in-group may be allowed or prohibited from further sharing or discussing the secret among themselves either implicitly or explicitly (e.g., via legal restrictions).

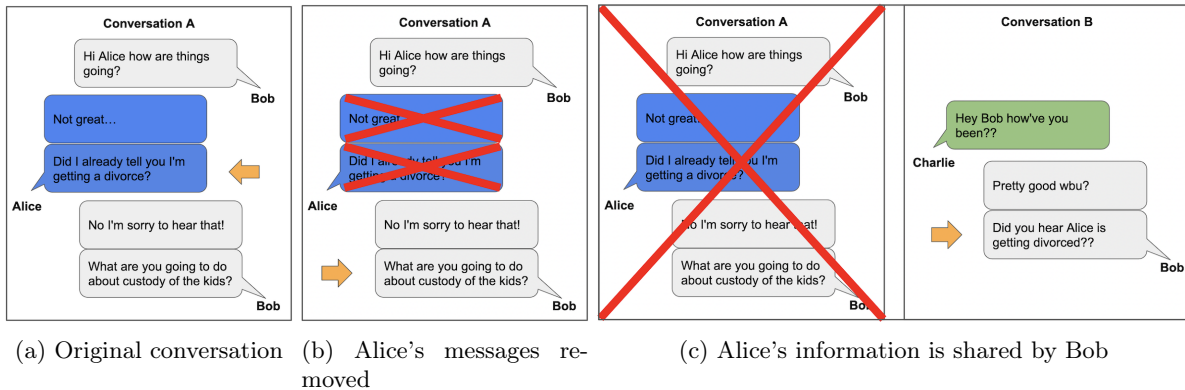


Figure 2: Illustration of the difficulties in removing private information from a dataset. Private information indicated by orange arrows: (a) The original conversation, where Alice shares her private information with Bob. (b) The conversation with all of Alice’s messages removed. Bob’s last message still includes her private information. (c) The whole original conversation is removed. Conversation B still contains Alice’s private information though she is not in the conversation.

4.1 Secrets are contextual

Respecting privacy requires being aware of the context in which information is shared [24, 71]. Instead of simply not “memorizing” private information, humans keep information private through complex judgments of appropriateness dependent on conversational and socio-cultural context. These judgments require information beyond the text of a conversation, making it impossible for an observer, human or computer, to make these same judgments absent this context. Revealing a piece of information to some people may be fine, while it may be a slight violation of privacy to reveal it to a broader audience, and a more severe violation still to make it completely public. These perceptions of privacy are important when considering potentially private information in textual training data. **The scope users mean to share their data in must be considered when deciding whether or not to use it in training.**

Privacy is not a binary variable. Information that is readily shared in one context may be private in another. The counselor texting Bob to help him cope with his mental health crisis (Table 1) may share details about his situation with other professionals or emergency services if there’s reason to believe Bob poses a risk to himself or others, but is otherwise prohibited from sharing what Bob texts. A specific identifying piece of information such as a phone number would be considered sensitive if it belongs to a private individual, and benign if it belongs to a public entity such as a company. More broadly, pieces of data can lie on a spectrum of privacy *levels* with different restrictions and expectations, between the two extremes of fully public (e.g., Wikipedia) or fully private (e.g., someone’s search history).

Training a language model for public use on data that was not intended for that level of publicity violates the original privacy expectations of that data. Nissenbaum’s contextual integrity [71] provides a framework for disambiguating which contexts information can be shared under and with whom. Under contextual integrity, there are five features (the data subject, sender, recipient, information type, and transmission principle) such that if any one is modified, the expectation of privacy could change. In practice, this context could be indicated in the form of social cues and norms, or through regulations (such as Health Insurance Portability and Accountability Act (HIPAA) for medical data) or non-disclosure agreements for corporate information. Under this framework, privacy is considered violated when information is shared outside an acceptable context, which also allows some concept of different degrees of privacy violation[86].

Language models do not understand context. In practice, building machine learning systems that are sufficiently aware of context to appropriately judge the privacy of a piece of information is challenging. Outside of privacy concerns, detecting implied context and reacting to it appropriately has been an active area of research for language models [95]. For example, work has been done to assess how appropriately chatbots respond to delicate situations, such as responding to sexual harassment [14], discussions of suicidal intent or other mental health issues [72, 66], and mentions of violence or physical danger [49, 66]. The results were less than encouraging, with Nobles et al. [72] finding that the majority of the time, chatbots responded in inappropriate ways in these situations, ranging from simply saying "can you repeat that" to giving actively harmful information. Another related line of work is context-aware, long-form, ethical and persona-based, response generation [62, 60, 43, 55], where the chatbot or dialog system is supposed to hold a conversation with previous context taken into consideration [50]. This context could be the persona of the user, or previous conversations. Although this task has advanced in the past few years, the proposed models are plagued by the same challenges as the chatbots and the problem is far from solved [41].

As another example, consider the case of Alice telling Bob about her divorce, as illustrated in Figure 2. Assuming we had a way to recognize that Alice getting a divorce is private information, we could remove Alice’s message to Bob about her divorce. However, this still leaves two more messages referencing the same sensitive information. Bob’s message to Charlie explicitly says that Alice is getting a divorce, and thus obviously refers to the same secret. Bob’s reply to Alice asking about the custody of her kids is more subtle. Understanding that this message is sensitive necessarily requires broader knowledge about the contexts in which asking about custody may occur, and more personalized knowledge about which context most likely applies to Alice. The content of Bob’s message can thus be considered just as sensitive as Alice’s original message, yet automatically identifying the sensitivity of Bob’s message may be much more challenging.

4.2 Secrets are hard to identify

There are many ways of articulating the same point and additional phrasings can be added as language continues to evolve. Private information communicated through language is no different. This can make it difficult to identify whether a piece of text corresponds to private information.

