

generating the anonymization dictionary comprises: using a statistical method or an artificial intelligence method to detect, in the activity record, a plurality of target entities to be anonymized; assigning an anonymized identity to each unique target entity of the plurality of target entities; generating dictionary entries for the plurality of target entities, wherein each dictionary entry comprises a target entity and a corresponding anonymized identifier comprising the anonymized identity for the target entity; generating, by the activity monitoring engine, an equivalence map, wherein generating the equivalence map comprises: making a determination that a resource is associated with a set of target entities of the plurality of target entities; storing, in the equivalence map, an identity relationship specifying that anonymized identities corresponding to the set of target entities are associated with the resource; replacing, by the activity monitoring engine, the plurality of target entities in the activity record with their anonymized identifiers from the anonymization dictionary to obtain an anonymized activity record; and storing, by the activity monitoring engine, the anonymized record.

12. The method of claim 11, further comprising detecting an additional target entity to be anonymized based on its consistent appearance uniquely with one of the plurality of target entities.

13. The method of claim 11, wherein making the determination that the resource is associated with the set of target entities comprises using statistical analysis to determine user accounts are associated with the resource when the user accounts access an email account more than the user accounts access other email accounts.

14. The method of claim 11, wherein making the determination that the resource is associated with the set of target entities comprises using cluster analysis to determine user accounts are associated with the resource based on utilization parameters.

15. The method of claim 11, wherein generating the equivalence map is performed prior to or after storing the anonymized record.

16. A computer system, comprising: the anonymization engine programmed to perform the method of claim 1 or 11; and a repository configured to store the anonymized activity record.

17. The system of claim 16, further comprising: a threat analysis engine programmed to analyze the anonymized activity record using a threat detection algorithm.

Description

BACKGROUND

[0001] Activity on a company's computing devices may be tracked in order to detect behaviors that may pose threats.

SUMMARY

[0002] In general, in one aspect, the invention relates to a method for processing activity records. The method includes obtaining an activity record, and generating an anonymization dictionary. Generating the anonymization dictionary includes detecting, in the activity record, a set of target entities to be anonymized, making a determination that a resource is associated with a subset of the target entities of the set of target entities, and after making the determination, assigning an anonymized identity to the subset of target entities, and generating an anonymization identifier for each target entity in the subset of target entities to obtain a set of anonymization identifiers, each including the anonymized identity. The method further includes processing the activity record using the anonymization dictionary to obtain an anonymized activity record and storing the anonymized activity record.

[0003] In general, in one aspect, the invention relates to a method for processing activity records. The method

includes obtaining an activity record, and generating an anonymization dictionary. Generating the anonymization dictionary includes detecting, in the activity record, a set of target entities to be anonymized, assigning an anonymized identity to each unique target entity of the set of target entities, and generating dictionary entries for the set of target entities. Each dictionary entry includes a target entity and a corresponding anonymized identifier including the anonymized identity for the target entity. The method further includes generating an equivalence map. Generating the equivalence map includes making a determination that a resource is associated with a subset of the target entities of the set of target entities and storing, in the equivalence map, an identity relationship specifying that the subset of the target entities is associated with the resource. The method also includes processing the activity record using the anonymization dictionary to obtain an anonymized activity record and storing the anonymized record.

[0004] In general, in one aspect, the invention relates to a system for processing activity records. The system includes an anonymization engine programmed to obtain an activity record, and generate an anonymization dictionary. Generating the anonymization dictionary includes detecting, in the activity record, a set of target entities to anonymize, making a determination that a resource is associated with a subset of the target entities of the set of target entities, and assigning an anonymized identifier to the subset of target entities. The anonymization engine is further programmed to process the activity record using the anonymization dictionary to obtain an anonymized activity record. The system also includes a repository configured to store the anonymized activity record.

[0005] In general, in one aspect, the invention relates to a system for processing activity records. The system includes an anonymization engine programmed to obtain an activity record, and generate an anonymization dictionary. Generating the anonymization dictionary includes detecting in the activity record, a set of target entities to be anonymized, assigning an anonymized identity to each unique target entity of the set of target entities, and generating dictionary entries for the set of target entities. Each dictionary entry includes a target entity and a corresponding anonymized identifier including the anonymized identity for the target entity. The anonymization engine is further programmed to generate an equivalence map. Generating the equivalence map comprises making a determination that a resource is associated with a subset of the target entities of the set of target entities, and storing, in the equivalence map, an identity relationship specifying that the subset of the target entities is associated with the resource. Generating the equivalence map further comprises processing the activity record using the anonymization dictionary to obtain an anonymized activity record. The system also includes a repository configured to store the anonymized activity record.

