

# Data Anonymization in AI and ML Engineering: Balancing Privacy and Model Performance Using Presidio

SURYA GANGADHAR PATCHIPALA

Director, Consulting Expert - Data, AI, ML Engineering, CGI Inc

*Abstract- Data anonymization plays a pivotal role in Artificial Intelligence (AI) & Machine Learning (ML) to ensure individual privacy even as it allows data-led insight. Large datasets in healthcare, finance, and many other industries carry highly sensitive information; consequently, privacy is a sensitive issue. Microsoft's open-source tool, Presidio, detects and eliminates personally identifiable information (PII) in structured and unstructured data. In this article, we examine the tradeoff between privacy and model performance, where Presidio's machine learning-friendly PII anonymization techniques for tokens, masks, and data perturbation allow good data to be protected while still providing utility. Moreover, the article looks at how Presidio helps an organization meet the requirements of privacy laws like GDPR, HIPAA, and CCPA and how it does it responsibly. Data anonymization has key challenges, such as loss of model accuracy and re-identification risks, which are discussed with insight into how Presidio helps mitigate them. Presidio takes this further by utilizing effective data anonymization to allow organizations to develop privacy-compliant, high-performing AI systems that respect privacy and run on data responsibly.*

*Indexed Terms- Data Anonymization, AI Privacy, Machine Learning, Presidio, Personally Identifiable Information (PII)*

## I. INTRODUCTION



Today's data is a critical resource, i.e., the 'new oil', because it can power innovation, shape decisions, and open up new possibilities. Data is behind everything from personalized recommendations in e-commerce to predictive models in healthcare diagnosing diseases. With AI and ML advancing more quickly than ever, organizations need more data than ever to train sophisticated algorithms to automate processes, improve accuracy, and deliver previously unattainable insights.

Yet, as A.I. systems are now permeating our lives in more and more ways (our cars, bank accounts, consumer behavior, and our very speech), we have begun to think about data privacy. Addressing how much personal and sensitive information is being collected escalates fears about how to guarantee the privacy of individuals to the advantage of the great potential of data for A.I. development. In this book, protecting sensitive information is a technical challenge and a moral and regulatory imperative.

Data anonymization is one of the best ways to fix these privacy problems. Data anonymization converts datasets so they cannot be retraceable to certain persons. The data is collected in this way to allow companies and researchers to use it for valuable insights while at the same time maintaining the

privacy of any personally identifiable information, names, social security numbers, or addresses. Take healthcare as an example; patient data is precious for training and building models and sensitive enough to be anonymized for smooth data flow.

However, anonymizing data has its issues for A.I. and ML engineers. One area that has become more and more important in the age of the GDPR in Europe or the California Consumer Privacy Act in the U.S. is anonymization. This term is a made-up one based on the already well-known term anonymize. Achieving high levels of accuracy in our models is often highly dependent on the richness and granularity of our data. The more specific data, the better a machine learning model can learn and the better it is to make predictions. If not done carefully, anonymization can remove information from the data that is useful for modeling degrading performance.

This presents a dilemma: But if we sacrifice performance to protect privacy, how does this data matter in the first place? What tools and policies must organizations require and trust to comply with privacy laws and build accurate and effective machine-learning systems? Presidio, an open-source solution from Microsoft, is where we come in.

This is the exact problem that Presidio was built to solve. Anonymization of text data that contains sensitive information is possible while preserving the integrity and utility of data using the techniques. Presidio is openly available, making it flexible and customizable, serving a broad range of use cases, including those in the financial and health industries and NLP applications. Natural Language Processing techniques can enable presidio to help engineers balance competing needs of privacy and performance by detecting and anonymizing data with PII in text data.

And it most certainly cannot be over emphasized! In the 'leak of the world,' financial and reputational damage is possible, so organizations need to do what they need to do for the sake of their sensitive information. But, simultaneously, they must stay competitive while developing these high-performing A.I. models that can unlock new insights and capabilities. Practical tools such as Presidio remove

the burden by anonymizing data, enabling organizations still to gain value from the data for machine learning purposes.

In this article, we'll look closer at this workhorse of AI and ML engineering – data anonymization – and elaborate on why finding the right balance between privacy and performance is important.

## II. WHAT IS DATA ANONYMIZATION?

Data anonymization is a technique to modify the data so it cannot be identified or retrace oneself to someone. Processing is a critical approach used to protect privacy especially with personally identifying information (PII) including names, addresses, phone numbers, and social security numbers, which allow identification of people. Anonymization achieves this by killing or changing those identifiers to such a degree that we can use the extracted data for analysis, research, or machine learning while protecting the privacy of the people in such data.

With the pressure to protect individuals' privacy mounting, there is an increasing imperative to leverage today's massive amounts of data for more insightful business decisions. For example, in healthcare, finance, or any industry where big amounts of data is gathered about sensitive information, data anonymisation is used to avoid misuse or exposure of that data in a way that leads to privacy foul ups or identity theft.

Here's how data anonymization works:

- Removing PII: The simplest anonymization method involves deleting personally identifiable information from datasets. It might mean removing names, phone numbers, email addresses, and other particular identifiers.
- Altering PII: In many cases, instead of eliminating data, anonymization involves altering the PII so that it may not be returned to a specific person. Similarly, we could replace a name with a random code or token that does not directly refer to the original person.