Form and Meaning: There are many ways to communicate any piece of information. Even private information that has an ascribed format (phone numbers, email addresses, credit card numbers, etc.) can appear in multiple forms. For example, numbers and symbols can be spelled out and content can be alluded to, eg: first initial last name at gmail dot com. Synonyms for words can be used, changing the appearance of text but not the meaning. Anaphoric and cataphoric references² present a further challenge to recognizing private information. If Bob were to instead say "Did you hear she’s getting a divorce?" to refer to Alice in Figure 2, that information would still be just as private and involve Alice just as much, but it is harder to automatically recognize.

Format-free pieces of private information, such as those referenced in Table 1 are even more difficult to identify and delineate. Figure 2 shows how drawing a boundary around the text that references a given secret can be difficult. If Alice’s divorce is sensitive information, then should mention of her custody battle be as well? What about future conversations where Alice is referred to as “single”? Imbuing a language model with enough societal context and awareness to recognize these connections appears challenging.

Repeated information can still be private information. Stating the same private information time and again does not make it less private. One example is a company credit card number. This number might be shared again and again within the company, but it remains private to people outside the company. Consider also the case of the Panama papers (Table 1) which contained leaked legal and financial documents detailing how shell companies were created for illegal ends, like tax evasion. Even though 300 journalists exchanged conversation about the Panama papers over the course of a year, the topic of these conversations is no less sensitive. If a language model had been trained on the journalists’ emails, the topic of the investigation (or individual names of suspects or sources) could have been memorized and leaked. De-duplicating the training

²Using a word/phrase to refer to a named entity that is named earlier or later (respectively) in a conversation. In the case of anaphoric references, this entity may also have been defined in a prior interaction.

dataset would not necessarily reduce the likelihood of this information from being learned, as the examples in the dataset are not necessarily near or exact duplicates of each other, but just happen to reference the same broad sensitive topic. Respecting privacy in this case again requires both the ability to recognize the private information and to gather all training data items that refer to this private information.

Language evolves, and so does private information. Changes in language or in social norms can shift the way in which people talk about secrets (or whether something is considered secret or not). For example, the word ‘queer’ was reclaimed by some members of the LGBTQ community beginning in the 1990s as part of the gay rights movement. A system that aims to automatically detect sensitive pieces of text would thus have to be aware of such shifts in linguistic meanings. Yet, changes in language can be swift, in particular on social media [82, 29] where movements, such as #MeToo or Black Lives Matter, can quickly and radically shift the meaning of words and phrases. Additionally, to evade automated censorship and demonetization methods that target specific keywords and phrases, specific topics are routinely re-represented with new words and phrases [47, 39, 90].

Beyond the evolution of language, secrets can also evolve. For example, while much of the content of the Panama papers investigation was highly sensitive and confidential while the investigation was ongoing, the findings were then made public. At the same time, the identities of the journalists’ sources should remain secret essentially forever.

While languages and secrets naturally evolve, language models are typically trained once on a static dataset. Over time, these datasets, and thus the language models trained on them, become less useful for understanding current language. In Section 5, we further explore how the use static datasets can present a challenge for privacy enhancing techniques such as data sanitization (Section 5.1).

4.3 In-groups are hard to identify

Just like finding the borders of a secret is ill-posed, identifying the group of users who are privy to a secret (the in-group) is equally challenging. Indeed, individual text fragments can contain information pertaining to many individuals or organizations at once. The decision of whether to share the secret with a given individual varies from secret to secret, thus the in-group for each secret in different contexts is different. Additionally, just like the secret itself, the in-group can change and grow as relationships continue to evolve in the real world. Thus, even setting a reasonable bound on the *size* of the in-group for each secret can be difficult. As we discuss in Section 5.2, the lack of such a bound poses a particular challenge for articulating meaningful guarantees with differential privacy.

Secrets can involve or be shared among many people. Natural language is *meant to be shared*. We use language to articulate and communicate our thoughts and our observations. At times, these thoughts and observations can also be about other people. Yet, many approaches to data privacy—in particular differential privacy—implicitly or explicitly assume that a user’s private information does not transcend the user’s own data (i.e., the user can protect their privacy simply by not sharing their own data). This assumption can be clearly violated in a variety of natural ways in which humans exchange textual information.

Consider the example described in Table 1 of an employee of a company who sends their wife’s SSN to another employee. We found an example of such an instance in an email from the Enron corpus [48]. While the employee’s wife might “own” her SSN, it now appears in the corpus of text written by the employee. Typically, nothing prevents one user from sharing another person’s private information (such sharing could be legally prohibited, or violate social trust, but these consequences do not mean sharing cannot occur). Thus it can be difficult to define a sole “owner” of a piece of private information.

What we say is often influenced by what others around us have said, which makes drawing dividing lines for privacy, much harder. For example, social media whisper networks, like those discussed in [36] are almost exclusively devoted to sharing private information about people not in the network. In this case, a person outside of the network would still have their private information shared, and to complicate things, the collective information about this person could come from hundreds of different people’s conversation data. Another high-profile example is in the shadow profiles Facebook created of individuals who did not have Facebook accounts. Without any personally volunteered data, Facebook was able to classify enough data

to attribute it to an individual. In Web-scraped datasets [76, 78] that are commonly used in training large language models [22] it is typically not possible to unambiguously map individual pieces of information to specific “owners”.

In-groups have no clear upper-bound. For any individual secret, we could attempt to identify the in-group of people who know the secret. Given such knowledge, we could attempt to remove all mentions of the secret from the entire group. Alternatively, it could be tempting to provide privacy guarantees that are (inversely) proportional to the size of this group (e.g., as in differential privacy), following the intuition that information that has been shared many times is less sensitive than information that has been shared more rarely.

Yet this intuition fails to hold in regard to some of the examples listed in Table 1, and there is no one number k where a piece of information shared with $\geq k$ users can reasonably be assumed to be “non-sensitive”. One individual might share their closest secrets with a handful of friends or family members in a group chat. Others may share the same topics to a broader audience in a support group forum or on their (private) social media page. Companies and governments routinely share sensitive information with hundreds or thousands of employees. And more than 300 journalists communicated in secret for over a year before disclosing their findings in the Panama papers [42]. All this information is definitely private, but within specific contexts is allowed to be shared with a potentially large group of individuals.