[0006] In general, in one aspect, the invention relates to a method for processing activity records. The method includes obtaining a set of activity records, providing at least one of the set of activity records to each of a set of workers, receiving, from each of the set of workers, a set of target entities, generating, using the sets of target entities, an anonymization dictionary, providing a copy of the anonymization dictionary to each of the set of workers, receiving, from each of the set of workers, at least one anonymized activity record generated using the copy of the anonymization dictionary, and storing the anonymized activity records.

[0007] In general, in one aspect, the invention relates to a method for processing activity records. The method includes obtaining a first set of anonymized activity records and a first local anonymization dictionary from a first endpoint agent, and obtaining a second set of anonymized activity records and a second local anonymization dictionary from a second endpoint agent. The method further includes storing the first set of anonymized activity records, the first local anonymization dictionary, the second set of anonymized activity records, and the second local anonymization dictionary, and performing a threat analysis using the first set of anonymized activity records, the first local anonymization dictionary, the second set of anonymized activity records, and the second local anonymization dictionary.

[0008] Other aspects of the invention will be apparent from the following description and the appended claims.

BRIEF DESCRIPTION OF DRAWINGS

[0009] FIGS. 1A-1D show systems in accordance with one or more embodiments of the invention.

including domain names, IP addresses, port numbers, host names, etc.) affiliated with an activity documented in the activity records. Further, in the event that the identity of the resource is required, the resource may be identified by an authorized viewer of the activity records. Authorized viewers may be, for example, employee supervisors, employees of the human resources department, company-internal and/or external security analysts, etc. Depending on the level of authorization, an authorized viewer may review the fully de-anonymized (i.e., original) activity records, or activity records that remain anonymized to various, configurable, degrees.

[0018] When a threat alert is issued by the threat analysis algorithms based on detected unsafe activity, a supervisor may be notified. The supervisor may access the anonymized activity record for a review of the activities that triggered the threat alert. Depending on the privileges of the supervisor, the activity record may be partially or fully de-anonymized, allowing the supervisor to assess the threat potential of the activities.

[0019] Anonymization of activity records may be used to protect the identities of resources including user's identities but also company-related information, such as, for example, company-internally used domain names and IP addresses. Such company-internal information may need to be anonymized, for example, in cases where activity records are shared with 3.sup.rd parties in order to avoid unintentional disclosure of company-internal information such as IT infrastructure details. Anonymized activity records may be shared with external 3.sup.rd parties, for example, in cases where a threat analysis is performed externally, e.g. by a service provider offering threat monitoring and analysis. Further, anonymized activity records may be shared with external developers of threat detection algorithms and/or anonymization algorithms.

[0020] FIG. 1A shows a system in accordance with one or more embodiments of the invention. The system may include one or more users (102) that may interact with computing devices A-N (104A-104N). Each computing device (104A-104N) may host an endpoint agent (106A-106N) that may generate activity records, based on a user's interaction with the computing device. The system in accordance with one or more embodiments of the invention further includes an activity monitoring engine (108) that may analyze activity records received from the endpoint agents A-N (106A-106N) for threats, as further described below, with reference to FIG. 1B and FIGS. 2A, 2B, 3A and 3B. The system in accordance with one or more embodiments of the invention may also include third party systems (110) that may contribute to the threat detection and analysis as further described below. Each of these components is described below.

[0021] In one or more embodiments of the invention, users (102) operate one or more of the computing devices A-N (104A-104N). The users may have user accounts that may, for example, grant access to certain resources such as content stored on connected company servers, software applications, etc. User accounts may be personalized based on the needs of the users. For example, a design engineer may have access to technical design resources such as mechanical parts libraries, while not being allowed to access sales data, whereas an employee of the human resources department may have access to personnel data while not being allowed to access technical design resources and sales data.

[0022] The computing devices (104A-104N) may be used by the users (102) to perform work-related tasks. The computing devices may, however, also be abused, for example, by users accessing data in an unauthorized manner, bypassing security measures, using pirated applications and/or media, copying sensitive information on external, removable storage media, etc. In addition, the computing devices may face company-external threats, caused, for example, by hacking attacks, and malware. A computing device (104A-104N) may be, for example, a mobile device (e.g., a laptop computer, smart phone, personal digital assistant, tablet computer, or other mobile device), a desktop computer, server, blade in a server chassis, or any other type of computing device or devices that includes at least the minimum processing power, memory, and input and output device(s) to perform one or more embodiments of the invention. A computing device may include one or more computer processor(s), associated memory (e.g., random access memory (RAM), cache memory, flash memory, etc.), one or more storage device(s) (e.g., a hard disk, an optical drive such as a compact disk (CD) drive or digital versatile disk (DVD) drive, a flash memory stick, etc.), and numerous other elements and functionalities, such as input and output device enabling a user to interact with the computing device. The computing device may further be connected to a network (e.g., the company's local area network (LAN), a wide area network (WAN) such as the Internet, mobile networks, or any other type of network via a network interface connection.

