



(12) **DEMANDE DE BREVET CANADIEN  
CANADIAN PATENT APPLICATION**

(13) **A1**

(86) **Date de dépôt PCT/PCT Filing Date:** 2022/03/14  
 (87) **Date publication PCT/PCT Publication Date:** 2022/09/22  
 (85) **Entrée phase nationale/National Entry:** 2023/09/11  
 (86) **N° demande PCT/PCT Application No.:** US 2022/020189  
 (87) **N° publication PCT/PCT Publication No.:** 2022/197606  
 (30) **Priorités/Priorities:** 2021/03/13 (US63/160,780);  
 2021/03/18 (US63/162,829)

(51) **Cl.Int./Int.Cl. H04L 9/40** (2022.01)  
 (71) **Demandeur/Applicant:**  
 DIGITAL REASONING SYSTEMS, INC., US  
 (72) **Inventeurs/Inventors:**  
 ACUFF, JENA, US;  
 CARL, BRANDON, US;  
 KAMATH, UDAY, US;  
 HUGHES, CORY, US  
 (74) **Agent:** MARKS & CLERK

(54) **Titre : ACTIONNEMENT D'ALERTE ET RETROACTION D'APPRENTISSAGE AUTOMATIQUE**  
 (54) **Title: ALERT ACTIONING AND MACHINE LEARNING FEEDBACK**

(57) **Abrégé/Abstract:**

Some aspects of the present disclosure relate to systems, methods, and computer-readable media for configuring a conduct surveillance system. In one example implementation, a computer implemented method includes: receiving at least one alert from a conduct surveillance system, where the at least one alert represents a potential violation of a predetermined policy, where the predetermined policy includes a scenario, a target population, and a workflow; determining whether each of the at least one alert represents an actual violation of the predetermined policy; calculating a metric based on the actual violations and the potential violations where the metric includes a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and changing at least one of the scenario, the target population, or the workflow based on the calculated metric.

**Date Submitted:** 2023/09/11

**CA App. No.:** 3211747

**Abstract:**

Some aspects of the present disclosure relate to systems, methods, and computer-readable media for configuring a conduct surveillance system. In one example implementation, a computer implemented method includes: receiving at least one alert from a conduct surveillance system, where the at least one alert represents a potential violation of a predetermined policy, where the predetermined policy includes a scenario, a target population, and a workflow; determining whether each of the at least one alert represents an actual violation of the predetermined policy; calculating a metric based on the actual violations and the potential violations where the metric includes a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and changing at least one of the scenario, the target population, or the workflow based on the calculated metric.

## ALERT ACTIONING AND MACHINE LEARNING FEEDBACK

### CROSS-REFERENCE TO RELATED APPLICATION

This application claims priority to and benefit of U.S. provisional patent application serial  
5 no. 63/160,780 filed March 13, 2021 and U.S. provisional patent application serial no. 63/162,829  
filed March 18, 2021, which are hereby fully incorporated by reference and made a part hereof.

### BACKGROUND

The present disclosure generally relates to monitoring communications for activity that  
10 violates ethical, legal, or other standards of behavior and poses risk or harm to institutions or  
individuals. The need for detecting violations in the behavior of representatives of an institution  
has become increasingly important in the context of proactive compliance, for instance. In the  
modern world of financial services, there are many dangers to large institutions from a compliance  
perspective, and the penalties for non-compliance can be substantial, both from a monetary  
15 standpoint and in terms of reputation. Financial institutions are coming under increasing pressure  
to quickly identify unauthorized trading, market manipulation and unethical conduct within their  
organization, for example, but often lack the tools to do so effectively.

Moreover, systems and methods for monitoring communications can be tuned or adapted  
to detect types of violations of behavior, or to increase the accuracy of those systems and methods.  
20 Advanced systems and methods for monitoring communication can be based on complicated  
models, including machine learning models and other techniques. Therefore it can be difficult for  
a user of a system or method to tune or adapt that system or method.

Thus, among other needs, there exists a need for effective identification of activity that  
violates ethical, legal, or other standards of behavior and poses risk or harm to institutions or  
25 individuals from electronic communications. Furthermore, there exists a need for effective ways  
to improve the identification of violation conditions and effective ways to configure systems to  
identify violation conditions. It is with respect to these and other considerations that the various  
embodiments described below are presented.

30

### SUMMARY

Embodiments of the present disclosure are directed generally towards methods, systems,  
and computer-readable storage medium relating to, in some embodiments, an intuitive review and

investigation tool designed to facilitate efficient and defensible reviews of electronic communications (e.g., messages) by analysts (sometimes referred to herein as “users”). In certain implementations, users are empowered to quickly target risk areas and increase accuracy of reviews as a result of a flexible workflow for evaluating alerts. Users can also have a more granular  
5 feedback loop for alerts as an input into reporting and model training (see FIG. 13, for example). Actioning can be performed on communications at the hit, alert, and/or message level to illustrate progress of reviewed communications and streamline business processes (see FIG. 13, for example).

A streamlined actioning workflow can allow users to easily close alerts and add relevant  
10 context (for example, person of interest and comments) to elevated alerts requiring further review (see FIG. 25, for example). Alerts assigned to a user can be accessed from a user dashboard, where the user can also see the total messages awaiting review (see FIG. 26, for example).

In one aspect, the present disclosure relates to a computer-implemented method, which, in one embodiment, receiving at least one alert from a conduct surveillance system, where the at least  
15 one alert represents a potential violation of a predetermined standard and where the conduct surveillance system generates the alerts in response to an electronic communication between persons matching a violation of a predetermined policy, where the predetermined policy includes a scenario, a target population, and a workflow; determining whether each of the at least one alert represents an actual violation of the predetermined policy; calculating a metric based on the actual  
20 violations and the potential violations where the metric includes a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and changing at least one of the scenario, the target population, or the workflow based on the calculated metric.

In some embodiments of the present disclosure, the scenario includes a machine learning  
25 classifier, and where determining whether the at least one alert represents an actual violation includes labeling the at least one alert and using the labeled at least one alert to train the machine learning classifier.

In some embodiments of the present disclosure, the metric is displayed to a user.

In some embodiments of the present disclosure, the scenario includes a lexicon, and where  
30 the lexicon represents one or more terms or regular expressions.

In some embodiments of the present disclosure, changing the scenario includes changing the lexicon by adding or removing terms or regular expressions from the lexicon.

In some embodiments of the present disclosure, the computer implemented method includes, in response to determining that the at least one alert represents an actual violation, actioning the alert.

5 In some embodiments of the present disclosure, actioning the alert includes receiving a user input from the user interface representing whether the at least one alert represents an actual violation.

In some embodiments of the present disclosure, the target population includes a domain exclusion list and where changing the target population includes changing the domain exclusion list.

10 In some embodiments of the present disclosure, the electronic communication includes metadata, the scenario includes rules for filtering the electronic communication based on the metadata, and where changing the scenario includes changing the rules for filtering the electronic communications based on the metadata.