Also, in machine learning, A.I., an important tool is anonymization; so large datasets are used datasets used to ing the data, and engineers use the data to do

analysis but under privacy laws and ethics. However, performing the process poorly can result in data that isn't useful for the intended purpose, especially if anonymized data is created from more complex datasets.

## 2.1. Principles of Anonymizing Data

### 2.1.1. Irreversibility:

One of the fundamental principles in data anonymization is that after the anonymization process is over, data should no longer be possible from which to reverse the anonymization process and get back the original information. If they remove or replace the original identifiers or sensitive details, nobody can re-identify the individual from the anonymously edited data, not even the cleverest of techniques. Anonymizing data makes it irreversible, so individuals' privacy is not compromised even if it is distributed widely.

### 2.1.2. Data Utility:

Of course, the goal isn't to make the data useless but to protect privacy while preserving as much utility as possible from data. An anonymized dataset can retain (hopefully) enough valuable information to support useful analysis, research, or model training. For instance, anonymized patient data used in medical research means that particular identities may have been removed. However, detecting patterns and trends there should still be possible.

### 2.1.3. Risk Reduction:

One of the critical aspects of the anonymization is reducing the risk of re-identification. Anonymized data, however, can sometimes be connected back to an individual through cross-referencing with another dataset. Efficient anonymization practices aim to reduce this risk by changing the data sufficiently so that virtually no one can reverse engineer it back into anything identifiable (or possibly build their small town).

In reality it's about getting the balance right between good use of the data and protecting privacy. On the one hand, insufficient anonymization can leave individuals re-identified, and on the other hand, too aggressive anonymization can lead to no longer useful data. Thus, it will require craftsmanship, managing the process, and adapting it to how the data will be used.

## III. THE NEED FOR DATA ANONYMIZATION IN A.I. & ML ENGINEERING

This includes anonymizing data during AI and machine learning (ML) engineering for some important reasons, all boiling down to protecting personal privacy, maintaining adherence to laws, and promoting ethical A.I. development. And since data is getting increasingly valuable to train models and make predictions, responsible data management has never been more important.

The most obvious, protection of an individual from their information being viewed is data anonymization. If you are not beating it dead, then many industries such as financial, customer service and healthcare collect heaps of sensitive data. Moreover, this data frequently contains personally identifiable information (PII), like names, medical records, or financial information. However, with anonymization, there is a high chance that this sensitive data will be revealed, resulting in more privacy suspension, identity theft, and unauthorized use of personal data. This information can be used for research, analysis, or training A.I. models while anonymizing data and protecting individuals.

Anonymization of this data is particularly important for industries such as healthcare. The information contained in patient records is very sensitive and can be misused if not anonymized properly. On the other hand, anonymized patient data permits the application of information technology models that diagnose disease, recommend treatment, and improve public health, but with zero privacy issues on behalf of a patient. Given the high stakes in healthcare, anonymizing became a must when using data responsibly.

Data anonymization and compliance with data protection regulations are another extremely important aspect. Governments worldwide have passed laws to protect people's privacy, such as telling organizations how to manage personal data responsibly. The General Data Protection Regulation (GDPR) is the law that sets the rules for how to gather, process, and store personal data on Earth, the territory in Europe. The most effective way by which organization can remove the

privacy is anonymization, as it protects the organization from privacy threats. On the other side, in the United States, the Health Insurance Portability and Accountability Act (HIPAA) regulates how organizations dealing with health care handle personal health information, and anonymization allows them to satisfy these standards without giving up on valuable use of medical information. There is likewise, however, the protection of personal data in California and the control that residents of California will hold over their information, as a case in point, for the California Consumer Privacy Act (CCPA). We make it possible to share data without violating those privacy regulations by anonymizing it.

On top of all this, data anonymization is ethically very important, too. With the rapid march of A.I. technology, so are the fears about ethical data use. Training A.I. and ML models require large datasets from real individuals. An incorrectly anonymized personal information could unintentionally violate an individual's privacy. Suppose an A.I. model is built with identifiable information in the data used. In that case, backlash can be very significant to the point if exposed that private details were used without proper safeguards. This is an ethical use of data, respecting people's privacy and limiting the potential to do damage.

Anonymization is another crucial reason A.I. and ML engineering must anonymize data to reduce the risk of re-identification. Anonymized, the data could still be cross-referenced with other data sets to re-identify individuals. It is even more worrying when you involve large datasets, where even a supposedly anonymous piece of information can be raised suspiciously as yielding a personal identity when interpreted together with other available data. While the risk of re-identification in A.I. and ML must be minimal, a privacy breach can cost a lot when dealing with sensitive data. An effective anonymization technique reduces this risk; the data is still secure and unrelated to certain individuals.

In addition, data anonymization enables data to be shared across research for A.I. Data can be anonymously shared between institutions, research teams, and even borders without violating privacy regulations. In medical research, sharing anonymized

patient data is especially critical, as it helps power next-generation A.I. models to predict disease, diagnose it better, or improve treatments. Anonymization enables organizations to collaborate and innovate in an increasingly data-driven world without exposing sensitive information, creating better, more accurate machine learning models yet preserving privacy.