threat analysis engine uses threat detection algorithms to determine whether an anonymized activity record includes indications of threats. The threat detection algorithm may evaluate activities stored in the anonymized activity records, and if an abnormal activity is detected, the threat detection algorithm may issue an alert. The threat detection algorithm may further quantify the risk resulting from abnormal activities. A high risk score may indicate an elevated risk thus warranting an alert, whereas a lower score may not necessarily trigger an immediate alert. The detection of abnormal activities may be based on a comparison with typical, i.e. expected activities. For example, a user activity stored in an activity record may be compared to typical behavior of a user, as per the user's role in the company. Examples for abnormal user behavior incompatible with the user's role include an engineer copying a customer list to an external storage device, a salesman copying confidential engineering records to an external storage device, etc. Alternatively, or additionally, the detection of abnormal user behavior may be based on a comparison with historical user behavior and/or data from known previous insider-threat cases. For example, a company's employee that primarily relied on the Internet to research suppliers' products but that recently started to use the Internet in order to contact direct competitors of the company may also be considered suspicious. Further, the threat detection may compare resource activity stored in an activity record with historical and/or typical resource activity. For example, a sustained access to a hard drive may be considered suspicious if the accessed hard drive has historically been mostly idle. Other methods for performing threat detection may be performed by the threat analysis engine without departing from the invention. Threat detection may, for example, in addition involve a human operator, e.g., a security expert, performing manual threat detection, and/or a manual review of threats detected by the threat analysis engine.

[0035] The threat analysis engine may include one or more application programming interfaces (APIs) to permit interaction with the third party systems (110), for example, to share anonymized activity records with 3.sup.rd parties, to access threat analysis algorithms and/or anonymization engines developed by 3.sup.rd parties, download them, and potentially set them up to replace currently used algorithms.

[0036] The threat analysis engine (122) may be executing on a single computing device, e.g., on the server that also executes the anonymization engine. In other implementations, the threat analysis engine may be executing on a separate computing device, either locally or remotely, e.g. hosted on a cloud server. Further, multiple threat analysis engines may be used to perform a distributed threat analysis.

[0037] In one or more embodiments of the invention the threat analyst interface may serve as a user interface providing access to various functions of the system for the detection of cyber-threats. The threat analyst interface may, for example, display alerts triggered by detected threats or potential threats. The threat analyst interface may further display threat-related information, which may, for example, also include the anonymized activity record that triggered the threat alert. The threat analyst interface may further include configurable filters that allow selective displaying of threats, potential threats, and threat-related information. For example, a filter may be used to display only threat-related information related to activities of a particular user, or a group of users. In addition, a filter may be configured to suppress alerts for abnormal activities where the risk score, computed by the threat analysis engine, does not exceed a set threshold. In one embodiment of the invention, the threat analyst interface may further display de-anonymized or partially de-anonymized versions of anonymized activity records. The amount of de-anonymization may depend on the operator's level of authorization, and may range from complete de-anonymization, for an operator that is fully authorized to view sensitive user data to no de-anonymization if the operator is only equipped with basic viewing privileges.

[0038] In one embodiment of the invention, the threat analyst interface (124) may further be used to configure various components of the activity monitoring engine (108). The threat analyst interface may be used, for example, to parameterize the anonymization engine (112) and the threat analysis engine (122). For example, the threat analyst interface may be used to configure the degree of anonymization performed by the anonymization engine (112), and/or to select and parameterize a particular threat detection algorithm to be used by the threat analysis engine (122). In addition, the threat analyst interface (124) may provide access to the anonymization correspondence repository (114), allowing the viewing and editing of anonymization correspondence repository content.

[0039] In one or more embodiments of the invention, the threat analyst interface (124) is a graphical user

interface (GUI), that is operatively connected to the anonymization engine (112) and the threat analysis engine (122), and that may further be operatively connected to the anonymization correspondence repository (114) and the anonymized activity record repository (120). The threat analyst interface (124) may execute on any computing device that provides the input and output interfaces necessary for an operator to interact with the threat analyst interface. A suitable computing device may be, for example, a desktop computer, a server, or a mobile device (e.g., a laptop computer, smart phone, personal digital assistant, tablet computer, or other mobile device).