15 In another aspect, the present disclosure relates to a system, which in one embodiment includes: at least one processor; at least one memory storing computer readable instructions configured to cause the at least one processor to perform functions for creating and/or evaluating models, scenarios, lexicons, and/or policies, where the functions include: receiving data associated with at least one of text data, model training, lexicons, scenarios, and policies, where the functions for creating and/or evaluating models comprise creating at least one scenario based  
20 on at least one of the models, lexicons, and non-language features; creating one or more policies mapping to the at least one scenario and a population; upon receiving an alert that a policy match occurs, triggering an alert indicating, to a user, that a policy match has occurred which requires a user action, where a policy corresponds to actions that violate at least one of a combination of signals and metrics, a population, and workflow.

25 In some embodiments of the present disclosure, the model training includes training at least one model configured to analyze the text data from one or more electronic communications between at least two persons.

In some embodiments of the present disclosure, the user action includes review and interaction by a user via a user interface.

30 In some embodiments of the present disclosure, the model training includes evaluating the model against established datasets.

In some embodiments of the present disclosure, the alert to the user is evaluated by the user and a corresponding user decision is made to confirm or deny accuracy of the alert.

In some embodiments of the present disclosure, the user decision is provided into a feedback loop, and where the feedback loop is configured to improve the model training.

In some embodiments of the present disclosure, the user decision is provided into the feedback loop and where the feedback loop is configured to improve the lexicons, scenarios, or policies.

In some embodiments of the present disclosure, the feedback loop is configured to change a lexicon.

In some embodiments of the present disclosure, changing the lexicon includes configuring the lexicon so that it includes or excludes terms or regular expressions.

In some embodiments of the present disclosure, the feedback loop is configured to measure the rate of false positives and to change one or more of the lexicons, scenarios, and policies based on the rate of false positives.

In some embodiments of the present disclosure, the scenario includes Boolean operators, and where the feedback loop is configured to change one or more of the Boolean operators.

In some embodiments of the present disclosure, the feedback loop is configured to monitor the rate of false positives over a period of time, and change one or more of the lexicons, scenarios, and policies based on the rate of false positives over the period of time.

In yet another aspect, the present disclosure relates to a non-transitory computer-readable medium storing instructions which, when executed by one or more processors, cause a computing device to perform specific functions. The functions performed include receiving at least one alert from a conduct surveillance system, where the at least one alert represents a potential violation of a predetermined standard and where the conduct surveillance system generates the alerts in response to an electronic communication between persons matching a violation of a predetermined policy, where the predetermined policy includes a scenario, a target population, and a workflow; determining whether each of the at least one alert represents an actual violation of the predetermined policy; calculating a metric based on the actual violations and the potential violations where the metric includes a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and changing at least one of the scenario, the target population, or the workflow based on the calculated metric.

The following provides a non-limiting discussion of some example implementations of various aspects of the present disclosure. Some aspects and embodiments disclosed herein may be utilized for providing advantages and benefits in the area of communication surveillance for

regulatory compliance. Some implementations can process all communications, including electronic forms of communications such as instant messaging (or “chat”), email, voice, and/or social network messaging to connect and monitor an organization’s employee communications for regulatory and corporate compliance purposes. Some embodiments of the present disclosure unify  
5 detection, user interfaces, behavioral models, and policies across all communication data sources, and can provide tools for compliance analysts in furtherance of these functions and objectives. Some implementations can proactively analyze users’ actions to identify breaches such as unauthorized activities that are against applicable policies, laws, or are unethical, through the use of natural language processing (NLP) models. The use of these models can enable understanding  
10 content of communications and map signals to behavioral profiles in order to locate high-risk individuals.

Other aspects and features according to the example embodiments of the present disclosure will become apparent to those of ordinary skill in the art, upon reviewing the following detailed description in conjunction with the accompanying figures.

15

### **BRIEF DESCRIPTION OF THE DRAWINGS**

Reference will now be made to the accompanying drawings, which are not necessarily drawn to scale.

FIGS 1A-1C illustrate methods according to various aspects of the present disclosure. FIG.  
20 1A illustrates a method for creating alerts based on a policy match according to one embodiment of the present disclosure. FIG. 1B illustrates a method of configuring a computer system to detect violations in a target dataset according to one embodiment of the present disclosure. FIG. 1C illustrates a method for increasing the accuracy of a conduct surveillance system according to one embodiment of the present disclosure.

25 FIGS. 2A-2C illustrate various aspects of the present disclosure. FIG. 2A illustrates various aspects of displayed events, properties, and communications data, in accordance with one or more embodiments of the present disclosure. FIGS. 2B and 2C illustrate various aspects of policies, including scenario, population, and workflow, in accordance with one or more embodiments of the present disclosure.

30 FIG. 3 is a diagram relating to workflows in accordance with one or more embodiments of the present disclosure.

FIG. 4 illustrates various aspects displayed to a user interface, including elements of a graphical user interface, with a sidebar, content, and aside areas, in accordance with one or more embodiments of the present disclosure.

FIG. 5 illustrates a visual view including various data representations beyond simple text,  
5 in accordance with one or more embodiments of the present disclosure.

FIG. 6 illustrates aspects of knowledge tasks, in accordance with one or more embodiments of the present disclosure.

FIG. 7 illustrates a profile view corresponding to a particular entity, in accordance with one or more embodiments of the present disclosure.

10 FIG. 8 illustrates a profile view and particularly labels an “aside” section of a displayed graphical user interface, in accordance with one or more embodiments of the present disclosure.

FIG. 9 illustrates various aspects of alerts, hits, and actions, in accordance with one or more embodiments of the present disclosure.

15 FIG. 10 illustrates various aspects of alert hit previews and list cards, in accordance with one or more embodiments of the present disclosure.

FIG. 11 illustrates various aspects of metrics and tabs, in accordance with one or more embodiments of the present disclosure.

FIG. 12 is a computer architecture diagram showing a general computing system capable of implementing one or more embodiments of the present disclosure described herein.

20 FIG. 13 is a flow diagram illustrating components and operations of a feedback loop and system in accordance with one embodiment of the present disclosure.

FIG. 14 illustrates various aspects of a model dashboard, in accordance with one or more embodiments of the present disclosure.

25 FIG. 15 illustrates various aspects of a model dashboard with statistics, in accordance with one or more embodiments of the present disclosure.

FIG. 16 illustrates various aspects of lexicon evaluation, in accordance with one or more embodiments of the present disclosure.

FIG. 17 illustrates various further aspects of lexicon evaluation including a confusion matrix, in accordance with one or more embodiments of the present disclosure.

30 FIG. 18 illustrates various further aspects of lexicon evaluation, in accordance with one or more embodiments of the present disclosure.

FIG. 19 illustrates various further aspects of lexicon evaluation, in accordance with one or more embodiments of the present disclosure.

FIG. 20 illustrates various aspects of scenarios in accordance with one or more embodiments of the present disclosure.

FIG. 21 illustrates various aspects of policy administration functionality in accordance with one or more embodiments of the present disclosure.

5 FIG. 22 illustrates a user interface for accessing a repository in accordance with one or more embodiments of the present disclosure.

FIGS. 23A-23F illustrates user interfaces for configuring a scenario in accordance with one or more embodiments of the present disclosure. FIG. 23A illustrates a user interface for viewing one or more datasets. FIG. 23B illustrates a user interface for labeling a dataset. FIG. 23C  
10 illustrates an annotation applied to a dataset and an interface for applying labels to a dataset. FIG. 23D illustrates a user interface for configuring a lexicon to be applied to the dataset. FIG. 23E illustrates a user interface for evaluating a lexicon. FIG. 23F illustrates a scenario created using the lexicon that was configured in the interface shown in FIG. 23E.