Also, data anonymization is important for freeing A.I. to be used safely in sensitive industries such as finance and government. These breaches could have a major consequence, affecting highly confidential information. Anonymized data can be used to train A.I. models to detect fraud in the financial sector, and government agencies can mine the anonymized data to improve public services while not compromising the individual's privacy. In high-risk industries, anonymization of these technologies is crucial to guard against data security vulnerabilities when embracing A.I. and ML.

#### IV. INTRODUCTION TO PRESIDIO

Presidio is an advanced, open-source tool developed by Microsoft that addresses a crucial challenge in modern data science: Automating the detection and anonymization of personally identifiable information (PII) within text-based data. As organizations in every industry work with more sensitive information, they must protect their data from compromise while leveraging it to inform A.I. and machine learning techniques. This problem becomes easier with Presidio as we simplify ensuring compliance with privacy regulations without losing access to its utility for analysis.

Unlike many generic anonymization tools, Presidio is highly configurable. Rather than a one-size-fits-all offering, Presidio lets organizations customize its functions to their particular privacy and data processing needs. With lots of variation and complexity in their datasets, businesses working in these areas are very well served by this flexibility, and it's important if you're navigating datasets in this realm to meet compliance with privacy laws like GDPR, HIPAA, or CCPA. However, they can maintain the data quality they are using to build their machine

learning models, natural language processing (NLP) applications, or other data-driven tasks with Presidio.

As such, one of Presidio's core strengths is its capacity to discover PII in unstructured text data — information that typically poses some of the most challenging data to process. Presidio uses Natural Language Processing (NLP) and Machine Learning techniques to scan text data and detect sensitive entities like names, phone numbers, addresses, and dates and more complex entities like credit cards or social security numbers. Presidio then presents several anonymization techniques that may mask the PII to golden value, redact the PII value, substitute the PII value with a token, or encrypt the PII value, leaving the user in control of deciding how to handle the PII value found in the message.

Presidio brings technical capabilities and is built with modern business needs in mind. This was because the language they used was the natural language they intended to speak. Moreover, its open-source nature gave organizations the advantage of extending the tool to their specific operational needs. Presidio provides a thorough and adaptable strategy for securing customer administration logs, therapeutic records, or money-related reports while upgrading quality for A.I. preparation, analysis, or business insight.

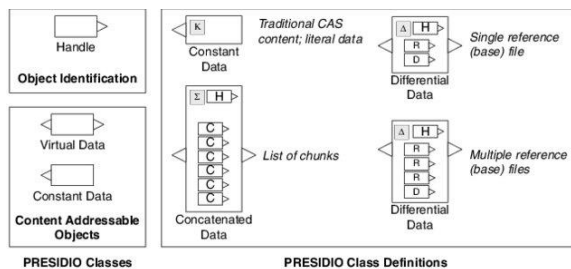


Figure 1. PRESIDIO data classes

## V. HOW PRESIDIO WORKS

Presidio uses natural language processing (NLP) and machine learning algorithms to automatically find and anonymize people associated with Personal Identifiable Information (PII) in text based data. It scans over unstructured text to search for names, social security numbers, credit card numbers and more different types of sensitive info. Once that sensitive data is detected, Presidio uses various anonymization

techniques to protect that sensitive data, still allowing utility and further analysis.

The process by which Presidio works can be broken down into two main stages:

### 5.1. Entity Recognition

The first thing that Presidio does is recognize entities (the sensitive information), where it scans the text for tokenized entities. Using NLP models and pre-trained machine learning algorithms, Presidio can automatically recognize entities such as:

- Names of people
- Phone numbers
- Social Security numbers
- Email addresses
- Credit card details
- Dates and addresses

These entities are also detected using text patterns, where information appears within context. It can recognize various entity types beyond just PII and can be configured to detect any domain-specific PII relevant to an organization's needs. This flexibility makes it useful in healthcare, finance, and customer service industries, where sensitive information is processed at any time.

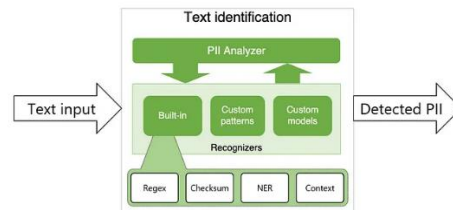


Fig 1. Presidio Analyzer

### 5.2. Data Anonymization

Once PII is identified, Presidio moves into the anonymization phase, where numerous techniques are applied to any removed PII, making it impossible to link it back to an individual. Anonymization methods are specific and can be customized to the organization's privacy and data utility needs. Some of the main techniques include:

- Masking: Information is hidden sensitively by replacing many characters with usual symbols (e.g., turn "John Doe" into "\*\*\*\* \*"). Masking is a popular means of redaction when you need part



of the data (such as a phone number) to be visible and available for processing. Still, the individual's identity needs to be obscured.