4.4 Human notions of privacy

In contrast with common ML privacy preserving mechanisms which focus on preventing models from memorizing private information, humans very clearly memorize sensitive information that they learn. Unlike LMs, we use learned conversational rules to gauge how appropriate or polite something is to share in a given context. One of the simplest proposed sets of rules explaining how we speak—Grice’s Maxims—are a set of four rules (together comprising the Cooperative Principle of Conversation) that describe “normal” conversation [35]. Of these maxims, the ones we use to keep private information to ourselves are “quantity” (say exactly the amount appropriate in a given context) and “relevance” (say only what is relevant to the current context). These maxims are easy to state but *heavily* context dependent, making them difficult to operationalize for technology.

Other conversational frameworks, like politeness theory or relevance theory [10, 89] also rely heavily on context, making their application to NLP systems challenging. At a minimum, these frameworks require prior knowledge about the people involved in the conversation, the socio-cultural context, and past conversations—sometimes with people not involved in the current interaction, who may not have contributed themselves to the same dataset. Given only text data, and none of this further information, it is often impossible to gather all the context necessary to judge if saying something will violate someone’s privacy. Furthering the idea that people memorize and use other methods to preserve privacy, previous work in psychology has shown that we are most likely to remember information that is either very in line with what we have seen before or very *different* from what we’ve seen before [34]. For example, when told a piece of surprising information that we know is supposed to be kept secret, we are likely to remember the information, but choose to not share it.

In summary, humans respect privacy in natural language not by failing to memorize secrets, but by forming a judgment on whether any given piece of information is appropriate, or not, to share with a given party in a given context (unless they share it by mistake, or, by malice, intentionally). Applying a similar approach to language models would require an intrinsic understanding of language and social contexts that goes beyond the capabilities of existing methods, as described in the next section.

5 A Critical Analysis of Privacy Technologies for Language Models

Natural language processing algorithms that aim to respect privacy either remove private information from the data (through text *sanitization* [3, 57, 19, 73]), or design learning algorithms that mitigate the risks of information leakage by not memorizing private information (through *differentially private* learning [28, 15, 1, 64]).

In this section, we evaluate the *claims* of these protection methods about preserving privacy, in the context of language data. Our approach is to lay out the assumptions that data sanitization and differential privacy (DP) make (either implicitly or explicitly). Then, we discuss how awareness of context, difficulty determining the borders of a secret and attributing it to individuals, and other privacy nuances (as extensively discussed in Section 4), can invalidate these assumptions. We discuss the kinds of privacy violations that each method would or would not protect against, and highlight that, given any specific definition for data, **data protection is not equivalent to privacy protection**. They do overlap in many cases where a unit of data contains all the private information about an individual. So, by removing it or not memorizing it (i.e., protecting it from being inferred), we protect the individual’s privacy. However, in general, privacy is much broader than data protection, and this is notably the case in natural language.

5.1 Data sanitization

Data sanitization claims to preserve privacy by removing private information. The critical assumptions are that it is possible to *formally specify private information*, and to *design efficient algorithms to identify and remove private information according to the provided specifications*. We evaluate how realistic these assumptions are, and question if data sanitization can preserve privacy in any meaningful way.

Based on the foundations of privacy in Section 4, we argue that private information expressed in text is difficult to specify and identify, and its removal (according to a given specification) is insufficient to preserve privacy in many situations. Text data can be written in many forms, and the borders of private information are indeterminate. This significantly narrows the application of data sanitization to limited cases where the secret is written according to a context-independent template (e.g., phone number written as consecutive digits).

Sanitization is insufficient because private information is context dependent, not identifiable, and not discrete. Most data sanitization methods are algorithms that use parsers and classification models to tag each word in an input text either based on defined patterns or already tagged data (where sensitive words are manually identified). These techniques work best for identifying well-formatted private information, such as social security numbers, and specific forms of medical note datasets [44, 19, 57, 92]. However, as we discuss in Section 4.2, even well-defined information can be written in many formats or alluded to indirectly. For example, identifying the social security number “*the first 2 digits are two two, and the remaining ones are three ...*” is much more challenging than identifying “223...”. So, even in cases where specifying private information is possible, their reliable identification might be very hard.

Further, identifying and removing non-specific private information, such as the case of Alice’s divorce and custody battle, or the entire discussion around the Panama papers, is significantly more challenging (if not impossible) for data sanitization schemes (which are based on classification models). In general, secrets have no borders, and identifying the scope of relevant information is beyond the capability of taggers and parsers. Besides, understanding sensitivity requires inferring the context, which is a very hard task for algorithms. First of all, there is no formal way to define context, and supervised machine learning models are nonrigorous, empirically inaccurate, and non-explainable methods to classify sensitive information. Secondly, the context related to a piece of text might not be present within the text, which makes understanding the context impossible even for humans. Third, since taggers and parsers require defining ahead of time what the “sensitive” categories are, this limits what information might be related to other sensitive information. Knowing that Alice’s custody battle is sensitive requires understanding that there would be no battle if there were no divorce and requires cultural context (Section 4.1) that is beyond (current) algorithms. Fourth, the context can change after data redaction, which consequently can change the sensitivity of text. So, any claim for data privacy based on sanitization is always outdated.

Changing the context of a piece of information can increase or decrease expectations to privacy. Bob may have a relatively small expectation of privacy when he makes a public social media post, but very high expectations when texting a crisis counselor. In this context, the act of sharing any data at all³ would be considered a privacy violation. Data sanitization completely ignores this, as it assumes information to be

³The prominent text helpline, the Crisis Text Line, recently admitted to doing exactly this for the purposes of helping a for-profit company train language models to improve customer service [53, 75].

discrete and treats privacy as a binary variable. This problem resembles the numerous failed attempts for anonymizing high-dimensional data by removing certain attributes [68, 31]. In the context of language data (with enormous number of dimensions), there is always a possibility of inferring sensitive information even if many pieces of text are redacted. This means that it is possible that either we fail to achieve an acceptable level of privacy through sanitization, or a hypothetically privacy-preserving data sanitization might result in removing almost all the text, rendering it useless: “sanitized data isn’t”.⁴

Data sanitization is useful in very limited cases. We argue that it is not possible to claim privacy using data sanitization algorithms: there is not a specification that would allow private information to be redacted from free-form text data because private text data is not easily identifiable and requires additional context to determine if the information should be redacted. However, data sanitization is a useful obfuscation method in the cases where pieces of context-independent, well-defined, static private information are to be removed from a text dataset.

