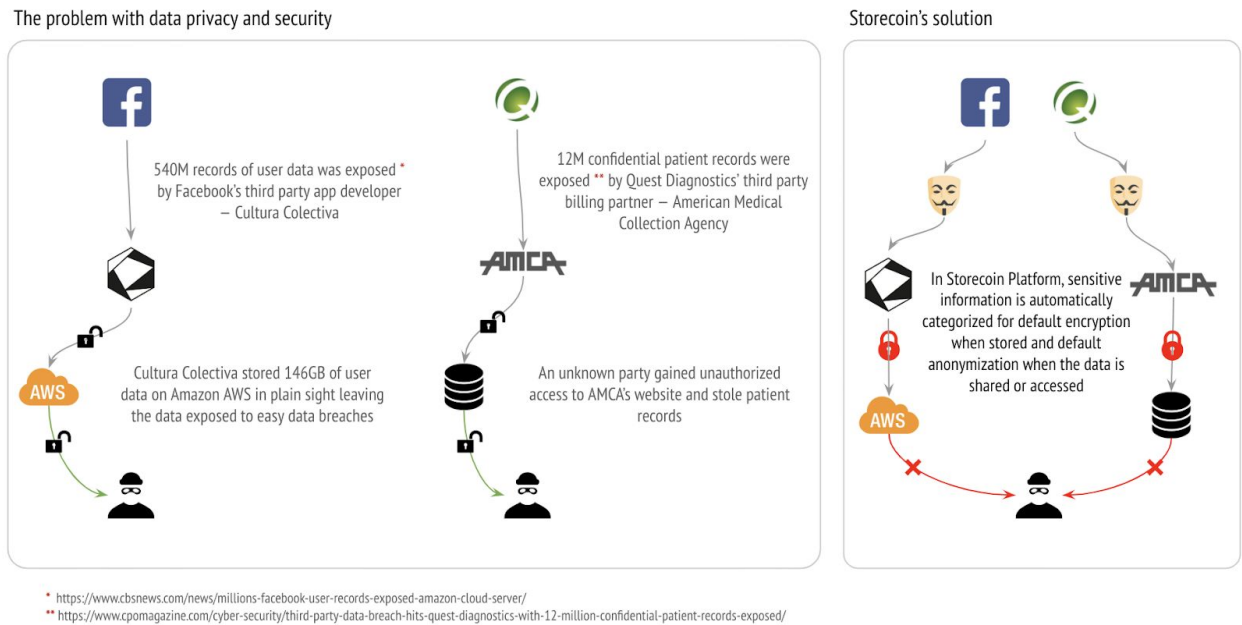# Storecoin platform — uses cases around data privacy

## Introduction

Storecoin Platform facilitates decentralization of data into programmable money called *datacoins*. The data created by the users of tokenized apps can be discovered, purchased, and used by any data buyer who is interested in the data. However, this poses a challenge when the data contains personally identifiable information or PII: how to protect the privacy of the data and its creator even when that data is traded on the platform? One option is enforcing restrictions on trading such data, but that's not practical as evidenced by ongoing data sharing and privacy violations [5] by prominent social networks, data exchanges, data brokers, and others. To date, our conversations about data privacy have been painted with too broad a brush stroke and have been generic to the point of irrelevance. When we ask questions about privacy and data, we need to ask what is the real threat to privacy that data trading represents. This leads us to personally identifiable information, which, if that's the case, means that our focus should be on **de-identifying** that data. Given this premise, a more practical approach would be to allow sharing or selling any information, including PII, but at the same time ensuring that the privacy of the data is protected. Storecoin Platform achieves de-identification with default encryption of sensitive information at rest and context-sensitive anonymization when sensitive information is accessed and traded. Storecoin's privacy and de-identification specification is available here.

In this document, we discuss the use cases around privacy-preserving data trading. For the purpose of this document, we select use cases that contain sensitive information such as the PII discussed above and discuss how sensitive information is treated while being traded on the platform.

# Privacy-preserved data-trading on the Storecoin Platform

Storecoin's approach to privacy preservation is illustrated in the following diagram using recent data exposures at Facebook and Quest Diagnostics as examples.