[0040] FIG. 1C shows an anonymization correspondence repository (114) in accordance with one or more embodiments of the invention. The anonymization correspondence repository includes an anonymization dictionary (116). The anonymization dictionary (116) includes resource profiles (150.1-150.Z). Resource profiles may be used to structure information that may be attributed to a particular resource, e.g., a particular user or a particular company. For example, a user may be a resource that has multiple email addresses for which the anonymization dictionary includes entries under the same resource profile, and/or a company may have multiple resources such as IP addresses, domain names, etc., for which the anonymization dictionary includes entries under the same resource profile. Separate resource profiles may therefore exist as necessary to group resources. Each resource profile may include entries, each entry including a target entity (152.1.1-152.2.2) and a corresponding anonymized identifier (154.1.1-154.2.2). An anonymized identifier (154.1.1-154.2.2) may include an anonymized identity (ID) (156.1, 156.2), an entity type (158.1.1-158.2.2) and instance ID (160.1.1-160.2.2). Each of these elements is described below. Target entity identifiers and corresponding anonymized identifiers where an affiliation with a particular resource profile is unknown may exist in the dictionary without being grouped under a common resource profile.

[0041] The pairs of target entity and anonymized identifier in the anonymization dictionary (116) may store information used for the anonymization and de-anonymization of activity records by establishing a relationship between the target entities (152.1.1-152.2.2) in the activity records to be anonymized, and the anonymized identifiers (154.1.1-154.2.2) used to replace the target entities during the anonymization of the activity records. A target entity may be a term in the activity record that may reveal the identity of the resource associated with the target entity or that may reveal sensitive company information (i.e., sensitive information about the company (or legal entity) with which the resource is associated), and which therefore may need to be anonymized in order to protect the resource's identity and/or protect sensitive company information. A target entity may be, for example, a user name, a login name, an email address, a company name, a partner name, an IP address, a domain name, a host name, etc. The anonymized identifier may be a descriptor that does not allow the identification of the resource or sensitive company information without looking up the relationship of anonymized identifier and target entity in the anonymization dictionary. The anonymized identifier may include an anonymized identity (ID) (156.1, 156.2), an entity type (158.1.1-158.2.2) and an instance ID (160.1.1-160.2.2). A unique anonymized ID (156.1, 156.2), specific to a resource profile, may be, for example, any type of string, a number, symbols, or any combination thereof. The anonymized ID may be randomly or systematically selected. For example, user identifiers used in the dictionary may be "USER1" "USER2" "USER3" etc., or "000" "001" "002", etc.

[0042] In one or more embodiments of the invention, an anonymized identifier (154.1.1-154.2.2) may further include an entity type (158.1.1-158.2.2) to classify the target entity. Classifications may include, for example, "user name", "email address", "company name", "domain name", "host name", "IP address", and "other" for other selected strings deemed sensitive data.

[0043] In one or more embodiments of the invention, the anonymization dictionary may include multiple target entities with the same entity type. Consider a scenario, where user John Smith uses the email accounts "j.smith@gmail.com", "john_smith@yahoo.com" and "john.smith@company.com". Each one of these email addresses, if stored in an activity record, may be detected as a target entity by the anonymization engine. Accordingly, the anonymization dictionary may include three entries for the three email addresses of user "John Smith". All three target entities are of the entity type "email address". Further, all three entries are organized in a single resource profile because they are affiliated with the same resource (user "John Smith").

[0044] Multiple entries having the same entity type may therefore be created under the same resource profile in

"email address". The anonymization dictionary, however, does not define an affiliation with a common resource (user "John Smith") for these entries. The affiliation with a common resource may therefore not be documented by the anonymization dictionary, in accordance with this particular embodiment of the invention. To perform an anonymization of a target entity detected in an activity record, the anonymization engine may search the anonymization dictionary for the target entity. The anonymization engine may then replace the target entity with the anonymized identifier that corresponds to the target entity in the anonymization dictionary. If the anonymization dictionary does not include an entry for the target entity, the anonymization engine may add such an entry to the anonymization dictionary.

[0051] One skilled in the art will recognize that the system for anonymizing activity records is not limited to the components shown in FIGS. 1A-1D. For example, various components, including the anonymization engine, the anonymization dictionary, the anonymized activity record repository and the threat analysis engine, may exist repeatedly, either locally, or distributed over multiple computing devices on-premises, off-premises and/or cloud-based, for example, in order to perform distributed anonymization of activity records. Further, anonymization may be performed to different degrees, either in parallel, or serially. In addition, activity records to be anonymized may not necessarily originate from a computing device equipped with an endpoint agent. For example, the activity records may be obtained via an application programming interface (API) from third party applications, or any other source that may provide activity records that require anonymization. Also, even though particular structures of the anonymization correspondence repository are shown in FIGS. 1C and 1D, the anonymization correspondence repository may be structured in any other way, as long as the anonymization dictionary included in the anonymization repository establishes relationships between target entities and anonymized identifiers.