FIG. 24 illustrates various aspects of actioning communications in accordance with one or  
15 more embodiments of the present disclosure.

FIG. 25 illustrates various aspects of actioning communications in accordance with one or more embodiments of the present disclosure.

FIG. 26 illustrates various aspects of actioning communications in accordance with one or  
20 more embodiments of the present disclosure.

## DETAILED DESCRIPTION

Although example embodiments of the present disclosure are explained in detail herein, it is to be understood that other embodiments are contemplated. Accordingly, it is not intended that  
25 the present disclosure be limited in its scope to the details of construction and arrangement of components set forth in the following description or illustrated in the drawings. The present disclosure is capable of other embodiments and of being practiced or carried out in various ways.

It must also be noted that, as used in the specification and the appended claims, the singular forms “a,” “an” and “the” include plural referents unless the context clearly dictates otherwise.

30 By “comprising” or “containing” or “including” is meant that at least the named compound, element, particle, or method step is present in the composition or article or method, but does not exclude the presence of other compounds, materials, particles, method steps, even if the other such compounds, material, particles, method steps have the same function as what is named.

In describing example embodiments, terminology will be resorted to for the sake of clarity. It is intended that each term contemplates its broadest meaning as understood by those skilled in the art and includes all technical equivalents that operate in a similar manner to accomplish a similar purpose. It is also to be understood that the mention of one or more steps of a method does  
5 not preclude the presence of additional method steps or intervening method steps between those steps expressly identified. Steps of a method may be performed in a different order than those described herein without departing from the scope of the present disclosure. Similarly, it is also to be understood that the mention of one or more components in a device or system does not preclude the presence of additional components or intervening components between those components  
10 expressly identified.

### Definitions

The following discussion provides some descriptions and non-limiting definitions, and related contexts, for terminology and concepts used in relation to various aspects and embodiments  
15 of the present disclosure.

An “event” can be considered any object with a fixed time, and an event can be observable data that happens at a point in time, for example an email, a badge swipe, a trade (e.g., trade of a financial asset), or a phone call (see also the illustration of FIG. 2A).

A “property” relates to an item within an event that can be uniquely identified, for example  
20 metadata (see also illustration of FIG. 2A).

A “communication” or “electronic communication” can be any event with language content, for example email, chat, a document, social media, or a phone call (see also illustration of FIG. 2A). An electronic communication may also include, for example, audio, SMS, and/or video. A communication may additionally or alternatively be referred to herein as, or with respect to, a  
25 “comm” (or “comms”), message, container, report, or data payload.

A “metric” can be a weighted combination of factors to identify patterns and trends (e.g., a number-based value to represent behavior or intent from a communication). Examples of metrics include sentiment, flight risk, risk indicator, and responsiveness score. A metric may additionally or alternatively be referred to herein as, or with respect to, a score, measurement, or  
30 rank.

A “post” can be an identifier’s contribution within a communication, for example a single email within a thread, a single chat post, a continuous burst of communication from an individual,

or a single social media post (see also illustration of FIG. 2A). A post can be considered as an individual's contribution to a communication.

A "conversation" can be a group of semantically related posts, for example the entirety of an email with replies, a thread, or alternatively a started and stopped topic, a time-bound topic, and/or a post with the other post (replies). Several posts can make up a conversation within a communication.

A "signal" can be an observation tied to a specific event that is identifiable, for example rumor language, wall crossing, or language of interest.

A "policy" can be a scenario applied to a population with a defined workflow. A policy may be, for instance, how a business chooses to handle specific situations, for example as it may relate to ongoing deal monitoring, disclaimer adherence, and/or anti money laundering (AML) monitoring. As used herein, a policy may additionally or alternatively be referred to as, or with respect to, a "KI" or "key indicator", or rules engine. As illustrated in FIGS. 2B and 2C, in some embodiments a policy can be comprised of three items: a scenario as a combination of signals and metrics (as an example of usage, using NLP signals and metrics to discover intellectual property (IP) theft language or behaviors); a population, as the target population over which to look for the scenario (e.g., sales team(s), department(s), or group(s) of persons); and workflow, as actions taken when a scenario triggers over a population (e.g., alert generation).

An "alert" can indicate to a user that a policy match has occurred which requires action (sometimes referred to herein with respect to "actioning" an alert), for example a scenario match. A signal that requires review can be considered an alert. As an example, an indication of intellectual property theft may be found in a chat post with language that matches the scenario, on a population that needs to be reviewed.

A "manual alert" can be an alert added to a communication from a user, not generated from the system. A manual alert may be used, for example, when a user needs to add an alert to language or other factors for further review.

A "hit" can be an exact signal that applies to a policy on events, for example an occurrence of the language "I'm taking clients with me when I leave", a behavior pattern change, and/or a metric change. As used herein, a hit may additionally or alternatively be referred to herein as, or with respect to, a "KI" ("key indicator"), event, and/or highlight.

A "review" can be the act of a user assigning actions on hits, alerts, or communications.

A "tag" can be a label attached to a communication for the purpose of identification or to give other information, for example a new feature set that will enable many workflow practices.

A “knowledge graph” can be a representation of all of the signals, entities, topics, and relationships in a data set in storage. Knowledge graphs can communications, some of which may contain alerts for a given policy. Other related terms may include a “knowledge base.” In some embodiments, a knowledge graph can be a unified knowledge representation.

5 A “personal identifier” can be any structured field that can be used to define a reference or entity, for example “jeb@jebbush.com”, “@CMcK”, “EnronUser1234”, or “(555) 336-2700” (i.e., a personal identifier can include email, a chat handle, or a phone number). As used herein, a hit may additionally or alternatively be referred to herein as, or with respect to, an “entity ID”.

10 A “mention” can be any descriptive string that is able to be referenced and/or extracted, for example “He/Him”, “The Big Blue”, “Enron”, or “John Smith”. Other related terms may include “local coreference.”

An “entity” can be an individual, object, and/or property IRL, and can have multiple identifiers or references, for example John Smith, IBM, or Enron. Other related terms may include profile, participant, actor, and/or resolved entity.

15 A “relationship” can be a connection between two or more identifiers or entities, for example “works in” department, person-to-person, person-to-department, and/or company-to-company. Other related terms may include connections via a network graph.

The following discussion includes some descriptions and non-limiting definitions, and related contexts, for terminology and concepts that may particularly relate to workflows in accordance to one or more embodiments of the present disclosure, some of which may be further understood by reviewing the diagram of FIG. 3.

25 A “smart queue” can be a saved set of search modifiers with an owner and defined time, for example, a daily bribery queue, an action pending queue, an escalation queue, or any shared/synced list. As used herein, a smart queue may additionally or alternatively be referred to herein as, or with respect to an action pending queue, analyst queue, or scheduled search.

A “saved search” can be a saved set of search modifiers with no owner, for example a monthly QA check, an investigation search, or an irregularly used search. As used herein, a saved search may additionally or alternatively be referred to herein as, or with respect to a search copy or a bookmark.