Data sanitization is currently widely adopted across industries as a data pre-processing step for removing Personally Identifiable Information (PII) or protected health information (PHI) by companies such as Microsoft, Paypal and Mastercard [5, 6, 4, 23, 96, 96, 32] and numerous start-ups (SkyFlow, Ground Labs, PII tools, MailTumble, etc.). Data sanitization can remove some specified information, and can help to reduce the privacy risks to some (unknown) extent. However, it cannot claim that it preserves privacy of individuals, as it has no formal definition for privacy which remains meaningful in the context of language data.

5.2 Differential Privacy

Differential privacy (DP) is a data protection measure designed to assure users that contributing their data to a dataset will not reveal much additional information about the user, when the result of a DP algorithm trained on the dataset is released. Put another way, the data protection guarantee offered by DP is that an adversary cannot easily distinguish whether any individual *record* was used in the computation:

Definition 1. *ϵ -Differential Privacy [28]. For a privacy loss parameter $\epsilon \geq 0$, a training algorithm A satisfies ϵ -DP if and only if for any pair of training datasets D and D' that differ in only one record, and any set of output models S : $\Pr[A(D) \in S] \leq e^\epsilon \Pr[A(D') \in S]$.*

While many applications benefit from this protection, we argue that language data cannot be partitioned to ensure that algorithms trained with DP meet the standard of privacy we put forth in Section 4: *to only emit private information to appropriate people in appropriate contexts*. This is because sensitive language data, as we have seen, cannot necessarily be attributed to one individual or group, whether or not their data is included in the dataset. Thus, while applications of some DP algorithms likely alleviate risks to privacy, they alone are insufficient for *guaranteeing* the absence of privacy violations in language models.

Differential privacy requires a unified definition for secret boundaries, which is very hard if not impossible to achieve for language data. The data protection guarantees of DP hold for any dataset D , and any content of the sensitive record. Thus, compared with data sanitization approaches, DP sidesteps the issue of determining the context or content of private data by providing a worst-case guarantee that applies to any data record. This enables applying DP algorithms in any setting where privacy is considered protected as long as each data record is protected.

However, the main issue with applying DP to language data arises in how we define the *boundaries* of private information. That is, how should we define what constitutes a data “record” in Definition 1. Prior work has considered various granularities, from individual tokens or words, to sentences and documents, or all of a user’s data [54, 64].

Identifying data records with individual words or sentences makes sense from a machine learning perspective, since training batches are often split at such a granularity. But the corresponding privacy guarantees are mostly inadequate since the removal of any individual word or sentence from the training data is insufficient

⁴C.f., “Anonymized data isn’t” – Cynthia Dwork

to hide most types of private information (except maybe a password or SSN that falls inside a single data record). It is thus much more appropriate to define DP with respect to all of a user’s data. Indeed, “user-level” DP is the way in which the original DP definition is intended to be interpreted [28]. In the context of language data, a user-level DP guarantee says that the trained model will be insensitive to the addition or removal of *all the data written by any individual user*. Yet, if we consider the examples in Table 1, it is clear that many types of private information cannot be erased from a dataset by the removal of a single user (even after assuming that a “user” in the system/network is associated with a unique “individual”). Indeed, text is a means of *communicating* information with others. Thus, removing all of a user’s messages is not sufficient to remove the private information from the training set, since others might reference the same information.⁵

Protecting a specific unit of data is not the same as protecting privacy. The issue we are highlighting above is that private information can span across data provided by multiple individuals. It is important to note, however, that the formal guarantees of DP hold regardless of such relationships in the dataset D . What is questioned here is what these data protection guarantees mean, semantically, for the protection of users’ privacy. Differential privacy can protect privacy to the extent that withholding one user’s data from the dataset can. Thus, it is useful for specific types of structured data, for example, when each individual contributes a record that contains sensitive attributes about them (e.g., whether they have been diagnosed with a particular disease). Or alternatively, when a user’s secrets are indeed restricted solely to text written by that user (e.g., an individual’s search history). These protections, however, cannot satisfy the full privacy expectations we discussed in Section 4 regarding natural language data, where private information is not bounded by data records (and can even be about individuals who do not contribute any data), and collections of snippets of text that cover different pieces of private information might overlap. So, withholding any specific unit of data from the dataset cannot *guarantee* protection of privacy.

The need for privacy does not diminish with in-group size. The protection guarantees of DP for groups of users diminish exponentially as the size of the group increases ($k\epsilon$ -DP for groups of size k). However, in practice, the fact that some information is shared among more individuals does not necessarily make it less sensitive. The sensitivity depends on the context and the reasons why the data provided by k individuals contains the same private information. Moreover, appropriately bounding the size k of a group that is privy to a secret is also hard. For example, a community of individuals might share secrets at the level of the whole community. In this case, DP does not provide any strong guarantees for protecting such secrets.

On privacy guarantees and promises. Ideally, we would like to achieve “secret-level” differential privacy, i.e., the algorithm is insensitive to addition or removal of any piece of private information (e.g., Alice’s divorce, or a company’s secret). But satisfying such a definition would require precisely understanding the context and boundaries of secrets, which is exactly a difficulty that DP aims to avoid.

In a typical instantiation of user-level DP, the privacy guarantees provided by the training algorithm are hard to match to the above ideal. While some pieces of information enjoy strong formal protection guarantees from being definitely contained in one user’s data (e.g., text in a user’s personal search history), others are only protected at an exponentially small level (e.g., sensitive information shared among a large group).