- **Tokenization:** This replaces sensitive data with unique (but nonsensitive) tokens ('John Doe' becomes 'User1234'). These can be reversed later if authorized systems need them, but external users will never see the PII.
- **Redaction:** Redaction is the most basic; at its core, redaction is blacking out or deleting the sensitive information from the dataset (e.g., "John Doe" becomes "[REDACTED]"). This is a useful method when there should be no evidence of the original PII.
- **Encryption:** Presidio encrypts PII securely to store it, keeping sensitive data private and keeping PII protected by only allowing authorized users and systems to read it, protecting the data from being readable by unauthorized parties.

With these techniques, however, Presidio can enable its customers to maintain the use of that data for analytics, machine learning, or natural language processing while at the same time complying with privacy regulations like GDPR, HIPAA, or CCPA. Businesses can choose how to anonymize data per the rules applicable to their use cases with the tool's flexibility.

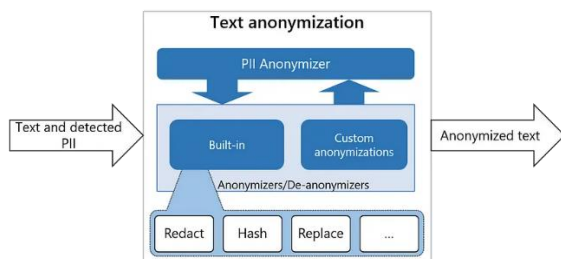


Fig 2. Presidio Anonymization

### 5.3 Customization and Flexibility

Presidio is customizable to let organizations modify what types of PII to detect and how to anonymize it. Thanks to this flexibility, the model fits many industries, including anonymizing patient records in healthcare or customer service chat logs. Further, organizations can train Presidio to identify their own PII of choice if they need to handle sensitive data differently.

## VI. BALANCING PRIVACY AND MODEL PERFORMANCE

The big problem in A.I. and ML is how to strike a proper balance between data privacy and model performance. Data anonymization is a preservation of privacy tool that has been shown to degrade the performance of A.I. models due to removing or obscuring features that may be important for ML processes. A.I. models take advantage of patterns in dataset patterns to learn and deliver accurate results. However, when anonymizing sensitive information, the richness of the data can suffer from a loss of accuracy or effectiveness in the model.

The problem is that anonymization usually involves altering or removing a dataset's personally identifiable information (PII) and other sensitive attributes. While these features are sensitive, they also serve as an important signal for machine learning algorithms. In cases like customer segmentation, fraud detection, personalized recommendations, and so on, anonymizing exact ages, locations, or transaction details might make it hard for the model to make predictions. In doing so, critical details are removed, leaving the model with fewer true patterns to learn from, and if they do exist in the raw data, it will need help understanding the nuanced patterns to achieve higher accuracy and predictive power.

Since this is a problem, Presidio has devised a way to let A.I. engineers customize the anonymization to choose which parts of the data they want to anonymize and how. Unlike traditional anonymization methods, their approach lets organizations choose their anonymization strategy to suit their models and datasets. Such a selective anonymization approach can be essential in maintaining data utility without deriving privacy principles.

Consider that instead of redacting the whole lot, Presidio allows for techniques such as partial Masking or tokenization, which maintain the structures and patterns of the data while not providing the actual details. If a model depends on the shape of email addresses or credit card numbers, for example, using that model to detect fraud, masking the data but leaving the shape intact means the model will continue learning from the data without compromising privacy.

Like sensitive names or I.D.s, tokenizing these items ensures that the relationships are still usable by the machine learning model, but they are no longer identifiable individuals.

This adds flexibility to help minimize the damage anonymization inflicts on model performance. The A.I. models can train efficiently by keeping the structure of the data intact and keeping the important pattern of the data, which is nonsensitive to analysts. All this lets organizations build models that make precise predictions while following privacy laws like GDPR, HIPAA, or CCPA.

Moreover, Presidio provides the ability to experiment with different degrees of anonymization and a smoother interface for trying out the utility privacy tradeoffs of model definitions for training. Different anonymization techniques can be engineered to find the right balance: enough privacy for the model not to be degraded too high. For example, in the case of a customer service dataset, masking the customer name in such a way that certain attributes like age range and location remain would help ensure enough data for a recommendation system to function while at the same time retaining the safety of the customer's identity.

## VII. HOW PRESIDIO ANONYMIZES (AND ANONYMIZED) ITS DATA

Presidio uses several robust techniques for anonymizing sensitive data to guarantee privacy protection and data utility. With PHE, organizations can use these techniques to tailor PII handling methods according to their privacy needs while ensuring data quality for analytics or AI model purposes. Below are the primary methods that Presidio employs for data anonymization:

### 7.1. Tokenization

Tokenization replaces sensitive data with tokens; tokens are unique, nonsensitive placeholders that do not convey direct meanings or values outside the system they are processed for. For example, tokens like "User123" or "Token456" can be used as a replacement (i.e., names, social security numbers, or credit card details). These tokens can only be linked to the original data if by a map table managed securely

so that it is impossible to find and reverse engineer the sensitive information from the token itself.

Given that the relationship of pieces of data needs to be preserved in a machine learning or analytics context, tokenization is especially useful. For instance, an A.I. model that studies customer behavior will hardly need the actual customers' names or I.D.s. Still, customers need to be identifiable and consistently track all their interactions. In this case, tokenization lets the model learn from the data without showing the sensitive PII.