30 The following discussion includes some descriptions and non-limiting definitions, and related contexts, for terminology and concepts that can relate to a graphical user interface (and associated example views as output to a user) that can be used by a user to interact with, visualize,

and perform various functionalities in accordance to one or more embodiments of the present disclosure.

A “sidebar” can be a global placeholder for navigation and branding (see, e.g., illustrations in FIG. 4.

5 “Content” as shown and labeled in, for example, FIG. 4, identifies where primary content will be displayed.

An “aside” as shown and labeled in, for example, FIG. 4, is a location for supportive components that affect the content or provide additional context. Further related aspects of “aside” are shown in the example of FIG. 8. An aside can be a column of components that support, define,  
10 or manipulate the content area.

A “visual view” as illustrated in, for example, FIG. 5, can include a chart, graph, or data representation that is beyond simple text, for example communications (“comms”) over time, alters daily, queue progress, and/or relationship metric(s). As used herein, visual views may additionally or alternatively be referred to herein as, or with respect to charts or graphs.

15 A “profile” can be a set of visuals filtered by an identifier or entity, for example by a specific person’s name, behavior analytics, an organization’s name, or QA department. As used herein, profiles may additionally or alternatively be referred to herein as, or with respect to relationship(s) or behavior analytics.

Now also referring to the diagram of FIG. 6, smart queues can enable teams to work to  
20 accelerate “knowledge tasks”. Signals that require review (i.e., alerts), comprise monitoring. These can be from external systems. Knowledge tasks can provide feedback via a “learning loop” into models.

Now also referring to the view in the illustration of FIG. 7, a particular profile view can provide insights such as behavioral insights to, for instance, an entity (here, a particular person).  
25 The profile can include a unified timeline with hits, and communications. Also, profiles can provide aggregates of/into entities, metrics, visuals, events, and relationships. As mentioned briefly above and as illustrated in FIG. 8, an aside can be a column of components that support, define, or manipulate the content area.

Now referring to the view in the illustrations of FIGS. 9 and 10, and as discussed in some  
30 detail above, an “alert” can be the manifestation of a policy on events, and a “hit” (or “alert hit”) can be the exact signal that applies to a policy on events. An “action” can be the label that is applied to: a single hit; all hits under an alert; or all hits on a message. A “list card” can be an

object that contains a summary of the content of a comm in the “list view”, which can be a list of events with communications that may have an alert.

Now referring to the view in the illustration of FIG. 11, as discussed in some detail above, a “metric” can be a weighted combination of factors to identify patterns and trends. A “tab” can be an additional view that can display content related to a current view, for example sibling content.

The following discussion includes some descriptions and non-limiting definitions, and related contexts, for terminology and concepts that may particularly relate to machine learning models and the training of machine learning models, in accordance with one or more embodiments of the present disclosure.

10 A “hit” can be an exact signal that applies to a policy on events, for example an occurrence of the language “I’m taking clients with me when I leave”, a behavior pattern change, and/or a metric change. As used herein, a hit may additionally or alternatively be referred to herein as, or with respect to, a “KI” (“key indicator”), event, and/or highlight.

A “pre-trained model” can be a model that performs a task but requires tuning (e.g., supervision and/or other interaction by an analyst or developer) before production. An “out of the box model” can be a model that benefits from, but does not require, tuning before use in production. Pre-trained models and out of the box models can be part of the building blocks for a policy. As used herein, a pre-trained model may additionally or alternatively be referred to herein as, or with respect to, “KI engines” or “models”.

20 In some embodiments, the present disclosure can provide for implementing analytics using “supervised” machine learning techniques (herein also referred to as “supervised learning”). Supervised mathematical models can encode a variety of different data aspects which can be used to reconstruct a model at run-time. The aspects utilized by these models may be determined by analysts and/or developers, for example, and may be fixed at model training time. Models can be retrained at any time, but retraining may be done more infrequently once models reach certain levels of accuracy.

### **Description of Example Embodiments of Present Disclosure**

A detailed description of various aspects of the present disclosure, in accordance with various example embodiments, will now be provided with reference to the accompanying drawings. The drawings form a part hereof and show, by way of illustration, specific embodiments and examples.

The following provides a non-limiting discussion of some example implementations of various aspects of the present disclosure

In some embodiments, the present disclosure is directed to a system for indicating to a user when a policy match has occurred which requires action by the user. The system can include a processor and a memory configured to cause the processor to perform functions for creating and/or evaluating models, scenarios, lexicons, and/or policies. As a non-limiting example, the processor and memory can be part of the general computing system illustrated in FIG. 12.

Embodiments of the present disclosure can implement the method illustrated in FIG. 1A. The instructions stored on the memory can include instructions to receive 102 data associated with text data, model training, lexicons, scenarios and/or policies. Creating and/or evaluating models can include creating a scenario based on the models, lexicons, and non-language features. It should be understood that the scenario can be based on any combination of models, lexicons, and non-language features. As a non-limiting example, the scenario can be based on a single model, but multiple lexicons and multiple non-language features.

As described herein, the model can correspond to a machine learning model. In some embodiments, the machine learning model is a machine learning classifier that is configured to classify text. Additionally, in some embodiments, the model training can include training models for analysis of text data from one or more electronic communications between at least two persons.

The present disclosure contemplates the machine learning training techniques known in the art can be applied to the data disclosed in the present disclosure for model training. For example, in some embodiments, the model training can include evaluating the model against established datasets. As another example, the model training can be based on a user input, for example a user input that labels the data.

The system can be configured to create 104 one or more policies mapping to the scenario and a population. In embodiments with more than one scenario and/or more than one policy, it should be understood that any number of scenarios and/or policies can be mapped to one another. As non-limiting examples, the system can be configured to map multiple scenarios to multiple policies, or multiple scenarios to the same policy or policies.

When the system receives an alert that a policy match occurs, the system can trigger 106 an alert indicating, to a user, that a policy match has occurred which requires action. The policy can correspond to actions that violate at least one of a combination of signals and metrics, a population, and workflow (referred to herein as a “violation”)

Additionally, the present disclosure contemplates that the alerts can be reviewed by the user or by a machine learning model. This review can include determining whether the alerts correspond to an actual violation, and can be used to change the scenario, or change any of the parts of the scenario (e.g. models, lexicons, and non-language features).

5 In some embodiments of the present disclosure, a user can review the data and perform an interaction using a user interface (e.g., a graphical user interface that is part of or operably connected to the computer system illustrated in FIG. 12). The action can include review and interaction by a user via a user interface, which is optionally part of the computing device in FIG. 12. As a non-limiting example, in some embodiments, the system can provide the alert to the user  
10 through the user interface, and then the user can confirm or deny the accuracy of the alert using the user interface. Based on the user input, the system can determine whether the alert was a true positive, true negative, false positive, or false negative. The system can use the information about the alerts, including whether the alert was a true positive, true negative, false positive, or false negative, as an input into the system to improve the operation of the system. This can be referred  
15 to as “feedback.” The present disclosure contemplates that the feedback can be an input into the machine learning model to improve the model training (e.g. the information about the alerts is “fed back” into the model to train the model further). Alternatively or additionally, the present disclosure contemplates that the feedback can be used to change other parameters within the scenario. For example, the feedback can be used to adjust the lexicon or non-language features of  
20 the scenario. This can include adding or removing terms from the lexicon, or adding/removing non-language features from the scenario.