This does not mean that information leakage is unbounded. In practice, the provable guarantees offered by DP algorithms are often estimated to be rather loose (i.e., the true leakage is less than what we can mathematically compute) [70, 79, 18]. Yet, the main premise of DP is precisely that it provides provable guarantees, compared to the ad-hoc heuristic guarantees of many other privacy preserving techniques. These strong guarantees have at times been interpreted as a “promise” to users [27], that their secrets will be protected regardless of their decision to share their data. As we have seen however, in the context of

⁵One could expand “user data” to encompass all *conversations* that a user has participated in (including all replies they received), as in Figure 2c. First, satisfying such a level of DP is technically more challenging in decentralized settings (e.g., in Federated Learning [64]) since a data record now spans multiple participants (network users). Second, such an increased granularity remains insufficient to protect knowledge of Alice’s divorce (Table 1) if this secret is further referenced in other conversations (such as between Bob and Charlie in Figure 2c).

language data, this promise loses most of its meaning. We could then ask, whether the formal underpinnings of DP necessarily make it the privacy notion of choice for training LMs, or whether other approaches could provide for more semantically meaningful (albeit possibly only heuristic) forms of privacy protection.

6 Summary and Discussion

Underlying all the challenges of training language models that understand and respect privacy is the complexity of human privacy norms. The vast literature attempting to define privacy and provide frameworks for assessing and understanding it demonstrates the nuance required to disambiguate between what might be similar scenarios. Private information can take many forms, continuously change, and be shared by and among groups whose members fluctuate according to changes in human relationships. In summary, **the boundaries of what data should be acceptable to use for a so-called “privacy preserving” language model are inherently fuzzy and context dependent.**

These challenges limit the applicability of existing techniques like data sanitization and differential privacy (Section 5). Yet, **these privacy-enhancing techniques are often presented as providing users with certain *guarantees* of privacy, which are *not* meaningful enough given the assumptions they make about what constitutes private information.** It is true that, from the perspective of an individual user, the application of any obfuscation technique can only benefit privacy, compared to not applying them (e.g. training a model with DP is better than training the same model without DP). Yet, when applying privacy-preserving techniques to the collection of *new* forms of data (as for training LMs on all types of text produced about every aspect of our lives), we need more realistic and rigorous privacy guarantees.

What alternatives do we have? One might argue that models trained solely on publicly available data, such as text scraped from the Web, alleviate privacy concerns. And indeed, this is the approach taken by many recent large LMs [11, 8, 78, 77]. Yet, publicly accessible does not mean public-intended: publicly shared data typically comes with an intended context of use, which language models could violate by memorizing data [13]. Furthermore, the lack of public discourse and understanding around what happens with collected text data makes informed consent difficult to collect. Ideally, we want LMs to be trained solely on data that is intended (or allowed) for public dissemination. In addition, such LMs could be further fine-tuned or personalized locally on a user’s non-public data, only if the model is going to be used by the same user. But disentangling data that is intended for public use, and obtaining appropriate user consent for its use remains challenging. We discuss these issues in more detail below.

Publicly accessible data is not public-intended. Data that is publicly accessible (e.g., on the Web) is not necessarily intended for unfettered public dissemination, and its use in LMs could still pose privacy risks. For example, publicly available data might not be released by the data subject, such as leaked or subpoenaed email datasets [48, 33], copy/pasting conversations to distribute, or doxing an individual. Posts on social media can also sometimes be made public inadvertently [63, 93]. Furthermore, online text can be deleted or modified. A language model trained on earlier versions of such data would thus inadvertently serve as a data archive. Finally, models trained on Web data might also surface new unintended ways for this public data to be searchable. The example given in Table 1 where an individual posted their contact information on their Github is an actual example of training data extracted from a LM [13].

Can users provide informed consent? Mostly not. Suppose that we asked users to opt-in to having parts or all of their data used out of context to train a language model. For example, one mobile chat client might tell users that it will deploy privacy-preserving LM training on their chat messages, and users and their friends can decide to use the service, or not. Moreover, users might even have the option of flagging individual messages as acceptable for use in training or not. We argue that even if such a consent mechanism were to exist, it would be challenging for users to reach an informed decision about the consequences of their actions.

To start, even experts on ML privacy currently only have a partial understanding of the risks of data memorization and extraction (Section 3), and about how well various defense mechanisms perform. As we argue in Section 5.2, even principled approaches such as differential privacy cannot provide privacy

guarantees that are directly interoperable with privacy *expectations* users might have for their text data. Moreover, individual users cannot properly consent to providing their sensitive information, since they are often not the only person holding that information. As we have illustrated in Section 5.2 and Table 1, sensitive information is routinely shared among many users, all of which would have to jointly consent to release or withhold that piece of data. Put differently, the responsibility to share or hide private information always lies with the entire group that has knowledge of the information. Without understanding how their data will be stored, processed, and disseminated, people are unable to give informed consent.

Private personalization. One approach that we view as a promising middle ground, and worthy of further exploration, is the development of LMs that are trained exclusively on data that is explicitly intended for public use, and further fine-tuned (or personalized) on users’ local (private) data. As long as the model is only used in the local context of the user, the main privacy risks to the user would be alleviated.

7 Conclusions


Our entire life is mediated through language, much of which is monitored and processed by technology. **No discussion of privacy is complete without a deep analysis of how language data is handled.** In this paper, we call for a rigorous understanding of privacy expectations, and for meaningful guarantees of privacy, in the context of language data. We highlight that data protection (with all its limitations) is not equivalent to privacy, existing so-called privacy-preserving methods do not provide reliable assurance about privacy, and users are not in a position to give consent for their data to be used for arbitrary computations. We argue that the only truly privacy preserving solution is to rely exclusively on data that is intended to be public.

8 Acknowledgments

The authors would like to thank David Mimno, Nicholas Carlini, Helen Nissenbaum, Vitaly Shmatikov, Noah Fiedel, Greg Yauney, Maria Antoniak, Federica Bologna, and Martin Strobel for discussions about different parts of this paper.