### 7.2. Masking

Another popular Presidio anonymization technique is Masking, which partially obscures sensitive data and hides certain details. Or, for example, if the credit card number is involved, the last four digits may be shown (e.g., "\*\*\*\* \* 1234"), and for a phone number, a part of the number may be masked (e.g., "(555) \*-\*\*").

However, this method is particularly useful when some data structure must be preserved for operational or analytical reasons. Systems can keep the data's original format with Masking, and the sensitive information does not come out completely. In cases where an external user or system requires very limited visibility into data for validation (e.g., a customer service agent might need to verify the last digits of a card without ever seeing the entire number), it's the solution to use.

### 7.3. Data Perturbation

Data perturbation involves adding random "noise" to data to degrade precision while keeping the global statistical properties. In this technique, the amount of information lost when individual-level details are obscured is minimized, while the amount of information remaining useful for analysis is maximized. For instance, a dataset containing individuals' ages might introduce slight variations, such as changing 29 to 30 or 45 to 47, making the age distribution less discernible while remaining on the same distribution.

For this reason, this method is employed in statistical analysis and engine learning tasks where the structure of the entire set of data has to be kept without taking

care of the exact data. However, using perturbation, data scientists can still run accurate analyses on the data with the key information anonymized, so long as personal data is protected well enough.

#### 7.4. Encryption

Presidio supports the encryption of highly sensitive data, converting readable data to an encoded format accessible only through appropriate decryption keys by authorized users. Data that is required to be stored or transmitted safely often undergoes encryption. If unauthorized users find their way into encrypted data, they can only read or use it once they have the decryption key.

For companies that process enormous amounts of information, encryption can be especially useful; financial institutions, healthcare providers, or the government, for example, often have reason to keep data confidential. Its end goal is to help secure sensitive data, for example, credit card numbers, clinical records, or close-to-home identifiers, aside from being available to authorized mortal under tight security protocols.

### VIII. ANONYMIZED DATA FOR BETTER MODEL PERFORMANCE

However, anonymizing data is crucial for privacy protection and can destroy key data features that hold the predictive power for an A.I. model. However, engineers can use multiple strategies that guarantee high model performance while using anonymized datasets. Feature engineering is one such strategy, where you pick, transform, and create new features out of the data so that they still retain some meaningful information even after you have anonymized the data. By funneling the learning of the models down to the most relevant features and keeping the data patterns intact, engineers can ensure that the models learn without using sensitive information.

Data augmentation is another important technique to expand the dataset by generating little perturbed copies of the existing data points. That can redress the blow to the quality of detail lost due to anonymization by providing models with more examples to learn from. You engineer anonymized datasets and augment

them to train processes better so the models can generalize better based on the available data.

In addition, we want to maintain the structure and distribution of data when anonymizing it. Sensitive information is removed or masked so that the model can learn effectively, even if sensitive information is not present because removing or Masking does not sufficiently remove information such that the overall patterns in the data are maintained. For example, suppose anonymization changes an individual's age. In that case, the resulting distribution of ages across the population in the dataset must be maintained for training tasks where demographic information makes a difference, such as healthcare outcomes prediction or customer segmentation.

With well-considered features of engineering, data augmentation, and data distribution, engineers can cope with anonymization and keep their machine learning models machine-learning models accurate and reliable.

#### Presidio and Handling Different Data Types

Presidio can anonymize structured and unstructured data and is flexible in how you want it to handle sensitive data across data types.

Presidio can anonymize structured data, such as tables or databases, and detect sensitive information like social security numbers, addresses, phone numbers, or credit card details. Each data point is neatly structured in rows and columns in such a dataset, making it convenient for Presidio to scan specific fields to look for personally identifiable information (PII) and to apply the right anonymization techniques like tokenization, Masking, or encryption. It is very useful in various industries like finance, where huge databases of transactional or account details of customers need to be anonymized without losing the analytical value of the dataset.

When it comes to unstructured data, such as documents, emails, and free-form notes, Presidio's Natural Language Processing (NLP) technologies apply their strengths. Processing and anonymizing unstructured data is difficult. Advanced NLP models in Presidio detect sensitive entities in natural language like names, addresses, dates, or any other context-



dependent PII. Presidio allows us to anonymize text-based data in environments like customer service with logs of different conversations, all of which contain proprietary customer information, and we need to remove that before training our A.I. models.

Presidio supports structured and unstructured data to anonymize the widest possible range of datasets, from transactional databases to complex text documents, and adhere to privacy regulations while preserving individual privacy. Companies can continue to use their data for machine learning, analytics, or business intelligence, regardless of the form of the original dataset, thanks to the flexibility of this offering.

## IX. HOW TO USE PRESIDIO IN A.I. AND ML USE CASES

Because it is so powerful, it is no surprise that Presidio has become a widely adopted tool across numerous industries, which allows it to anonymize such sensitive data while continuing to allow the use of that data for AI and machine learning (ML) purposes. Presidio guarantees data privacy and ensures that organizations can fulfill the requirements based on regulations while the data is fully leveraged to its potential. Below are some key use cases where Presidio has been particularly effective:

### 9.1. Healthcare

The healthcare industry must protect patient privacy with the use of sensitive, medical records and personal health information. HIPAA forces Health Care Providers to remove identifiable patient data from the data they wish to share or analyze, so that this data cannot create identifiable links which may breach privacy or violate the privacy of a patient. Anonymizing these records, Presidio helps healthcare organizations remove sensitive data like patient names, medical record numbers, and addresses while also detecting this activity.