As a non-limiting example, a scenario has a pre-trained machine learning model, a target lexicon of regular expressions and text, and a target set of non-language features that includes metadata. In this example, the scenario can be configured to identify communications that  
25 correspond to the machine learning model and lexicon, where the metadata shows that the communication is from a time span of the previous two years. The system can then produce alerts by determining whether each of the communications in the dataset is a policy match with the scenario. The user can review the communications that are a policy match with the scenario, and determine whether each communication is a violation, and input those results into the system.  
30 Then, based on those results, the system can be configured to change the scenario to improve the effectiveness of the scenario. This can include maximizing or improving certain measures of accuracy such as the ROC curve described herein, the true positive rate, precision, recall, or confusion matrix. As a non-limiting example, this can include changing the scenario to target

metadata in a shorter timeframe, e.g., by changing it from two years to one year. The system and/or the user can then use one or more of the measures of accuracy (e.g., the true positive rate) to see if the measure of accuracy has improved after changing the scenario. By monitoring the accuracy of the scenario as the scenario is changed, it is possible to tune the scenario to improve the measures of accuracy. Again, these are merely non-limiting examples of techniques for measuring the error rate, and it will be understood to one of skill in the art that any techniques for measuring error rate that are known in the art can be used in combination with the system and methods disclosed herein.

Embodiments of the present disclosure can also include computer implemented methods for configuring a computer system to detect violations in a target dataset. With reference to FIG. 1B, the method 120 can include receiving 122 data associated with an electronic communication. The received data can include text data, and optionally metadata that are associated with one or more communications. As a non-limiting example, the data can include a set of emails, text messages, transcribed phone conversations, or combinations thereof. This data can also include “metadata” that can correspond to any information about the communication that is not found in the text itself.

At step 124, the received data can be labeled. As described throughout the present disclosure, labeling can include applying a label indicating whether the one or more communications that are part of the data correspond to a violation. Labeling can also include determining whether the received data includes a segment of target language, and applying a label to the parts of the data that contain that segment of target language. As a non-limiting example, this can include labeling certain communications in the dataset that contain the target language.

At step, 128, a machine learning model can be created based on the data. As described elsewhere in the present disclosure, this machine learning model can be a machine learning classifier that is configured to classify text. As a non-limiting example, the present disclosure contemplates that the model training can include evaluating the model against established datasets. As another non-limiting example, the model training can include training at least one model configured to analyze text data from one or more electronic communications between at least two persons. Additionally, it should be understood that the machine learning model can be any of the other machine learning models described herein, or known in the art.

At step 126, a lexicon can be created for the scenario. As described throughout the present disclosure, the lexicon can represent one or more terms or regular expressions. Optionally, at step

126, the lexicon can be imported partially or completely from a database, or chosen from a list of pre-generated lexicons by a user.

At step 130, a scenario can be created using the machine learning models and the lexicon, where the scenario can represent a violation condition (e.g. a violation of an ethical policy, regulatory policy, rule, law etc., as described in the other examples herein). The user can create  
5 the scenario by specifying the model or models that are used, as well as the lexicon or lexicons that are used.

In some embodiments, the scenario can be created 130 using components other than just a machine learning model and lexicon. For example, the scenario can include a filter, where the filter  
10 can be configured to exclude or include at least part of the dataset based on the data in the dataset. This can include filtering based on data such as metadata. Again, it should be understood that metadata can refer to any of the properties of a communication that are stored in the data, non-limiting examples of which are the time sent, time received, type of communication, etc.

The user or system can also specify how the models and lexicons are joined together.  
15 Again, as a non-limiting example, the scenario can combine one or more models and lexicons using Boolean logic (e.g. AND, OR, NOT, NOR). It should be understood that other logical systems and other logical operators can be used in combination with the method disclosed herein.

Optionally, in some embodiments, the scenario can be created based on feedback from actions the user has taken in response to previous alerts (described herein as “actioning” the alerts).  
20 This can include providing a user decision or user action into a feedback loop that is configured to improve the model training. As a non-limiting example, this user decision can include confirming or denying the accuracy of the alert. In some embodiments, the feedback loop can be configured to improve the lexicons, scenarios, or policies. As yet another non-limiting example, the feedback loop can be configured to change the lexicon, and changing the lexicon can include changing the  
25 lexicon so that it includes or excludes terms or regular expressions. As another non-limiting example, the scenario can include one or more Boolean operators, and the feedback loop can be configured to change one or more of those Boolean operators. Furthermore, in some embodiments of the present disclosure, the feedback loop can be configured to measure the rate of false positives between the actual and potential violations identified by the system, and change one or more of  
30 the lexicons, scenarios, and policies based on the rate of false positives. The feedback loop can also be configured to measure the rate of false positives over a period of time, and change one or more of the lexicons, scenarios, and policies based on the rate of false positives over the period of time.

It should be understood that the rate of false positives is intended only as a non-limiting example, and that the feedback loop can be configured to change the scenario, lexicons, and policies based on other measurements of error, accuracy, etc. As a non-limiting example, example, based on the actioning, the system can be configured to add or remove lexicons or models from  
5 the scenario.

At step 132, the computer system (e.g. the computer system of FIG. 12) can be configured to detect violation conditions in a target dataset using the scenario. This can include storing the scenario in a computer readable medium, receiving additional data for review, and determining whether the additional data contains communications that match the scenario (i.e. that are a “policy  
10 match”).

In some implementations, the scenario can be configured to allow for a user to easily configure the scenario. The system can be configured to prevent a user from changing the machine learning model, but enable the user to change parameters other than the model. This can allow the user to change the scenario and the type of communications identified by the scenario, without  
15 requiring knowledge of the machine learning model, or requiring that the model undergo retraining before use. In some embodiments of the present disclosure, techniques that can be used to reduce the error rates or increase the accuracy other than changing the model itself can be referred to as the “augmentation layer.” Non-limiting examples of techniques that can be included in the augmentation layer include lexicons, domain exclusion lists, and rules-based filter using metadata  
20 (e.g., filtering out alerts based on number of participants or message directionality). The present disclosure contemplates that any or all of the techniques in the augmentation layer can be adjusted based on the dataset.

Furthermore, the present disclosure contemplates that the scenario can be stored in a computer readable medium, for example the memory illustrated in FIG. 12. Similarly, in some  
25 embodiments of the present disclosure, more than one scenario can be stored in one or more computer readable medium. The one or more scenarios can be compared to one another, and the system can create an output based on the comparison. As a non-limiting example, the output based on the comparison can show what parts of the scenario are different, or what parts of the scenario have stayed the same, between the two scenarios. As a non-limiting example, this could include  
30 displaying that two scenarios include the same lexicon, but include different models, or different Boolean operators. The output including the difference between the first and second scenario can also include information about the versions of the two scenarios.

Additionally, some embodiments of the present disclosure are directed to a computer-implemented method 140 for increasing the accuracy of a conduct surveillance system.

With reference to FIG. 1C, the method can include receiving 142 at least one alert from a conduct surveillance system. As used in the present disclosure, a “conduct surveillance system” can refer to a tool for reviewing and investigating communications. Again, the alerts can represent a potential violation of a predetermined standard. The conduct surveillance system can generate the alerts in response to an electronic communication between persons matching a violation of a predetermined policy. As described in greater detail elsewhere in the present disclosure, the predetermined policy can include a scenario, a target population, and a workflow.