References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [2] Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar, and Pasin Manurangsi. Large-scale differentially private bert. *arXiv preprint arXiv:2108.01624*, 2021.
- [3] Tuomas Aura, Thomas A Kuhn, and Michael Roe. Scanning electronic documents for personally identifiable information. In *Proceedings of the 5th ACM workshop on Privacy in electronic society*, pages 41–50, 2006.
- [4] Joseph Austin, Shahir Kassam-Adams, Jason A LaBonte, and Paul J Bayless. Self-contained system for de-identifying unstructured data in healthcare records, August 1 2019. US Patent App. 16/255,443.
- [5] Andreas Balzer, David Mowatt, and Muiris Woulfe. Obfuscating information related to personally identifiable information (pii), November 17 2020. US Patent 10,839,104.
- [6] Andreas Balzer, David Mowatt, and Muiris Woulfe. Protecting personally identifiable information (pii) using tagging and persistence of pii, January 5 2021. US Patent 10,885,225.

- [7] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big?  In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 2021.
- [8] Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. <https://doi.org/10.5281/zenodo.5297715>, 2021.
- [9] Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, 2015.
- [10] Penelope Brown, Stephen C Levinson, and Stephen C Levinson. *Politeness: Some universals in language usage*, volume 4. Cambridge university press, 1987.
- [11] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [12] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th {USENIX} Security Symposium ({USENIX} Security 19)*, pages 267–284, 2019.
- [13] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. *arXiv preprint arXiv:2012.07805*, 2020.
- [14] Amanda Cercas Curry and Verena Rieser. #MeToo Alexa: How conversational systems respond to sexual harassment. In *Proceedings of the Second ACL Workshop on Ethics in Natural Language Processing*, pages 7–14, New Orleans, Louisiana, USA, June 2018. Association for Computational Linguistics.
- [15] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- [16] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [17] Mia Xu Chen, Benjamin N Lee, Gagan Bansal, Yuan Cao, Shuyuan Zhang, Justin Lu, Jackie Tsay, Yinan Wang, Andrew M Dai, Zhifeng Chen, et al. Gmail smart compose: Real-time assisted writing. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2287–2295, 2019.
- [18] Rishav Chourasia, Jiayuan Ye, and Reza Shokri. Differential privacy dynamics of langevin diffusion and noisy gradient descent. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [19] Franck Dernoncourt, Ji Young Lee, Ozlem Uzuner, and Peter Szolovits. De-identification of patient notes with recurrent neural networks. *Journal of the American Medical Informatics Association*, 24(3):596–606, 2017.
- [20] Shrey Desai and Greg Durrett. Calibration of pre-trained transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 295–302, 2020.
- [21] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [22] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, and Matt Gardner. Documenting the english colossal clean crawled corpus. *ArXiv*, abs/2104.08758, 2021.

- [23] Jennifer L Donovan, Gary Adler, and James Holladay. Management systems for personal identifying data, and methods relating thereto, January 12 2021. US Patent 10,891,359.
- [24] Paul Dourish. What we talk about when wetalk about context. *Personal and Ubiquitous Computing*, pages 19–39, 2004.
- [25] Greg Durrett and Dan Klein. Neural CRF parsing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 302–312, Beijing, China, July 2015. Association for Computational Linguistics.
- [26] Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*, pages 1–19. Springer, 2008.
- [27] Cynthia Dwork. The promise of differential privacy a tutorial on algorithmic techniques. In *2011 IEEE 52nd Annual Symposium on Foundations of Computer Science, D (Oct. 2011)*, pages 1–2. Citeseer, 2011.
- [28] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [29] Jacob Eisenstein, Brendan O’Connor, Noah A Smith, and Eric P Xing. Diffusion of lexical change in social media. *PloS one*, 9(11):e113114, 2014.
- [30] Vitaly Feldman. Does learning require memorization? a short tale about a long tail. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pages 954–959, 2020.
- [31] Simson Garfinkel, John M Abowd, and Christian Martindale. Understanding database reconstruction attacks on public data. *Communications of the ACM*, 62(3):46–53, 2019.
- [32] Aris Gkoulalas-Divanis, Paul R Bastide, Xu Wang, and Rohit Ranchal. Utility-preserving text de-identification with privacy guarantees, October 28 2021. US Patent App. 16/860,857.
- [33] Margaret Goss. *Temporal News Frames and Judgment: The Hillary Clinton Email Scandal*. PhD thesis, Carnegie Mellon University, 2020.
- [34] Andrea Greve, Elisa Cooper, Roni Tibon, and Richard N. Henson. Knowledge is power: Prior knowledge aids memory for both congruent and incongruent events, but in different ways. *Journal of Experimental Psychology: General*, 148(2):325–341, February 2019.
- [35] H. P. Grice. Logic and conversation. In Peter Cole and Jerry L. Morgan, editors, *Syntax and Semantics: Vol. 3: Speech Acts*, pages 41–58. Academic Press, New York, 1975.
- [36] Bridget Haire, Christy E. Newman, and Bianca Fileborn. Shitty Media Men. In Bianca Fileborn and Rachel Loney-Howes, editors, *#MeToo and the Politics of Social Change*, pages 201–216. Springer International Publishing, Cham, 2019.
- [37] Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Liang Zhang, Wentao Han, Minlie Huang, Qin Jin, Yanyan Lan, Yang Liu, Zhiyuan Liu, Zhiwu Lu, Xipeng Qiu, Ruihua Song, Jie Tang, Ji-Rong Wen, Jinhui Yuan, Wayne Xin Zhao, and Jun Zhu. Pre-trained models: Past, present and future. *AI Open*, 2021.
- [38] Peter Henderson, Koustuv Sinha, Nicolas Angelard-Gontier, Nan Rosemary Ke, Genevieve Fried, Ryan Lowe, and Joelle Pineau. Ethical challenges in data-driven dialogue systems. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 123–129, 2018.
- [39] Chaya Hiruncharoenvate, Zhiyuan Lin, and Eric Gilbert. Algorithmically Bypassing Censorship on Sina Weibo with Nondeterministic Homophone Substitutions. *Proceedings of the International AAAI Conference on Web and Social Media*, 9(1):150–158, 2015. Number: 1.