The problem with data privacy and security

540M records of user data was exposed * by Facebook's third party app developer — Cultura Colectiva

Cultura Colectiva stored 146GB of user data on Amazon AWS in plain sight leaving the data exposed to easy data breaches

12M confidential patient records were exposed ** by Quest Diagnostics' third party billing partner — American Medical Collection Agency

An unknown party gained unauthorized access to AMCA's website and stole patient records

Storecoin's solution

In Storecoin Platform, sensitive information is automatically categorized for default encryption when stored and default anonymization when the data is shared or accessed

* https://www.cbsnews.com/news/millions-facebook-user-records-exposed-amazon-cloud-server/
** https://www.cpomagazine.com/cyber-security/third-party-data-breach-hits-quest-diagnostics-with-12-million-confidential-patient-records-exposed/

Protecting user's privacy is hard because of different data security and retention policies among cooperating partners. Storecoin addresses this problem with:

- automatic identification and categorization of a default set of personally identifiable information
- automatic encryption of such sensitive information to protect against data theft, and
- automatic and context-sensitive anonymization of sensitive information to protect data privacy even when it is traded.

# Use cases [11]

The use cases discussed in this document are taken from real world scenarios as shown in the references. For each use case, we discuss how the de-identification schemes employed in the Storecoin Platform address privacy concerns.

# Use case 1 — privacy preservation while still providing researchers with the necessary data

*A national government project in central Europe was seeking to identify prisons that had populations that were at high risk for outbreaks of certain diseases so that they could intervene. They found that certain lifestyle traits, specifically a history of intravenous drug usage, piercings, and tattoos, had a high positive correlation with this disease. This lifestyle information was not codified and only existed in free form text notes. Their first solution was to manually mask or blur the records and supply the remaining information to the researchers. But it failed to achieve privacy objectives. Specific prisoners could often be identified. Their second solution was to use manual free form text data mining tools to extract only certain keywords, removing the entire record, and only supplying those keywords and the prison location. This proved successful. Their current plan is to use automated tools to identify key phrases, transform those into project-specific codified values, and then only supply that information along with the prison identifier to the researchers.*

## Analysis

- *This lifestyle information was not codified and only existed in free form text notes* — This poses a threat to the privacy of prison inmates, if the data is ever stolen. In Storecoin Platform, the lifestyle information is categorized as a form of PII and hence it is encrypted at rest.
- *Their first solution was to manually mask or blur the records and supply the remaining information to the researchers* — This is pseudonymization. While some information is hidden from the researchers, they are still able to re-identify the inmates based on the information shared with them. Storecoin allows app developers determine the exposure for re-identification, so they can ensure the right exposure without unintentionally exposing private information.
- *to extract only certain keywords, removing the entire record, and only supplying those keywords and the prison location. This proved successful* — This is a combination of anonymization and pseudonymization, with a right balance to expose just the right amount of information. Storecoin allows developers to run the data produced in their apps through various (combination of) de-identification schemes to learn how they can balance right exposure of user's data.
- *Their current plan is to use automated tools to identify key phrases, transform those into project-specific codified values, and then only supply that information along with the prison identifier to the researchers* — The project-specific codification uses contextualized anonymization. The prison identifier allows the officials to re-identify the inmates to proceed with the recommendations made by the researchers. Both of these are facilitated in Storecoin Platform.

## Example

The following table shows an example of this use case on Storecoin Platform. It compares the prisoner information before and after de-identification. The data in JSON format is used for simplicity. De-identification can work for data in any format. The specific scheme used for de-identifying each data field is described inline below.