[0052] FIGS. 2A, 2B, 3A and 3B show flowcharts in accordance with one or more embodiments of the invention. While the various steps in the flowcharts are presented and described sequentially, one of ordinary skill will appreciate that some or all of these steps may be executed in different orders, may be combined or omitted, and some or all of the steps may be executed in parallel. In one embodiment of the invention, the steps shown in FIGS. 2A, 2B, 3A and 3B may be performed in parallel with any other steps shown in FIGS. 2A, 2B, 3A and 3B without departing from the invention.

[0053] FIGS. 2A and 2B show methods for anonymizing activity records. During the anonymization of activity records, target entities to be anonymized may be identified, and subsequently these target entities may be replaced by anonymized identifiers. The anonymized identifiers used for replacing the target entities may be obtained from an anonymization dictionary. If an entry for the target entity does not exist in the anonymization dictionary, or if the anonymization dictionary itself does not exist, an entry in the anonymization dictionary or the anonymization dictionary itself, respectively, may be generated prior to replacing the target entities by the anonymized identifier. What constitutes a target entity to be replaced may be configurable. For example, only user identity-related target entities, e.g., user names, email addresses, etc. may be anonymized if the resulting anonymized activity records remain within the company. Alternatively, company resource-related target entities in general, e.g., user names, email addresses, domain names, IP addresses, etc. may be anonymized, for example if the resulting anonymized activity records are shared with external 3.sup.rd parties. Accordingly, the methods described in FIGS. 2A and 2B may be executed repeatedly, serially or in parallel, for the same activity records to obtain different degrees of anonymization. The method may be executed whenever an activity record or a set of activity records is obtained from an endpoint agent. Alternatively, the method may only be executed when activity records from a particular endpoint agent or a group of endpoint agents are received, while activity records from other endpoint agents are ignored. In one or more embodiments of the invention, certain steps of the methods described in FIGS. 2A and 2B may be performed in a distributed manner. The details of distributed versus local execution of these steps are discussed below.

[0054] In the following discussion of FIGS. 2A and 2B, FIG. 2A describes an embodiment of the invention where a dictionary may be generated in a first pass over the activity record(s), and where subsequently, the activity record(s) may be anonymized in a second pass over the activity record(s), whereas FIG. 2B describes an embodiment of the invention where the generation of the dictionary and the anonymization of the activity record(s) may be performed in a single pass over the anonymization record(s). In one embodiment of the

invention, the methods may be executed independently each time an activity record or a set of activity records is received, i.e., subsequent receipt of an additional activity record may result in renewed execution of the methods, independent from the previous execution. Repeated execution of the methods may therefore result in the generation of separate, independent anonymization correspondence repositories that may only be applicable to the activity record(s) from which they are derived. In an alternative embodiment, subsequent execution of the methods for a newly received user activity record may be dependent upon previous execution of the methods for a previously received user activity record. Accordingly, the anonymization correspondence repository, established when processing previously received activity records, may be used to anonymize a subsequently received activity record. During the processing of the subsequently received activity record, additional entries may be added to the anonymization correspondence repository, if the subsequently received activity record includes target entities not yet included in the anonymization correspondence repository.

[0055] Turning to FIG. 2A, in Step 200, the anonymization engine of the activity monitoring engine obtains an activity record from an endpoint agent. The activity record may be obtained, for example, by the endpoint agent pushing the activity record as it becomes available, or by the activity monitoring engine polling the endpoint agent for new activity records. Activity records may be obtained continuously, as they are generated by an endpoint agent, or they may be obtained in batches, for example, in scenarios where the endpoint agent accumulates activity records and provides them to the anonymization engine at fixed time intervals. In one embodiment of the invention, the activity record may be access-protected for the transmission from the endpoint agent to the anonymization engine, e.g., using encryption.

[0056] In Step 202, a determination is made about whether an anonymization correspondence repository exists. If no anonymization dictionary exists, the method may proceed to Step 204.

[0057] In Step 204, the activity record, obtained in Step 200, is scanned for target entities to be anonymized. Various methods for detecting target entities may be employed. The following includes a description of exemplary methods for detecting target entities. The exemplary methods are not intended to limit the scope of the invention.

Exemplary Method 1--Information About Target Entities May be Provided

[0058] For example, the IT administrator may provide a list of user names, email addresses, company names, host names, and/or domain names, etc. In addition, or alternatively, arbitrary strings, deemed sensitive, may be provided. Accordingly, the scanning of an activity record may be performed based on the provided target entities.