In some embodiments of the present disclosure, the scenario can include a machine learning classifier. Additionally, in some embodiments of the present disclosure, the scenario can include a lexicon. Again, as described herein, the lexicon can represent one or more terms or and regular expressions. A non-limiting example of a term that can be included in the lexicon is a string of one or more text characters, (e.g. a word).

At step 144, the system can determine whether determining whether each of the at least one alert represents an actual violation of the predetermined policy. As a non-limiting example, if the predetermined policy can configured to detect the dissemination of confidential information. This could represent a violation of a law, regulation, or internal policy. But a communication identified by the predetermined policy as a potential violation may not represent an actual violation of the underlying law, regulation or policy (i.e. a false positive).

In some embodiments of the present disclosure, determining whether each alert represents an actual violation of the policy is referred to as “actioning” the alert. This can include determining whether each of the at least one alert represents an actual violation of the policy, law, or ethical standard that the policy/scenario that generated the alert is configured to detect. Actioning the alert can include displaying the alert to a user and receiving a user input from a user interface representing whether the alert represents an actual violation of the policy.

In some embodiments of the present disclosure, the scenario can include a machine learning classifier and determining whether the at least one alert represents an actual violation can include labeling the alert and using the labeled alert to train the machine learning classifier. As another non-limiting example, the present disclosure contemplates that labeling can include labeling alerts as “good” “bad” and “neutral.” Optionally, a “good” alert is an alert that is considered to correctly identify a violation (e.g. a compliance risk), a “bad” alert is an alert that does not correctly identify a violation (i.e. a false positive), and a “neutral” alert is an alert that is

not a true or false positive. This can include alerts where there is ambiguity, or insufficient information to determine whether an alert is correct at the time that it is reviewed.

At step 146, the system calculates a metric based on the actual violations and the potential violations where the metric can include a number of false positives in the at least one alert or the  
5 number of false negatives in the at least one alert. In some embodiments of the present disclosure, the system can display the metric to the user of the system.

At step 148, the system can change the scenario, the target population, and/or the workflow based on the calculated metric. If the scenario used by the system includes one or more lexicons, changing the scenario can include adding or removing one or more terms or regular expressions  
10 from the lexicon(s). In some embodiments of the present disclosure, the target population includes a domain exclusion list and changing the target population includes changing the domain exclusion list.

The present disclosure also contemplates that, in some embodiments, the electronic communication can include metadata, and the scenario can include rules for filtering the  
15 communication based on the metadata. When the scenario includes rules for filtering the communication based on the metadata, changing the scenario can include changing the rules for filtering the communications based on the metadata.

The following describes actioning, and in particular aspects of how a user can “action” an alert, in accordance with some embodiments of the present disclosure. Some examples are  
20 illustrated in FIG. 9. In some embodiments, according to actioning at the hit level, a user can action a single hit under an alert, enabling more accurate reviews and granular feedback loops (see FIG. 13, for example). For actioning at the alert level, all hits under an alert can also be actioned. Also, to action at the communication level, all hits on a communication can be actioned. Thus, actioning, or an action can refer to a label that is applied to a single hit, all hits under an alert,  
25 and/or all its on a message.

In some embodiments, when marking a message as having been reviewed, the reviewed status is part of the default status list. If a new status list is created, then the reviewed status will not be available unless it is manually added. In some embodiments, when escalating a hit, alert, or message, people of interest can be assigned. In some embodiments, multiple communications  
30 can be actioned from the list. Actioning from the list view applies resolved status changes to hits containing an open status.

Regarding assignments, according to some embodiments, a hit, alert, or message can be assigned to another user, which will be displayed and accessible from their dashboard. Group assignments can also be done within escalation workflow. For instance, LDAP (lightweight directory access protocol) groups can be assigned during the escalation workflow for a hit, alert, and/or message. To change alert status, in some embodiments the hit statuses for a particular alert can be overwritten, once all hits are resolved or unresolved. In some embodiments, if any hits for a particular alert remain open, the alert actions may only apply to that open hit.

Regarding actioning status configurations, in some embodiments, functional permissions are available for actioning, thereby controlling a user's ability to action single messages or multiple messages at once. In some embodiments, a case management API includes actioning at the hit level in addition to actions at the alert and message level. Regarding assignment of manual alerts, in some embodiments, a manual alert can be assigned at the message level for an individual message. Manual alerts can be distinguished from system-generated alerts via the person icon in the alert pill. In some embodiments, to support supervision workflow, alerts may be segregated.

Now particularly referring to the diagram of FIG. 13, in some embodiments, according to actioning at the hit level, alerts and corresponding actioning can serve as an input into reporting and improved model training. A system is shown (and its operational flow), according to one embodiment of the present disclosure. The system provides for a feedback loop such that the results of reviewed alerts can be fed back to components of the system that are used for further training of models (as well as creating and evaluating lexicons, creating scenarios, and creating policies).

In some embodiments of the present disclosure, the system shown in FIG. 13 can be configured to implement one or more of the methods described with reference to FIGS. 1A-1C. As shown in FIG. 13, the system can include modules for creating and evaluating models and lexicons 1302, modules for creating scenarios 1304 and modules for creating policies 1306. In some embodiments of the present disclosure, these three modules can be used alone or in combination to perform the methods described with respect to FIG. 1B. These three modules 1302, 1304, and 1306 can be collectively referred to as "cognition studio" or a "scenario builder." Optionally, the repository 1308 can be used to store scenarios and/or information about the alerts, models, or labeled data that are described with reference to FIGS. 1A-1C above.

Similarly, as shown in FIG. 13, the system can include modules for generating alerts 1310, reviewing alerts 1312, and labeling hits 1314. In some embodiments of the present disclosure, these modules can be configured to perform part or all of the methods illustrated and described

with reference to FIGS. 1A and 1C. Additionally, FIG. 13 also illustrates a feedback path 1316 for how labeled hits can be fed back into “cognition studio” to further improve the scenarios created. Optionally, the present disclosure contemplates that the feedback illustrated in FIG. 13 is the feedback described above with reference to FIGS. 1A and 1C.

5           With reference to FIG. 22, a user interface for accessing a repository is shown (e.g. the cognition repository 1308 illustrated in FIG. 13) is shown. The user interface can allow a user to browse, search, import, and export models, lexicons, scenarios, and any of the other data stored in the repository. The exported models, lexicons, scenarios, or other data can be referred to as “artifacts.”

10           With reference to FIGS. 23A-23F, user interfaces for configuring a scenario according to embodiments of the present disclosure are shown. FIG. 23A illustrates a user interface for viewing one or more datasets. FIG. 23B illustrates a user interface for labeling a dataset. FIG. 23C illustrates an annotation applied to a dataset and an interface for applying labels to a dataset. FIG. 23D illustrates a user interface for configuring a lexicon to be applied to the dataset. FIG. 23E  
15 illustrates a user interface for evaluating a lexicon. FIG. 23F illustrates a scenario created using the lexicon that was configured in the interface shown in FIG. 23E.