| Prisoner information in clear text | Prisoner information after de-identification |
|---|---|
| ```{     prison_id: "<ID of prison location>",     prisoner_id: "<ID of the prisoner>",     name: "<Name of the prisoner>",     gender: "<male/female>",     age: <age>,     race: "<Race of the prisoner>",     lifestyle_info: "<sexual orientation,                     drug usage, tattoos,                     piercings, etc.>",     # Other information about the prisoner.     . . . }``` | ```{     # prison ID is retained as is.     prison_id: "<ID of prison location>",      # Directory replacement (DR) is used.     # This helps with identifying the prisoner     # later if needed.     prisoner_id: "<Mapped ID>",      # Name is masked. It could have been     # removed as well.     name: "**********",      # Gender is retained since researchers     # may need that information. However,     # this information could be masked if     # for example, there was only one     # female prisoner in this prison location,     # which identifies her indirectly. So,     # we can use contextualized anonymization     # (CA) also here.       gender: "<male/female>",      # DR is used to provide a range for the     # age.      age: <age low-high range>,      # CA is used to ensure that identity     # cannot be inferred through this     # information. So, this field can either     # be retained or masked.      race: "<Retained or masked>",      # CA and DR are used. CA is used to     # extract key phrases and DR is used to     # codify the key phrases. For example:     # Sexual orientation:     # 00 - Straight     # 01 - LGBT     # 02 - etc.     # Drug usage:     # 00 - clean     # 01 - Name of the drug used.     # 02 - etc.      # Except for mapped key phrases, all other``` |

```
                                          # information is removed. They are shown
                                          # as masked below for clarity.
                                          lifestyle_info: "<**** 01, *****
                                                           ** 07, **** 03,
                                                           00 **** , etc.>",

                                          # Other information is similarly
                                          # de-identified.
                                          . . .
                                          }
```

# Use case 2 — clinical trial

*A clinical trial is being planned that will involve independent reviewers of patient records to assess the response to an experimental drug. It may be necessary to inform patients of unusual findings. The trial sponsors set up a trial manager that will receive information from the physicians. The trial sponsor will perform the de-identification of the records, substituting clinical trial IDs for the original identifiers, obscuring dates, and redacting other non-clinical information. They chose to use a trial manager rather than ask the various patient physicians to perform de-identification based on the complexity of the trial requirements. The patients, physicians, and the trial sponsor agreed to allow a de-identification team access to the original patient data. The de-identification team and their systems are kept separate from the clinical trial results analysis. Only the de-identification team knows the relationship between clinical trial IDs and patient IDs. In the event that a significant finding is made by the review team, they communicate the finding to the de-identification team. The de-identification team contacts the patient's physician with the finding. The patient's physician examines the record and communicates with the patient. The physician informs the de-identification team that the patient has been informed. The de- identification team informs the review team, so that the review team can confirm that their ethical duty to ensure that the patient is informed has been met.*

## Analysis

This example involves 5 parties.

1. Drug manufacturer who wants to perform a clinical trial for an experimental drug.
2. Independent reviewers, who conduct the review of the experimental drug on a set of patients.
3. Patients who participate in the clinical trial.
4. Physicians who oversee informing patients of any unusual findings.
5. A trial manager who mediates between independent reviewers and physicians.

In this case, we have patient information on one side and trial results on the other side. Parties 1 and 2 are prohibited from identifying the patients directly and parties 3 and 4 are prohibited from

accessing the trial data. Trial manager is made responsible for de-identifying the patient information with the ability to re-identify them if needed.

- *substituting clinical trial IDs for the original identifiers, obscuring dates, and redacting other non-clinical information* — Original identifier substitution uses directory replacement approach with contextualized anonymization, so trial IDs can be customized for this particular clinical trial. The directory replacement allows for retrieving patient ID back when the patient needs to be communicated. Obscuring dates and redaction use masking and blurring.
- Use of a trial manager demonstrates a process where a third party may be used for protecting the data produced by the two sides of the trial from each other.
- In order for physicians to inform patients about unusual findings, re-identification is necessary. In this case, the trial data for the affected patient is made available without revealing the trial ID to the physicians.
- This use case demonstrates cooperation among multiple parties where their respective private data are de-identified and yet, they can perform their responsibilities without losing data precision.

## Example

In this use case, there are two sets of data that need to be de-identified — the patient information and trial data. The patient information is de-identified before sharing it with the trial researchers and similarly the trial data is de-identified before sharing it with physicians.