Exemplary Method 2--Information About Target Entities May be Inferred from Information That is Accessible to the Anonymization Engine

[0059] For example, the anonymization engine may parse the local directories of servers and users' computing devices for user profiles that may reveal information about target entities. Consider, for example, a computing device using a Microsoft Windows.RTM. operating system. In such a system, information about the users of the system may be obtained by inspecting the "\Users" directory, and/or the "\Documents and Settings" directory which may include login names of the users of the system. User names in the format "domain\username" may further allow the extraction of the domain name.

Exemplary Method 3--Information About Target Entities may be Derived

[0060] Consider a scenario where a user "John Smith" is already known. The known user name "John Smith" may be used to derive potential additional target entities. For example, a variety of potential email identifiers may be predicted. These may include, for example, "john.smith", "j.smith", "j_smith", etc.

Exemplary Method 4--Information About Target Entities may be Extracted by Analyzing Structural Characteristics Within an Activity Record, Thus Enabling the Identification of Target Entities Based on Typical,

anonymization dictionary in the anonymization correspondence repository. Completion of Step 210 marks the completion of the second pass, i.e., the anonymization of the activity record.

[0069] In one embodiment of the invention, the anonymization engine may rank entity types in a particular order. The ranking may be based on the degree of information a particular entity type provides. For example, the entity type "user name" may provide more information about a target type than the entity type "email address" because user names are detected only for users that have a user account allowing them to log on to a computing device of the company, whereas an email address may be an email address of a company-internal or external user. Accordingly, the information provided by a target entity "user name" may be more specific than the information provided by a target entity "email address". In one embodiment of the invention, the highest-ranked entity type of the entry may be inserted along with the anonymized identifier, rather than the entity type associated with the target entity being replaced. Consider, for example, a scenario where an anonymization dictionary entry exists that includes a user name and four email addresses. Accordingly, five target entities exist (one for the user name, and four for the email addresses). When a target entity is replaced using this anonymization dictionary entry, the target is replaced by the anonymized identifier and the entity type "user name", even if the replaced target entity is of the entity type "email address", based on the highest ranking of the entity type "user name".

[0070] For illustrative examples of the replacement of target entities with the corresponding anonymized identifiers and entity types, see FIG. 4 and the associated description below.

[0071] Continuing with the discussion of FIG. 2A, in Step 212, a determination is made about whether target entities to be anonymized are remaining. Target entities may be remaining if the anonymization dictionary did not include the entries necessary to resolve the remaining target entities. This may occur, for example, if the anonymization correspondence repository is based on an initial activity record or set of activity records, and when a newly received activity record that includes target entities that do not exist in the anonymization correspondence repository, is being processed. If target entities are remaining, the method may proceed to Step 214.

[0072] In Step 214, the anonymization correspondence repository is updated with the remaining target entities. The addition of an entry is described in detail below, with reference to FIGS. 3A and 3B.

[0073] Continuing with the discussion of FIG. 2A, in Step 216, the remaining target entities to be anonymized are replaced with the corresponding anonymized identifiers which were added to the anonymization correspondence repository in Step 214.

[0074] In certain scenarios, the target entities to be anonymized may only be remaining after the execution of Step 210 in cases where the method described in FIG. 2A is used to anonymize an activity record using an anonymization correspondence repository that was generated based on an earlier activity record. In embodiments of the invention where separate anonymization correspondence repositories are generated for subsequently received activity records, Steps 212-216 may not apply, i.e., the method may always proceed from Step 210 directly to Step 218.

[0075] Returning to Step 212, if a determination is made that no target entities to be anonymized are remaining, the method may proceed to Step 218.

[0076] In Step 218, the anonymized activity record is stored. In one or more embodiments of the invention, the anonymized activity record may be stored to the previously described anonymized activity record repository. The threat analysis engine may access the anonymized activity records stored in the anonymized activity record repository in order to analyze the stored activities for indications of threats, in accordance with one or more embodiments of the invention.

[0077] The following description of FIG. 2B covers an embodiment of the invention where the generation of the dictionary and the anonymization of the activity record(s) may be performed subsequently in a single pass over

name", and "other". Different methods may be employed to determine the entity type of a target entity. The following lists exemplary methods that may be used to identify the entity type of the target entity.

Exemplary Method 1

[0109] For the entity type "user name", the determination may be made based on information provided by the administrator. For example, the administrator may indicate that "John Smith" and "Joe Smith" are user names. Alternatively, and/or in addition, user names may be obtained by parsing the local directories of servers and/or user's computing devices for user profiles such as the "\Users" directory, and/or the "\Documents and Settings" directory.