          Through the use of a series of functional tools for creating and evaluating lexicons, creating scenarios, and creating policies (labelled collectively in the diagram as “Cognition Studio”), a user (such as a data scientist) user can create a model (e.g., perform training of a model) in cognition  
20 studio for evaluation against established datasets. The user can then create scenarios based on the model(s), lexicons, and non-language features (NLF). Next, the user can create polic(ies) which map to the scenario(s) and population.

          Following the steps collectively labeled under “Cognition Studio”, a user such as a business analyst publishes the scenario(s) to a data repository labeled in the diagram of FIG. 13 as  
25 “Cognition Repository”. The repository can be a data storage device that provides for version-controlled storage of all models, lexicons, scenarios, and policies, and which can allow for labeling of active or draft versions. A user such as a system administrator can select relevant scenario(s) and can select a target population. The user can also select target communication types (e.g., chat, email, etc.) and channels (e.g., chat applications, email servers, etc.), and mark the policy as active.

30           The system according to some embodiments can then use a new active policy or policy version against all newly ingested electronic communications to generate alerts as appropriate (see label in FIG. 13 of “Cognition Logic”) and as described in further detail above. Next, in operations collectively labeled as “Alert” in the diagram, a user such as an analyst (e.g., compliance

representative, etc.) can review the generated alerts and label each hit according to, for instance, escalation workflow in which a true positive is identified. The labeled hits can then be used as feedback to the “Cognition Studio” for supervised improvement of the aspects discussed above with respect to these components and respective functions.

5

### **Example Computing System Architecture**

FIG. 12 is a computer architecture diagram showing a general computing system capable of implementing one or more embodiments of the present disclosure described herein. A computer may be configured to perform one or more functions associated with embodiments illustrated in, and described with respect to, one or more of FIGS. 1-11 and 13-26. It should be appreciated that the computer may be implemented within a single computing device or a computing system formed with multiple connected computing devices. For example, the computer may be configured for a server computer, desktop computer, laptop computer, or mobile computing device such as a smartphone or tablet computer, or the computer may be configured to perform various distributed computing tasks, which may distribute processing and/or storage resources among the multiple devices.

As shown, the computer includes a processing unit, a system memory, and a system bus that couples the memory to the processing unit. The computer further includes a mass storage device for storing program modules. The program modules may include modules executable to perform one or more functions associated with embodiments illustrated in, and described with respect to, one or more of FIGS. 1-11 and 13-26. The mass storage device further includes a data store.

The mass storage device is connected to the processing unit through a mass storage controller (not shown) connected to the bus. The mass storage device and its associated computer storage media provide non-volatile storage for the computer. By way of example, and not limitation, computer-readable storage media (also referred to herein as “computer-readable storage medium” or “computer-storage media” or “computer-storage medium”) may include volatile and non-volatile, removable and non-removable media implemented in any method or technology for storage of information such as computer-storage instructions, data structures, program modules, or other data. For example, computer-readable storage media includes, but is not limited to, RAM, ROM, EPROM, EEPROM, flash memory or other solid state memory technology, CD-ROM, digital versatile disks (“DVD”), HD-DVD, BLU-RAY, or other optical storage, magnetic cassettes, magnetic tape, magnetic disk storage or other magnetic storage devices, or any other

medium which can be used to store the desired information and which can be accessed by the computer. Computer-readable storage media as described herein does not include transitory signals.

According to various embodiments, the computer may operate in a networked environment  
5 using connections to other local or remote computers through a network via a network interface unit connected to the bus. The network interface unit may facilitate connection of the computing device inputs and outputs to one or more suitable networks and/or connections such as a local area network (LAN), a wide area network (WAN), the Internet, a cellular network, a radio frequency network, a Bluetooth-enabled network, a Wi-Fi enabled network, a satellite-based network, or  
10 other wired and/or wireless networks for communication with external devices and/or systems.

The computer may also include an input/output controller for receiving and processing input from a number of input devices. Input devices may include, but are not limited to, keyboards, mice, stylus, touchscreens, microphones, audio capturing devices, or image/video capturing devices. An end user may utilize such input devices to interact with a user interface, for example  
15 a graphical user interface on one or more display devices (e.g., computer screens), for managing various functions performed by the computer, and the input/output controller may be configured to manage output to one or more display devices for visually representing data.

The bus may enable the processing unit to read code and/or data to/from the mass storage device or other computer-storage media. The computer-storage media may represent apparatus in  
20 the form of storage elements that are implemented using any suitable technology, including but not limited to semiconductors, magnetic materials, optics, or the like. The program modules may include software instructions that, when loaded into the processing unit and executed, cause the computer to provide functions associated with embodiments illustrated in, and described with respect to, one or more of FIGS. 1-11 and 13-26. The program modules may also provide various  
25 tools or techniques by which the computer may participate within the overall systems or operating environments using the components, flows, and data structures discussed throughout this description. In general, the program module may, when loaded into the processing unit and executed, transform the processing unit and the overall computer from a general-purpose computing system into a special-purpose computing system.

30

### CONCLUSION

The various example embodiments described above are provided by way of illustration only and should not be construed to limit the scope of the present disclosure. Those skilled in the art will readily recognize various modifications and changes that may be made to the present disclosure without following the example embodiments and applications illustrated and described  
5 herein, and without departing from the true spirit and scope of the present disclosure.

## CLAIMS

What is claimed is:

1. A computer-implemented method, comprising:
  - 5 receiving at least one alert from a conduct surveillance system, wherein the at least one alert represents a potential violation of a predetermined standard and wherein the conduct surveillance system generates the alerts in response to an electronic communication between persons matching a violation of a predetermined policy, wherein the predetermined policy comprises a scenario, a target population, and a workflow;
  - 10 determining whether each of the at least one alert represents an actual violation of the predetermined policy;
  - calculating a metric based on the actual violations and the potential violations wherein the metric comprises a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and
  - 15 changing at least one of the scenario, the target population, or the workflow based on the calculated metric.
2. The computer implemented method of claim 1, wherein the scenario comprises a machine learning classifier, and wherein determining whether the at least one alert represents an actual violation comprises labeling the at least one alert and using the labeled at least one  
20 alert to train the machine learning classifier.
3. The computer implemented method of claim 1 or 2, wherein the metric is displayed to a user.
- 25 4. The computer implemented method of any one of claims 1-3, wherein the scenario comprises a lexicon, and wherein the lexicon represents one or more terms or regular expressions.
- 30 5. The computer implemented method of any one of claims 1-4, wherein changing the scenario comprises changing the lexicon by adding or removing terms or regular expressions from the lexicon.