| Patient information in clear text | Patient information after de-identification |
|---|---|
| <pre>{<br>  hospital_id: "<ID of the hospital>",<br>  physician_id: "<ID of the physician>",<br>  patient_id: "<ID of the patient>",<br>  name: "<Name of the patient>",<br>  gender: "<male/female>",<br>  age: <age>,<br>  weight: <weight>,<br>  insurance_info: {<br>              <Insurance provider<br>              details>,<br>              },<br>  health_info: {<br>          <medical information such<br>          as prescriptions, medical<br>          conditions, test results,<br>          prior surgeries, allergies,<br>          etc.><br>          },<br><br>  payment_info: {<br>          <Payment information such<br>          as credit card/bank details,<br>          etc.></pre> | <pre>{<br>  # Hospital ID is retained as is.<br>  hospital_id: "<ID of the hospital>",<br><br>  # DR is used to protect the identify.<br>  # It is also used to re-identify the<br>  # physician, if they need to be contacted<br>  # with the unusual findings about the<br>  # patient.<br>  physician_id: "<ID of the physician>",<br><br>  # DR is used to protect the identify.<br>  # It is also used to re-identify the<br>  # patient, if they need to be contacted<br>  # with the unusual findings about them.<br>  patient_id: "<ID of the patient>",<br><br>  # Masked. It could be removed as well.<br>  name: "******",<br>  # Retained as is. This may be crucial<br>  # for the trial.<br>  gender: "<male/female>",<br><br>  # DR with a range.</pre> |

```
            },
   # Other information about the patient.
   . . .
}
```

```
age: <age low-high range>,

# Retained as-is.
weight: <weight>,

# Removed.
insurance_info: {
            <Insurance provider
            details>,
            },

# CA, DR, and Masking are used.
# CA is used to extract only the
# information that is needed for the
# trial.
# DR is used map the required medical
# information to codify what the trial
# needs.
# Masking is used to mask unwanted
# information. Such data could be removed
# as well if it is safe to do so.
health_info: {
            . . .
            },

# Removed. This information is not needed.
payment_info: {
            <Payment information such
            as credit card/bank details,
            etc.>
            },
# Other information about the patient.
. . .
}
```

The trial data is similarly de-identified to protect the identify of the researchers, any confidential information about the trial itself, etc. before sharing that information with the physicians.

# Use case 3 — Sharing PII with publishers and ad-exchange system

*There's been a whole lot of "implied consent" going on for some time now when it comes to collecting and sharing PII. That's because the online ad exchange system is flooded with players at different levels. There are not only the "publishers", selling advertising space based on user's personal data to "advertisers" hungry to get in front of the right eyeballs, but there are all sorts of third parties in between, brokering deals between advertisers and publishers. And once the information on whose eyeballs are up for sale goes out into the ad exchange ecosystem, there's no way to control the leakage of that data to numerous other parties. Even the advertisers who don't win the bid for the advertising space available still have some amount of access to the data they were invited to bid on, otherwise, how would they know if they wanted to buy those "eyeballs"? And that data is then used to update user profiles in various databases across the industry.*

*The central issue is data privacy "leak" as user's data propagates through multiple players in the ad exchange ecosystem. How this leak can be prevented in the light of the EU's General Data Protection Regulation*

*(GDPR)? There's even more concern industry-wide with the to-be-finalized ePrivacy Regulation, which is still being negotiated by the EU government, but which threatens to place an even heavier emphasis on consent for legal processing of personal data, casting a wider net than the GDPR.*

## Analysis

This use case illustrates the complexity of protecting user's privacy even with regulations such as the EU's GDPR. Since most services on the internet are either free or discounted because of the ad revenue, sharing PII will continue despite the regulations. The need for explicit consent from users for data collection has destroyed user experience by presenting users with an unceasing set of permissions that don't actually help with protecting their privacy.

- *Once the information on whose eyeballs are up for sale goes out into the ad exchange ecosystem, there's no way to control the leakage of that data to numerous other parties* — user data is shared among different players in the ad exchange ecosystem, so conformance to GDPR and other regulations is not sufficient. Data retention policies may be different with different players and may change over time, so there is no easy way to "take back" sensitive information after it is shared.
- *And that data is then used to update user profiles in various databases across the industry* — the curated and linked data across multiple players is never deleted or unlinked, thus posing greater risk to user's privacy. A piece of potentially useless data (such as a vacation photo) may become a privacy hazard if it is combined with the information from other players (such as name, address, income, etc.) because now the person in the photo can be identified definitively.
- *The central issue is data privacy "leak" as user's data propagates through multiple players in the ad exchange ecosystem* — the privacy leak may be unintended and the cooperating partners may not even be aware of this leak, so preventing new leaks and stopping any existing leaks are nearly impossible.

This use case demonstrates that regulations will not entirely prevent privacy leaks. What's more, the cost of complying with them makes it unfair for smaller players who cannot bear that cost, resulting in large players becoming larger.