- [40] Shlomo Hoory, Amir Feder, Avichai Tendler, Alon Cohen, Sofia Erell, Itay Laish, Hootan Nakhost, Uri Stemmer, Ayelet Benjamini, Avinatan Hassidim, et al. Learning and evaluating a differentially private pre-trained language model. In *Proceedings of the Third Workshop on Privacy in Natural Language Processing*, pages 21–29, 2021.
- [41] Dirk Hovy and Diyi Yang. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 588–602, 2021.
- [42] International Consortium of Investigative Journalists. About the Panama Papers investigations. <https://www.icij.org/investigations/panama-papers/pages/panama-papers-about-the-investigation/>, 2016.
- [43] Liwei Jiang, Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Maxwell Forbes, Jon Borchardt, Jenny Liang, Oren Etzioni, Maarten Sap, and Yejin Choi. Delphi: Towards machine ethics and norms. *arXiv preprint arXiv:2110.07574*, 2021.
- [44] Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [45] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [46] Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through memorization: Nearest neighbor language models. In *International Conference on Learning Representations*, 2019.
- [47] Soomin Kim, Changhoon Oh, Won Ik Cho, Donghoon Shin, Bongwon Suh, and Joonhwan Lee. Trkic G00gle: Why and How Users Game Translation Algorithms. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):344:1–344:24, October 2021.
- [48] Bryan Klimt and Yiming Yang. The enron corpus: A new dataset for email classification research. In *European Conference on Machine Learning*, pages 217–226. Springer, 2004.
- [49] Ahmet Baki Kocaballi, Juan C. Quiroz, Dana Rezazadegan, Shlomo Berkovsky, Farah Magrabi, Enrico Coiera, and Liliana Laranjo. Responses of Conversational Agents to Health and Lifestyle Prompts: Investigation of Appropriateness and Presentation Structures. *Journal of Medical Internet Research*, 22(2):e15823, February 2020. Company: Journal of Medical Internet Research Distributor: Journal of Medical Internet Research Institution: Journal of Medical Internet Research Label: Journal of Medical Internet Research Publisher: JMIR Publications Inc., Toronto, Canada.
- [50] Latitude. Ai dungeon, 2019.
- [51] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021.
- [52] Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Goldberg, and Byron C. Wallace. Does bert pretrained on clinical notes reveal sensitive data?, 2021.
- [53] Alexandra S. Levine. Suicide hotline shares data with for-profit spinoff, raising ethical questions, Jan 2022.
- [54] Daniel Levy, Ziteng Sun, Kareem Amin, Satyen Kale, Alex Kulesza, Mehryar Mohri, and Ananda Theertha Suresh. Learning with user-level privacy. *arXiv preprint arXiv:2102.11845*, 2021.

- [55] Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003, Berlin, Germany, August 2016. Association for Computational Linguistics.
- [56] Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large language models can be strong differentially private learners. *arXiv preprint arXiv:2110.05679*, 2021.
- [57] Pierre Lison, Ildikó Pilán, David Sánchez, Montserrat Batet, and Lilja Øvrelid. Anonymisation models for text data: State of the art, challenges and future directions. 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)*, pages 4188–4203, 2021.
- [58] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [59] Yunhui Long, Vincent Bindschaedler, and Carl A. Gunter. Towards measuring membership privacy. *ArXiv*, abs/1712.09136, 2017.
- [60] Junyu Lu, Xiancong Ren, Yazhou Ren, Ao Liu, and Zenglin Xu. Improving contextual language models for response retrieval in multi-turn conversation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '20*, pages 1805–1808, New York, NY, USA, 2020. Association for Computing Machinery.
- [61] Saeed Mahloujifar, Huseyin A Inan, Melissa Chase, Esha Ghosh, and Marcello Hasegawa. Membership inference on word embedding and beyond. *arXiv preprint arXiv:2106.11384*, 2021.
- [62] Bodhisattwa Prasad Majumder, Harsh Jhamtani, Taylor Berg-Kirkpatrick, and Julian McAuley. Like hiking? you probably enjoy nature: Persona-grounded dialog with commonsense expansions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9194–9206, Online, November 2020. Association for Computational Linguistics.
- [63] Alice E Marwick and danah boyd. Networked privacy: How teenagers negotiate context in social media. *New Media & Society*, 16(7):1051–1067, November 2014. Publisher: SAGE Publications.
- [64] H Brendan McMahan, Daniel Ramage, Kunal Talwar, and Li Zhang. Learning differentially private language models without losing accuracy. *arXiv preprint arXiv:1710.06963*, 2017.
- [65] Jamshed Memon, Maira Sami, Rizwan Ahmed Khan, and Mueen Uddin. Handwritten optical character recognition (ocr): A comprehensive systematic literature review (slr). *IEEE Access*, 8:142642–142668, 2020.
- [66] Adam S. Miner, Liliana Laranjo, and A. Baki Kocaballi. Chatbots in the fight against the COVID-19 pandemic. *npj Digital Medicine*, 3(1):1–4, May 2020. Bandiera_abtest: a Cc_license_type: cc_by Cg_type: Nature Research Journals Number: 1 Primary_atype: Comments & Opinion Publisher: Nature Publishing Group Subject_term: Epidemiology;Population screening Subject_term_id: epidemiology;population-screening.
- [67] Sasi Kumar Murakonda and Reza Shokri. Ml privacy meter: Aiding regulatory compliance by quantifying the privacy risks of machine learning. *Workshop on Hot Topics in Privacy Enhancing Technologies (HotPETs)*, 2020.
- [68] Arvind Narayanan and Vitaly Shmatikov. Robust de-anonymization of large sparse datasets. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*, pages 111–125. IEEE, 2008.
- [69] Milad Nasr, Reza Shokri, and Amir Houmansadr. Comprehensive privacy analysis of deep learning: Passive and active white-box inference attacks against centralized and federated learning. In *2019 IEEE symposium on security and privacy (SP)*, pages 739–753. IEEE, 2019.