Exemplary Method 2

[0110] For the entity type "email address", the determination may be made based on the characteristic format of an email address (email_identity@email_provider).

Exemplary Method 3

[0111] For the entity types "company name", and other strings deemed sensitive data (e.g. name of collaboration partners, i.e., partner names), the determination may be made based on information provided by the administrator, i.e., a list of terms including the associated entity type may be provided.

Exemplary Method 4

[0112] For the entity type "domain name" and "IP address", the determination may be made based on the characteristic formats of domain names and IP addresses. For example, an IPv4 address may be represented by 4 bytes separated by the "." sign. Domain names may be represented in a format such as "domain\username" or "username @ domain".

Exemplary Method 5

[0113] For the entity type "host name", the determination may be made based on the endpoint agents using a specific field for documenting the host name when generating an activity record. Those skilled in the art will appreciate that the invention is not limited to the aforementioned methods for determining entity types associated with target entities.

[0114] Continuing with the discussion of FIG. 3A, in Step 308, an instance ID is determined for the target entity. As previously discussed with reference to FIG. 1C, the instance IDs assigned to entries generated for different target entities may be unique within a resource profile, or they may be globally unique within the anonymization dictionary. In one embodiment of the invention, the instance ID may be included in the anonymized ID. In this case, the anonymized identifier may not include a separate entry dedicated to the instance ID. Further, in an alternative embodiment of the invention where a matching resource profile is not resolved for target entities to be added to the anonymization dictionary, there may be no instance IDs, and Step 308 may therefore be skipped.

[0115] In Step 310, the anonymized ID (assigned in Step 304), the entity type (assigned in Step 306), and the instance ID (assigned in Step 308, if used by the implementation of the anonymization dictionary) are combined to form an anonymized identifier, previously described with reference to FIG. 1C.

[0116] In Step 312, an entry including the target entity and the corresponding anonymized ID, established in Step 310, are stored in the anonymization dictionary.

[0117] In Step 314, a determination is made about whether target entities to be added to the anonymization dictionary are remaining. If target entities to be added to the anonymization dictionary are remaining, the method may return to Step 300.

[0118] FIG. 3B describes a method for generating an anonymization correspondence repository, in accordance with one or more embodiments of the invention, where an equivalence map is used to group entries for target entities that can be tracked back to a single resource (e.g. the same user, or the same company). The equivalence map may be included in the anonymization correspondence repository, in addition to the anonymization dictionary.

[0119] Turning to FIG. 3B, in Step 350, a target entity to be added to the anonymization dictionary in the anonymization correspondence repository is selected. If the method described in FIG. 3B is called from the method described in FIG. 2A, multiple target entities to be added to the anonymization dictionary may exist. In this case, the first target entity may be selected initially, and after completion of Steps 352-358 for the first target entity, the next target entity may be selected, etc. If the method described in FIG. 3B is called from the method described in FIG. 2B, only one target entity to be added to the anonymization dictionary may exist. In this case, Steps 352-358 may only be executed once.

[0120] In Step 352, an anonymized identity (ID), previously described in detail with reference to FIG. 1C, is generated for the target entity to be added to the anonymization dictionary. A new anonymized ID, different from all other anonymized IDs in the anonymization dictionary, may be assigned.

[0121] In Step 354, the entity type is determined for the target entity. The entity type for a detected target entity may be, for example, "user name", "email address", "company name", "domain name", "IP address", "host name", and "other". Different methods, as previously described with reference to Step 306 of FIG. 3A may be employed to determine the entity type of a target entity.

[0122] In Step 356, the anonymized ID (assigned in Step 352) and the entity type (assigned in Step 354) are combined to form an anonymized identifier, previously described with reference to FIG. 1D.

[0123] In Step 358, an entry including the target entity and the corresponding anonymized ID, established in Step 356, are stored in the anonymization dictionary.

[0124] In Step 360, a determination is made about whether target entities to be added to the anonymization dictionary are remaining. If target entities to be added to the anonymization dictionary are remaining, the method may return to Step 350. If a determination is made that no target entities to be added to the anonymization dictionary are remaining, the method may proceed to Step 362.

[0125] In Step 362, an equivalence map is generated in the anonymization correspondence repository. The equivalence map may be used to group entries for target entities and the corresponding anonymization identifiers in the anonymization dictionary based on resources with which the target entities are affiliated. First, for each entry in the anonymization dictionary the methods described in Step 302 of FIG. 3A may be used in order to determine a resource associated with the target entity of the entry. Those skilled in the art will appreciate that the invention is not limited to these aforementioned methods for identification of a resource associated with a target entity. Further, the aforementioned methods may be used in combination to identify a resource associated with a target entity. Next, once the resources associated with the target entities have been identified, an identity relationship may be established in the equivalence map for each resource with multiple associated entries in the anonymization dictionary, thus linking entries for target entities that are associated with the resource (e.g. a user or a legal entity). The identity relationships may then be stored in the equivalence map.