## Example

The following example illustrates how de-identification prevents privacy leaks and yet, serves the needs of all the players in the ecosystem. For simplicity, the de-identified information is shown between a pair of cooperating partners, but the same process exists between all partners. It can also be noted that once a user's information is de-identified at one level further de-identification is unnecessary.

For simplicity, the user information is shown to be collected by one publisher. In practice, the information may have been gathered by multiple players. In any case, as long as de-identification process is followed between the partners, privacy leaks can be prevented.

This example also shows that multiple strategies exist for de-identification based on the need for sharing the information.

| User information in clear text (at a publisher, who collects this information) | User information after de-identification |
|---|---|
| <pre>{
   name: "<Name of the user>",
   address: "<address of the user>",
   contact_info: {
                <phone numbers,
                email addresses,
                social network accounts,
                linkedin, etc.>
              },
   gender: "<male/female>",
   age: <age>,
   SSN: "<National ID>",
   drivers_id: "<Government issued ID>",
   credit_card_info: {
                   <credit cards with
                    expiry dates, etc.>
                 },
   financial_info: {
                <Bank accounts,
                brokerage accounts,
                401K, etc.>,
                },
   health_info: {
              <medical information such
              as prescriptions, medical
              conditions, health risks,
              etc.>
            },
   work_info: {
              <Companies worked for,
              salary information,
              titles, etc.>
            },
   online_info: {
              <IP addresses, cookies,
              devices, history of sites
              visited, search history,
              etc.>
            },
   # Other information about the user.
   . . .
}</pre> | <pre>{
   # **Masked**, **blurred**, or **DR**'d depending on
   # the need. If there is a need for
   # re-identifying the user, **DR** is used.

   name: "~~<Name of the user>~~",

   # **Masked**, **blurred**, or **DR**'d depending on
   # the need.

   address: "~~<address of the user>~~",

   # **Masked**, **blurred**, or **DR**'d depending on
   # the need.

   contact_info: ~~{~~
                ~~<phone numbers,~~
                ~~email addresses,~~
                ~~social network accounts,~~
                ~~linkedin, etc.>~~
              ~~}~~,
   # Can be shared as is.
   gender: "<male/female>",

   # **DR**'d with an age range.
   age: <age low-high range>,

   # **Masked**, or **blurred**.
   SSN: "~~<National ID>~~",

   # **Masked**, or **blurred**.
   drivers_id: "~~<Government issued ID>~~",

   # **Masked**, **blurred**, or **removed**. However,
   # it is also possible to share "codified"
   # information about credit worthiness
   # of the user also. If that's required
   # a **CA** is used for computing the score,
   # for example.
   credit_card_info: ~~{~~
                   ~~<credit cards with~~
                    ~~expiry dates, etc.>~~
                 ~~}~~,


   # Same as credit_card_info above.
   financial_info: {
                <Bank accounts,
                brokerage accounts,</pre> |

```
                                 401K, etc.>,
                             },

        # Masked, blurred, or removed. However,
        # it is also possible to share "codified"
        # information about the health score
        # of the user also. If that's required
        # a CA is used for computing the
        # health score, for example.

        health_info: {
                     <medical information such
                     as prescriptions, medical
                     conditions, health risks,
                     etc.>
                     },

        # Same as health_score above.
        work_info: {
                     <Companies worked for,
                     salary information,
                     titles, etc.>
                     },

        # May be shared as is.
        online_info: {
                     <IP addresses, cookies,
                     devices, history of sites
                     visited, search history,
                     etc.>
                     },
        # Other information about the user.
        . . .
}
```

The use cases discussed in this section touch different domains and different data types, but it can be observed that de-identification, especially when used with contextualized anonymization, can address data privacy adequately in all cases. So, this approach can be used safely even when PII is shared. Coupled with encryption, the privacy of the users of Storecoin Platform will be protected at rest as well as when their information is shared.