- [70] Milad Nasr, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, and Nicholas Carlini. Adversary instantiation: Lower bounds for differentially private machine learning. *arXiv preprint arXiv:2101.04535*, 2021.
- [71] Helen Nissenbaum. *Privacy in context*. Stanford University Press, 2009.
- [72] Alicia L. Nobles, Eric C. Leas, Theodore L. Caputi, Shu-Hong Zhu, Steffanie A. Strathdee, and John W. Ayers. Responses to addiction help-seeking from Alexa, Siri, Google Assistant, Cortana, and Bixby intelligent virtual assistants. *npj Digital Medicine*, 3(1):1–3, January 2020. Bandiera_abtest: a Cc_license_type: cc_by Cg_type: Nature Research Journals Number: 1 Primary_atype: Research Publisher: Nature Publishing Group Subject_term: Epidemiology;Rehabilitation Subject_term_id: epidemiology;rehabilitation.
- [73] Beau Norgeot, Kathleen Muenzen, Thomas A Peterson, Xuancheng Fan, Benjamin S Glicksberg, Gundolf Schenk, Eugenia Rutenberg, Boris Oskotsky, Marina Sirota, Jinoos Yazdany, et al. Protected health information filter (philter): accurately and securely de-identifying free-text clinical notes. *NPJ digital medicine*, 3(1):1–8, 2020.
- [74] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le. Specaugment: A simple data augmentation method for automatic speech recognition. *arXiv preprint arXiv:1904.08779*, 2019.
- [75] Keith Porcaro. The real harm of crisis text line’s data sharing, Feb 2022.
- [76] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [77] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-poukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osin-dero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher, 2021.
- [78] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- [79] Swaroop Ramaswamy, Om Thakkar, Rajiv Mathews, Galen Andrew, H Brendan McMahan, and Françoise Beaufays. Training production language models without memorizing user data. *arXiv preprint arXiv:2009.10031*, 2020.
- [80] Ahmed Salem, Apratim Bhattacharya, Michael Backes, Mario Fritz, and Yang Zhang. Updates-Leak: Data set inference and reconstruction attacks in online learning. In *29th USENIX Security Symposium (USENIX Security 20)*, pages 1291–1308. USENIX Association, August 2020.
- [81] Ahmed Salem, Yang Zhang, Mathias Humbert, Mario Fritz, and Michael Backes. ML-leaks: Model and data independent membership inference attacks and defenses on machine learning models. *ArXiv*, abs/1806.01246, 2018.

- [82] Dave Sayers. The mediated innovation model: A framework for researching media influence in language change. *Journal of Sociolinguistics*, 18(2):185–212, 2014.
- [83] Virat Shejwalkar, Huseyin A Inan, Amir Houmansadr, and Robert Sim. Membership inference attacks against nlp classification models. In *NeurIPS 2021 Workshop Privacy in Machine Learning*, 2021.
- [84] Weiyan Shi, Aiqi Cui, Evan Li, Ruoxi Jia, and Zhou Yu. Selective differential privacy for language modeling. *arXiv preprint arXiv:2108.12944*, 2021.
- [85] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 3–18. IEEE, 2017.
- [86] Daniel J. Solove. A taxonomy of privacy. *University of Pennsylvania Law Review*, January 2006.
- [87] Congzheng Song and Ananth Raghunathan. Information leakage in embedding models. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pages 377–390, 2020.
- [88] Congzheng Song and Vitaly Shmatikov. The natural auditor: How to tell if someone used your words to train their model. *ArXiv*, abs/1811.00513, 2018.
- [89] Dan Sperber and Deirdre Wilson. *Relevance: Communication and cognition*, volume 142. Citeseer, 1986.
- [90] Yingnian Tao. Who should apologise: Expressing criticism of public figures on Chinese social media in times of COVID-19. *Discourse & Society*, 32(5):622–638, September 2021. Publisher: SAGE Publications Ltd.
- [91] Om Dipakbhai Thakkar, Swaroop Ramaswamy, Rajiv Mathews, and Françoise Beaufays. Understanding unintended memorization in language models under federated learning. In *Proceedings of the Third Workshop on Privacy in Natural Language Processing*, pages 1–10, 2021.
- [92] Julien Tourille, Matthieu Doutreligne, Olivier Ferret, Aurélie Névéol, Nicolas Paris, and Xavier Tannier. Evaluation of a sequence tagging tool for biomedical texts. In *proceedings of the Ninth International Workshop on Health Text Mining and Information Analysis*, pages 193–203, 2018.
- [93] Sabine Trepte. The Social Media Privacy Model: Privacy and Communication in the Light of Social Media Affordances. *Communication Theory*, 31(4):549–570, November 2021.
- [94] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [95] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from language models, 2021.
- [96] David Williams. Systems and methods for automatically scrubbing sensitive data, April 29 2021. US Patent App. 16/665,959.
- [97] Ashima Yadav and Dinesh Kumar Vishwakarma. Sentiment analysis using deep learning architectures: a review. *Artificial Intelligence Review*, 53(6):4335–4385, 2020.
- [98] Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. Differentially private fine-tuning of language models. *arXiv preprint arXiv:2110.06500*, 2021.

- [99] Santiago Zanella-Béguelin, Lukas Wutschitz, Shruti Tople, Victor Rühle, Andrew Paverd, Olga Ohri-menko, Boris Köpf, and Marc Brockschmidt. *Analyzing Information Leakage of Updates to Natural Language Models*, pages 363–375. Association for Computing Machinery, New York, NY, USA, 2020.
- [100] Ruiqi Zhong, Dhruva Ghosh, Dan Klein, and Jacob Steinhardt. Are larger pretrained language models uniformly better? comparing performance at the instance level. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3813–3827, Online, August 2021. Association for Computational Linguistics.
- [101] Shoshana Zuboff. *The age of surveillance capitalism: The fight for a human future at the new frontier of power*. Profile books, 2019.