Presidio enables healthcare providers and research institutions to continue to harness anonymized patient data to power machine-learning models, enabling medical research, diagnosis, and treatment. For instance, by anonymizing medical data such as data sharing with the AI, training a model that can predict a disease outbreak, develop personalized treatment

plans, or detect how well a specific type of medication works. This is needed to ensure patient confidentiality is not violated while advancing healthcare.

### 9.2. Finance

Since public money is involved, financial services and other sectors deal with huge amounts of sensitive financial data, including bank account numbers, credit card details, and transaction histories. With regulations like GDPR and CCPA, which are centered on protecting personal financial information, banks, and financial institutions must ensure that this data is properly anonymized before they can use it to analyze it or share it with other parties.

As a leading provider of data anonymization and sensitive data masking, Presidio is an integral part of the ADVA platform, helping financial institutions anonymize sensitive data while still allowing critical AI-driven applications such as fraud detection and risk assessment to run. Presidio hides sensitive details, like a credit card number, by using Masking or tokenization to prevent A.I. models from viewing data and allowing them to find fraudulent transactions or assess customer credit risk without revealing a person's financial information. Such models help banks and financial institutions maintain privacy compliance using sophisticated machine learning models that protect the institution and its valued customers.

### 9.3. Customer Service and NLP

Businesses that do customer service work rely on the power of NLP models to enhance customer experience and automate service with chatbots or virtual assistants. However, customer service data, like chat logs or email conversations, contains sensitive information like names, addresses, and numbers. That information needs to be cleansed through anonymization to meet privacy regulations.

Anonymizing conversations and textual data with their NLP makes Presidio extremely good at helping brands train customer service models without exposing personal details. Take chatbots operated within banking and e-commerce, for instance: they can be trained on anonymized conversation to learn how to answer customer queries better or predict customer needs while never revealing any sensitive customer

data. Protecting user privacy improves the quality of AI-driven customer support.

#### X. PRESIDIO COMPLIANCE & PRIVACY REGULATIONS

This is what makes Presidio so specific – it is designed to enable organizations to comply with rigorous privacy regulations such as the General Data Protection Regulation (GDPR) in Europe, the Health Insurance Portability and Accountability Act (HIPAA) in the United States, and the more recent California Consumer Privacy Act (CCPA). A consequence of these regulations is that all personally identifiable information (PII) is limited by regulation on how it is collected, processed, and stored, and organizations have to take measures to protect individual privacy. There can be hefty fines and reputational damage for noncompliance.

Presidio gives organizations great control over how data-sensitive data is detected, anonymized, and handled in a way consistent with privacy policies without sacrificing data value. For instance, GDPR stresses the necessity of data minimization and pseudonymization, which reflects the only type of information collected so that it can be processed about a particular aim and the obligation to remove or prevent the recognition of personal data when possible. Presidio automatically identifies and anonymizes PII across rows (patient demographics, for example) and columns (email fields, fax numbers, etc.), which include text-based datasets (emails, medical records, etc.).

Presidio complies with HIPAA, the law that protects how healthcare data is used in the U.S. so that health data can be anonymized to keep patient data protected even while research, analysis, or A.I. model training takes place. HIPAA's privacy rule is that healthcare providers must remain anonymous with patient names, medical I.D.s, and contact details; however, anonymization of data allows healthcare providers to use data to improve healthcare outcomes further.

In addition to rights to know about, and the right to have, personal data deleted, the CCPA provides additional rights for consumers to learn the extent to which data about them is being collected – and the

CCPA also includes de-limitation rights. By proving that CCPA's requirements can be enforced without sacrificing customer data protection, Presidio can make this happen. Presidio's dishonest anonymization feature reduces the chance of exposing PII in these cases and helps businesses utilize anonymized data for knowledge and machine learning while honoring the buyer's rights.

At its core, Presidio lets organizations concretely and usefully meet compliance standards to handle sensitive data safely and within the needed data quality for AI and machine learning tasks.

##### 10.1. Limitations of Data Anonymization in A.I.

Presidio is a terrific tool for data anonymization. Still, the process has some built-in limits regarding machine learning, and A.I. One big limitation is that anonymization can sometimes reduce the contextual richness of the data, which is why it's less useful for A.I. models. Most machine learning algorithms require detailed feature-rich datasets for learning general patterns, making predictions, and improving over time. Anonymizing or removing sensitive information may cast some of the key signals upon which the data rely and thus may reduce the model's accuracy or predictive power.

Taking the healthcare example further, anonymizing patient information like exact age and patient medical histories could lead to losing some important patterns required by disease outcome prediction or treatment models. Anonymizing transaction details alter the accuracy of fraud detection models, as their behavior patterns might not help identify the frauds.

Additionally, data can always be re-identified, and anonymized data could be de-anonymized by merging it with other datasets (s). Attacks using sophisticated techniques, such as linking anonymized data with public records (or other datasets), could cause someone to reverse the anonymization and identify an individual. That part is particularly ominous if anonymized data sets have unique combinations of attributes that, when cross-checked, identify a person. For example, even if a dataset no longer holds names or addresses, a unique demographic combo (such as gender, age, and ZIP code) can re-identify a dataset.