## Use case 4 — Data aggregation and consumer privacy in Fintech [13]

*WASHINGTON — Federal regulators gathered Wednesday to discuss the fintech chartering process and some of the biggest challenges deterring the emerging industry from entering the banking space.*

*During the event, which covered a lot of ground beyond fintech, Treasury Secretary Steven Mnuchin and McWilliams agreed that they need to address data aggregation at banks and how consumers control their own data, as the FDIC is beginning to study the issue.*
*Federal Deposit Insurance Corp. Chairman Jelena McWilliams and Controller of the Currency Joseph Otting Otting said separately that there was agency consensus on modernizing the law, revisiting CRA assessment areas and determining how lending data is being collected, for example.*

*McWilliams said the FDIC has begun studying data aggregation and how data is shared between banks and third parties, including with fintech vendors.*

*"We need to take a look at this," she said, clarifying it would be a "truly preliminary" study, not a rulemaking. "There are privacy concerns and cyber concerns. And from the FDIC's perspective, the third-party vendor management is crucial."*

*Specifically, McWilliams said the FDIC is looking at who owns the data between parties as well as how much consumers have a right to their own data and whether more data should be shared.*

*Earlier in the day, McWilliams interviewed the Treasury's Mnuchin, who said that data aggregation was a critical matter but one that ought to be addressed through the private sector.*

*"This is a complicated issue and I would say my view from a consumer standpoint is, it should be very clear and very simple if your data is being shared, who it's being shared with," Mnuchin said. "In general, I like where there are private solutions as opposed to a government solutions."*

## Analysis

- *… need to address data aggregation at banks and how consumers control their own data* — this use case presents a scenario in which data aggregation among cooperating banks and their partners is inevitable and yet, consumers can somehow control their own data without sacrificing their privacy. Should this be addressed with regulations or can technology help with solving this issue?
- *FDIC has begun studying data aggregation and how data is shared between banks and third parties, including with fintech vendors* — data aggregation is necessary to provide personalized and differentiated services to customers, but how can the banks ensure that the privacy of their users is not sacrificed when they share the data with their partners?
- *There are privacy concerns and cyber concerns. And from the FDIC's perspective, the third-party vendor management is crucial* — the privacy concerns arise from the fact that the banks and their partners may have completely different data retention and privacy policies, so banks don't have much control when the data leaves their system. The cyber concerns arise from the fact that any sensitive information may not be encrypted end to end in the partner chain, thus exposing user data in the weakest link of the chain. Vendor management is needed to ensure that these policies are aligned throughout the chain, but it is very hard to achieve given the dynamic nature of the participants in the chain.
- *FDIC is looking at who owns the data between parties as well as how much consumers have a right to their own data and whether more data should be shared* — data ownership after it leaves one

partner to another in the chain is hard to define, given ever changing policies on data encryption, retention, and privacy. Is there a better solution that is agnostic to individual partner's privacy policies?

- *from a consumer standpoint is, it should be very clear and very simple if your data is being shared, who it's being shared with* — data sharing is opaque in today's system. Consumers don't have where some of their personal information is originated and who have access to that information? How can consumers be confident that their privacy is not compromised by others in the data aggregation chain?

## Example

The following example addresses the concerns raised above with de-identification and encryption. Any bank, which collects consumer data (*source* of data collection) encrypts the PII by default. This addresses cyber concerns discussed above. The data aggregation between banks and partners makes use of de-identification — specifically contextualized anonymization — so even PII can be shared between the partners without sacrificing user privacy and quality of data aggregation. The example illustrates how de-identification prevents privacy leaks and yet, serves the needs of all the players in the ecosystem. For simplicity, the de-identified information is shown between a pair of cooperating partners, but the same process exists between all partners. It can also be noted that once a user's information is de-identified at one level further de-identification is unnecessary.

For simplicity, the user information is shown to be collected by one bank. In practice, the information, full or in part, may have been gathered by one or more banks or their partners. In any case, as long as de-identification process is followed between the partners, privacy leaks can be prevented.

This example also shows that multiple strategies exist for de-identification based on the need for sharing the information.