[0126] In one embodiment of the invention, the equivalence map may be generated subsequent to completion of the anonymization dictionary for an activity record or a set of activity records. Alternatively, the equivalence map may be generated asynchronously, for example in a background process, at scheduled times, caused by certain trigger events, e.g., upon request by an application or analysis engine that requires information stored in the equivalence map, and/or under certain conditions, e.g. during times of low system load. Further, the equivalence map may be updated at regular intervals or upon availability of additional information that may allow the identification of a resource associated with a particular target entity, where the information previously

available was not sufficient to perform the identification of the associated resource. In another embodiment of the invention, no equivalence map is generated, i.e. Step 362 is skipped entirely.

[0127] FIG. 4 shows an example of an anonymization of an activity record performed in accordance with one or more embodiments of the invention. The example shown in FIG. 4 is not intended to limit the invention. The upper panel shows the sample activity record to be anonymized, and the lower panel shows the resulting anonymized activity record. In the activity record to be anonymized (upper panel) all target entities detected by the anonymization engine are marked by dashed rectangles. A total of seven target entities were detected: The first target entity is the user name "mark.thawley". The second target entity is the email address "pavel.sherbakov@ggmmaaiill.com", i.e. the email identity is "pavel.sherbakov". The third target entity is again the user name "mark.thawley". The fourth target entity is again the email address "pavel.sherbakov@ggmmaaiill.com". The fifth target entity is the domain name "DS\WIN81RS". The sixth target entity is the IP address "192.168.1.20". The seventh target entity is again the user name "mark.thawley".

[0128] In the example shown in FIG. 4, the anonymization engine is configured to anonymize all of the above target entities. The format of the entries replacing the target entities after anonymization is [{"1-entity type-anonymized identifier"}]. The anonymized identifier for the user mark.thawley is "9", with the entity type being "USER", thus indicating a user with a corporate user account. The anonymized identifier for email identity "pavel.sherbakov" is "5", with the entity type "EMAIL", indicating an email identifier. Note that the entity type would have been "USER" for the email identity "pavel.sherbakov" if a corporate user account under that name existed. The anonymized identifier for the domain name "DS\WIN81RS" is "6" with the entity type being "DOMAIN", indicating a domain name. The anonymized identifier for the IP address "192.168.1.20" is "2" with the entity type being "AIP", indicating an IP address.

[0129] Accordingly, the target entities in the activity record are replaced as shown in the anonymized activity record (lower panel) of FIG. 4.

[0130] In one or more embodiments of the invention, anonymized activity records, obtained using the methods described above, may be de-anonymized in order to identify resources associated with anonymized identifiers in the anonymized activity records. A de-anonymization may be necessary, for example, when the threat analysis engine issues a threat alert. In order to assess the threat and to take threat-mitigating action, it may be necessary to identify the resource associated with the activities in the anonymized activity records that have caused the alert. In one or more embodiments of the invention, the anonymization dictionary may be used to perform the de-anonymization, i.e., the translation from anonymized identifier(s) to resource(s). The de-anonymization may be performed as needed, i.e., depending on the nature of the threat alert, the de-anonymization may only be partially performed, for example for certain entity types. For example, if the threat appears to originate from an employee, only user names may be de-anonymized. In another case, where the nature of the threat alert is unknown, all entity types may need to be de-anonymized, i.e., a complete de-anonymization of all anonymized identifiers in an entire anonymized activity record may be performed.

[0131] Embodiments of the invention may enable the anonymization of activity records, thereby protecting resource identities, without limiting the analysis of the activities included in the activity records for threats. Threat detection may therefore be performed using the anonymized activity records. Further, the anonymization of the activity records, by classifying the entity type for each target entity being anonymized, adds structure to the activity records, which may improve and/or facilitate threat detection. In addition, various configurable levels of anonymization may allow to the activity monitoring to satisfy requirements of various different scenarios, ranging from anonymization used for company-internal threat monitoring to anonymization prior to publicly sharing the anonymized activity records, where they may be used by developers to design and validate threat detection algorithms and anonymization engines. Overall, the anonymization, performed in accordance with embodiments of the invention, may facilitate compliance with company-internal guidelines and national laws requiring the protection of user identities.

[0132] While the invention has been described with respect to a limited number of embodiments, those skilled in the art, having benefit of this disclosure, will appreciate that other embodiments can be devised which do not