Differential privacy and anonymization are also another limitation. Differential privacy techniques that add noise to datasets and reduce the risk of re-identification exist, but implementing these techniques can be very hard and often decreases the overall utility of the data. Too much noise can be added to protect privacy, making the dataset less helpful for training machine learning models, thereby losing precision and effectiveness.

## CONCLUSION

Privacy protection requires data anonymization, which is even more important in the rapidly growing fields of AI and machine learning (ML), where huge amounts of personal and sensitive data are processed daily. Yet, equally as important is the problem of ensuring that anonymized data still has enough utility to train high-performing models. However, when performed casually, the anonymization process can remove valuable information, resulting in a degradation of model accuracy and effectiveness.

Presidio addresses this challenge head-on with a powerful anonymization solution that allows organizations to retain sensitive information while maintaining high data quality in line with AI and ML demands. Balancing the tradeoff of privacy vs. the need for good quality data allows engineers to be flexible with handling both structured and unstructured data and have customizable anonymization techniques. In healthcare, this anonymizes patient records – in banking, it handles financial transactions; in service industries, customer interactions must adhere to important regulations like GDPR, HIPAA, and CCPA.

As data privacy is becoming more important than ever, organizations looking to develop ethical, responsible, and privacy-compliant A.I. systems must use tools like Presidio. By using smart anonymization, businesses can keep innovating with data-driven insight while protecting people's privacy and meeting the needs of evolving data protection laws.

## REFERENCES

[1] Chen, L. (2022, October 7). PII anonymization made easy by Presidio - Towards Data Science.

Medium.  
<https://towardsdatascience.com/building-a-customized-pii-anonymizer-with-microsoft-presidio-b5c2ddfe523b>

- [2] Benchmarking Advanced Text Anonymisation Methods: A comparative study on novel and traditional approaches. (n.d.). Retrieved from <https://arxiv.org/html/2404.14465v1>
- [3] Latanya Sweeney. “k-anonymity: a model for protecting privacy”. In: *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10.5 (Oct. 1, 2002), pp. 557–570. ISSN: 0218-4885. DOI: 10.1142/S0218488502001648. URL: <https://doi.org/10.1142/S0218488502001648> (visited on 11/21/2022) (cit. on p. 3).
- [4] Vpnreports. (2022b, July 18). Is Anonymous Data Really Anonymous? VPN Reports. <https://www.vpnreports.com/is-anonymous-data-really-anonymous/>
- [5] A. Machanavajjhala et al. “L-diversity: privacy beyond k-anonymity”. In: *22nd International Conference on Data Engineering (ICDE’06). 22nd International Conference on Data Engineering (ICDE’06)*. ISSN: 2375-026X. Apr. 2006, pp. 24–24. DOI: 10.1109/ICDE.2006.1 (cit. on p. 3).
- [6] Ninghui Li, Tiancheng Li, and Suresh Venkatasubramanian. “t-Closeness: Privacy Beyond k-Anonymity and l-Diversity”. In: *2007 IEEE 23rd International Conference on Data Engineering. 2007 IEEE 23rd International Conference on Data Engineering*. ISSN: 2375-026X. Apr. 2007, pp. 106–115. DOI: 10.1109/ICDE.2007.367856 (cit. on p. 3).
- [7] Cynthia Dwork. “Differential Privacy: A Survey of Results”. In: *Theory and Applications of Models of Computation*. Ed. by Manindra Agrawal et al. *Lecture Notes in Computer Science*. Berlin, Heidelberg: Springer, 2008, pp. 1–19. ISBN: 978-3-540-79228-4. DOI: 10.1007/978-3-540-79228-4\_1 (cit. on p. 3).
- [8] Brijesh Mehta et al. “Towards privacy preserving unstructured big data publishing”. In: *Journal of Intelligent & Fuzzy Systems* 36.4 (Jan. 1, 2019). Publisher: IOS Press, pp. 3471–3482. ISSN: 1064-1246. DOI: 10.3233/JIFS-181231. URL:

- <https://content.iospress.com/articles/journal-of-intelligent-and-fuzzy-systems/ifs181231> (visited on 11/04/2022) (cit. on p. 3).
- [9] Art. 4 GDPR – Definitions. General Data Protection Regulation (GDPR). URL: <https://gdpr-info.eu/art-4-gdpr/> (visited on 11/04/2022) (cit. on p. 5).
- [10] Batet, Montserrat and David Sánchez. 2018. Semantic disclosure control: Semantics meets data privacy. *Online Information Review*, 42(3):290–303. <https://doi.org/10.1108/OIR-03-2017-0090>
- [11] Batet, Montserrat and David Sánchez. 2020. Leveraging synonymy and polysemy to improve semantic similarity assessments based on intrinsic information content. *Artificial Intelligence Review*, 53(3):2023–2041. <https://doi.org/10.1007/s10462-019-09725-4>
- [12] Mendels, Omri. 2020. Custom NLP approaches to data anonymization. *Towards Data Science*. <https://towardsdatascience.com/nlp-approaches-to-data-anonymization-1fb5bde6b929>. Accessed: 2022-06-06.
- [13] Lison, P., Pilán, I., Sánchez, D., Batet, M., & Øvrelid, L. (2021). Anonymisation Models for Text Data: State of the Art, Challenges and Future Directions. In *Association for Computational Linguistics & International Joint Conference on Natural Language Processing, Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing* (pp. 4188–4203). <https://aclanthology.org/2021.acl-long.323.pdf>
- [14] You, L.L. & Pollack, K.T. & Long, Darrell. (2005). Deep Store: an archival storage system architecture. *Proceedings - International Conference on Data Engineering*. 804- 815. 10.1109/ICDE.2005.47.
- [15] Vpnreports. (2022, July 18). Is Anonymous Data Really Anonymous? *VPN Reports*. <https://www.vpnreports.com/is-anonymous-data-really-anonymous/>
- [16] Pierre Lison, Ildikó Pilán, David Sánchez, Montserrat Batet, and Lilja Øvrelid, Anonymisation Models for Text Data: State of the Art, Challenges and Future Directions (2021). *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*
- [17] Pierangela Samarati and Latanya Sweeney, Protecting Privacy when Disclosing Information: k-Anonymity and its Enforcement through Generalization and Suppression (1998). Technical report, SRI International
- [18] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith, Calibrating Noise to Sensitivity in Private Data Analysis (2006). *Theory of Cryptography*
- [19] Franck Deroncourt, Ji Young Lee, Ozlem Uzuner, and Peter Szolovits, De-identification of patient notes with recurrent neural networks (2017). *Journal of the American Medical Informatics Association*
- [20] Alistair EW Johnson, Lucas Bulgarelli, and Tom J Pollard, De-identification of free-text medical records using pre-trained bidirectional transformers (2020). *Proceedings of the ACM Conference on Health, Inference, and Learning*
- [21] Y. Pei, Y. Liu, N. Ling, Y. Ren and L. Liu, "An End-to-End Deep Generative Network for Low Bitrate Image Coding," 2023 IEEE International Symposium on Circuits and Systems (ISCAS), Monterey, CA, USA, 2023, pp. 1-5, doi: 10.1109/ISCAS46773.2023.10182028.
- [22] Zhu, Y. (2023). Beyond Labels: A Comprehensive Review of Self-Supervised Learning and Intrinsic Data Properties. *Journal of Science & Technology*, 4(4), 65-84.
- [23] B. Edwards, "Openai introduces gpt-4 turbo: Larger memory, lower cost, new knowledge," *Ars Technica*, November 2023. [Online]. Available: <https://arstechnica.com/information-technology/2023/11/openai-introduces-gpt-4-turbo-larger-memory-lower-cost-new-knowledge/>
- [24] R. Khawaja, "Best large language models (llms) in 2024," may 2024. [Online]. Available: <https://datasciencedojo.com/blog/best-large-language-models/>
- [25] H. N. B, "Confusion matrix, accuracy, precision, recall, f1 score," *Medium*, 2019. [Online].

- Available: <https://medium.com/analytics-vidhya/confusion-matrix-accuracy-precision-recall-f1-score-ade299cf63cd>
- [26] D. R. Almeida, "Synthetic data generation (part 1)," Apr 2024. [Online]. Available: <https://cookbook.openai.com/examples/sdgl>
- [27] "Managing environments," 2017. [Online]. Available: <https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-environments.html>
- [28] A. Stam and B. Kleiner, "Data anonymization: Legal, ethical, and strategic considerations," FORS Guide No. 11, Version 1.1, Lausanne, 2020, last update January 2022.
- [29] Krishna, K. (2022). Optimizing query performance in distributed NoSQL databases through adaptive indexing and data partitioning techniques. International Journal of Creative Research Thoughts (IJCRT). <https://ijcrt.org/viewfulltext.php>.
- [30] Krishna, K., & Thakur, D. (2021). Automated Machine Learning (AutoML) for Real-Time Data Streams: Challenges and Innovations in Online Learning Algorithms. Journal of Emerging Technologies and Innovative Research (JETIR), 8(12).
- [31] Murthy, P., & Thakur, D. (2022). Cross-Layer Optimization Techniques for Enhancing Consistency and Performance in Distributed NoSQL Database. International Journal of Enhanced Research in Management & Computer Applications, 35.
- [32] Murthy, P., & Mehra, A. (2021). Exploring Neuromorphic Computing for Ultra-Low Latency Transaction Processing in Edge Database Architectures. Journal of Emerging Technologies and Innovative Research, 8(1), 25-26.
- [33] Mehra, A. (2024). HYBRID AI MODELS: INTEGRATING SYMBOLIC REASONING WITH DEEP LEARNING FOR COMPLEX DECISION-MAKING. Journal of Emerging Technologies and Innovative Research (JETIR), Journal of Emerging Technologies and Innovative Research (JETIR), 11(8), f693-f695.
- [34] Thakur, D. (2021). Federated Learning and Privacy-Preserving AI: Challenges and Solutions in Distributed Machine Learning. International Journal of All Research Education and Scientific Methods (IJARESM), 9(6), 3763-3764.