| User information in clear text (at a bank, who collects this information) | User information after de-identification |
|---|---|
| ```{    name: "<Name of the consumer>",    address: "<address of the consumer>",    contact_info: {               <phone numbers,               email addresses, etc.>             },    gender: "<male/female>",    marital_status:    "<married/single/widow(er)/               etc>",    age: <age>,    SSN: "<National ID>",``` | ```{    # Masked, blurred, or DR'd depending on    # the need. If there is a need for    # re-identifying the user, DR is used.    name: "<Name of the consumer>",    # Masked, blurred, or DR'd depending on    # the need.    address: "<address of the consumer>",    # Masked, blurred, or DR'd depending on``` |

```
    drivers_id: "<Government issued ID>",
    income: <annual income>,
    job_title "<job title>",

    # The following data are created by
    # this bank about this consumer.

    financial_info: {
                    <
                        Bank account number,
                        account type,
                        account balance,
                        transactions,
                        credit standing,
                        payment history,
                        etc.
                    >,
                    },

    # Other information about the consumer.
    . . .
}
```

```
    # the need.

    contact_info: +
                    <phone numbers,
                    email addresses, etc.>
                +,
    # Can be shared as is. But can be
    # Masked or blurred also.
    gender: "<male/female>",

    # Can be shared as is. But can be
    # Masked or blurred also.
    marital_status:
    "<married/single/widow(er)/
                    etc>",

    # DR'd with an age range.
    age: <age low-high range>,

    # Masked, or blurred.
    SSN: "<National ID>",

    # Masked, or blurred.
    drivers_id: "<Government issued ID>",

    # DR'd with an income range.
    income: <income low-high range>,

    # Can be shared as is. But can be
    # Masked or blurred also.
    job_title "<job title>",

    # Individual fields are de-identified
    # differently.
    financial_info: {
                    <
                    # Masked, blurred, or DR'd
                        Bank account number,
                    # Probably shared as-is.
                        account type,
                    # DR'd with a range
                        account balance range,
                    # Removed or Masked
                        transactions,
                    # Probably shared as-is.
                        credit standing,
                    # Probably shared as-is
                        payment history,
                        etc.
                    >,
                    },

    # Other information about the consumer.
    . . .
}
```

The de-identification strategy used depends on whether re-identification of the consumer is required at a later date. Notice also that multiple strategies can be used depending on specific use cases, so no generalized and regid strategies are needed. For example, the data for a specific consumer may use a completely different set of de-identification strategies from another consumer.

Since the information is de-identified, privacy and data retention policies of the partner don't matter as the information cannot be used by the partner to identify the consumer directly. At the same time, data aggregation goals are achieved without losing accuracy. If the partner later produces aggregated information and shares them back with the bank, it can choose to de-identify certain data to protect its own privacy.

## Summary

- Prohibiting sharing/selling PII is not practical. But at the same time Storecoin Platform cannot ignore data privacy.
- Data privacy is ensured via de-identification. First, Storecoin Platform allows app developers categorize the app data, so any PII can be identified as such. Once a piece of data is identified as PII, it is automatically encrypted at rest and de-identified when the data is requested.
- Anonymization and pseudonymization are turned on by default. This means, data privacy is turned on by default. An app can however, turn off this default behavior if its use cases demand that. Such apps are flagged as not protecting user's privacy by Storecoin Platform, so that the users are well informed.
- Apps can use different anonymization techniques depending on the data. The same data may be anonymized differently by different apps depending on their customer base. Even within a given app, the data may be de-identified differently based on the request types and use cases.
- Not all information in the platform can be shared or traded. Certain *control* information are sensitive in nature and their access is protected by not categorizing them. Only categorized data can be discovered and traded.

## References

1. https://piwik.pro/blog/what-is-pii-personal-data/
2. https://www.experian.com/blogs/ask-experian/what-is-personally-identifiable-information/
3. http://www.ncsl.org/research/telecommunications-and-information-technology/data-disposal-laws.aspx
4. https://www.ftc.gov/system/files/documents/public_events/1223263/p155407privacyconmislove_1.pdf
5. http://fortune.com/2018/03/23/facebook-data-scandal-carolyn-everson/
6. http://www.nist.gov/
7. https://eugdpr.org/

8.  https://iapp.org/news/a/looking-to-comply-with-gdpr-heres-a-primer-on-anonymization-and-pseudonymization/
9.  https://en.wikipedia.org/wiki/De-identification
10. https://www.adweek.com/digital/security-firm-finds-millions-of-facebook-data-files-were-stored-on-amazons-public-cloud-servers/
11. http://ihe.net/uploadedFiles/Documents/ITI/IHE_ITI_Handbook_De-Identification_Rev1.1_2014-06-06.pdf
12. https://www.cmu.edu/iso/governance/guidelines/data-classification.html
13. https://www.americanbanker.com/list/fintech-charters-cra-and-data-sharing-fdics-mcwilliams-occs-otting-weigh-in
14. https://en.wikipedia.org/wiki/Data_breach