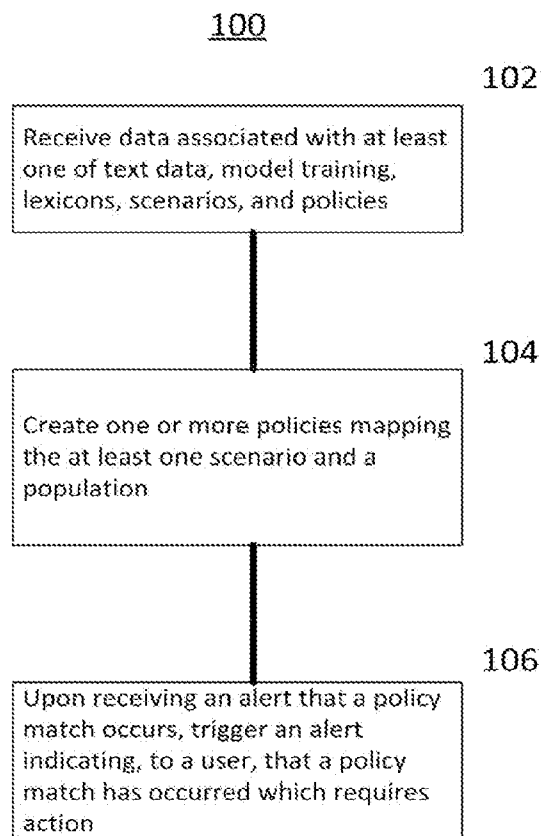
6. The computer implemented method of any one of claims 1-5, wherein, in response to determining that the at least one alert represents an actual violation, actioning the alert.
7. The computer implemented method of claim 6, wherein actioning the alert comprises receiving a user input from the user interface representing whether the at least one alert represents an actual violation.
8. The computer implemented method of any one of claims 1-7, wherein the target population comprises a domain exclusion list and wherein changing the target population comprises changing the domain exclusion list.
9. The computer implemented method of any one of claims 1-8, wherein the electronic communication comprises metadata, the scenario comprises rules for filtering the electronic communication based on the metadata, and wherein changing the scenario comprises changing the rules for filtering the electronic communications based on the metadata.
10. A system, comprising:
- at least one processor;
  - at least one memory storing computer readable instructions configured to cause the at least one processor to perform functions for creating and/or evaluating models, scenarios, lexicons, and/or policies, wherein the functions include:
    - receiving data associated with at least one of text data, model training, lexicons, scenarios, and policies, wherein the functions for creating and/or evaluating models comprise creating at least one scenario based on at least one of the models, lexicons, and non-language features;
    - creating one or more policies mapping to the at least one scenario and a population;
    - upon receiving an alert that a policy match occurs, triggering an alert indicating, to a user, that a policy match has occurred which requires a user action, wherein a policy corresponds to actions that violate at least one of a combination of signals and metrics, a population, and workflow.

11. The system of claim 10, wherein the model training comprises training at least one model configured to analyze the text data from one or more electronic communications between at least two persons.
- 5 12. The system of claim 10 or 11, wherein the user action comprises review and interaction by a user via a user interface.
13. The system of any one of claims 1-12, wherein the model training comprises evaluating the model against established datasets.
- 10 14. The system of any one of claims 1-13, wherein the alert to the user is evaluated by the user and a corresponding user decision is made to confirm or deny accuracy of the alert.
- 15 15. The system of claim 14, wherein the user decision is provided into a feedback loop, and wherein the feedback loop is configured to improve the model training.
16. The system of claim 15, wherein the user decision is provided into the feedback loop and wherein the feedback loop is configured to improve the lexicons, scenarios, or policies.
- 20 17. The system of claim 16, wherein the feedback loop is configured to change a lexicon.
18. The system of claim 17, wherein changing the lexicon comprises configuring the lexicon so that it includes or excludes terms or regular expressions.
- 25 19. The system of claim 15, wherein the feedback loop is configured to measure the rate of false positives and to change one or more of the lexicons, scenarios, and policies based on the rate of false positives.
- 30 20. The system of claim 15, wherein the scenario includes Boolean operators, and wherein the feedback loop is configured to change one or more of the Boolean operators.

21. The system of claim 16, wherein the feedback loop is configured to monitor the rate of false positives over a period of time, and change one or more of the lexicons, scenarios, and policies based on the rate of false positives over the period of time.
- 5 22. A non-transitory computer-readable medium storing instructions which, when executed by at least one processor of a computer, perform functions that include:
- receiving at least one alert from a conduct surveillance system, wherein the at least one alert represents a potential violation of a predetermined standard and wherein the conduct surveillance system generates the alerts in response to an electronic communication between
  - 10 persons matching a violation of a predetermined policy, wherein the predetermined policy comprises a scenario, a target population, and a workflow;
  - determining whether each of the at least one alert represents an actual violation of the predetermined policy;
  - calculating a metric based on the actual violations and the potential violations wherein
  - 15 the metric comprises a number of false positives associated with the at least one alert or the number of false negatives associated with the at least one alert; and
  - changing at least one of the scenario, the target population, or the workflow based on the calculated metric.

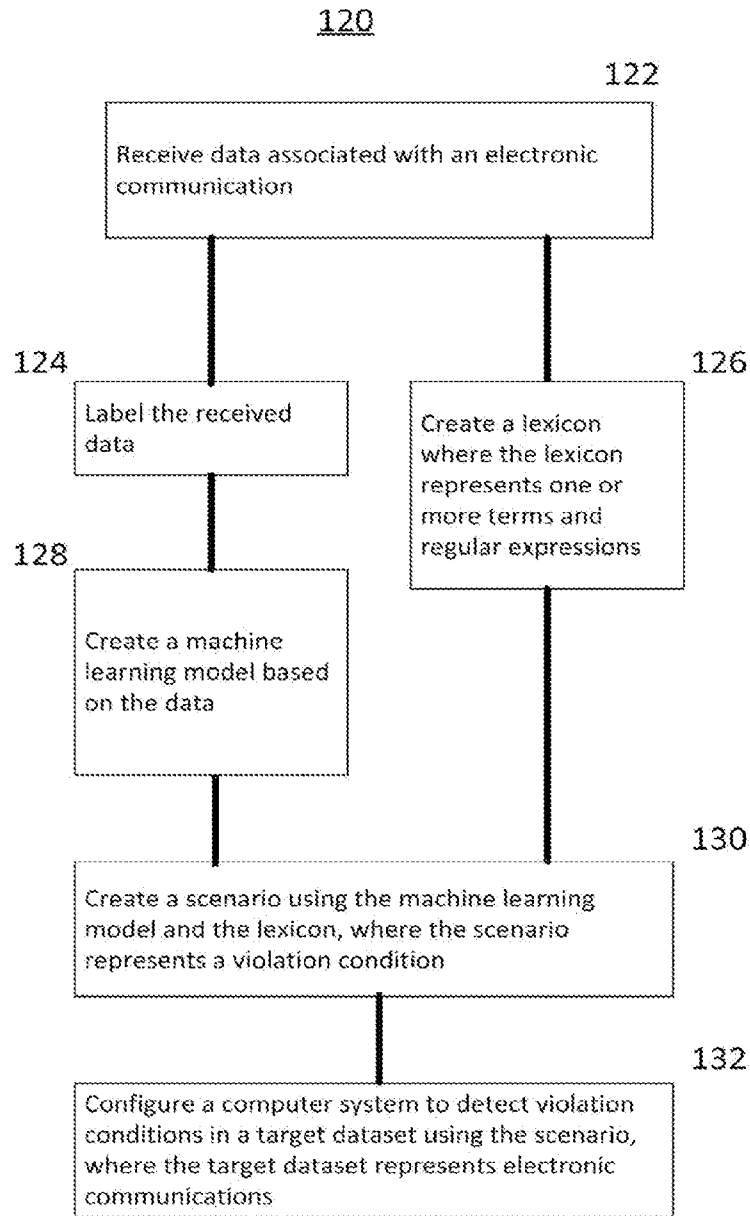
20

1/34

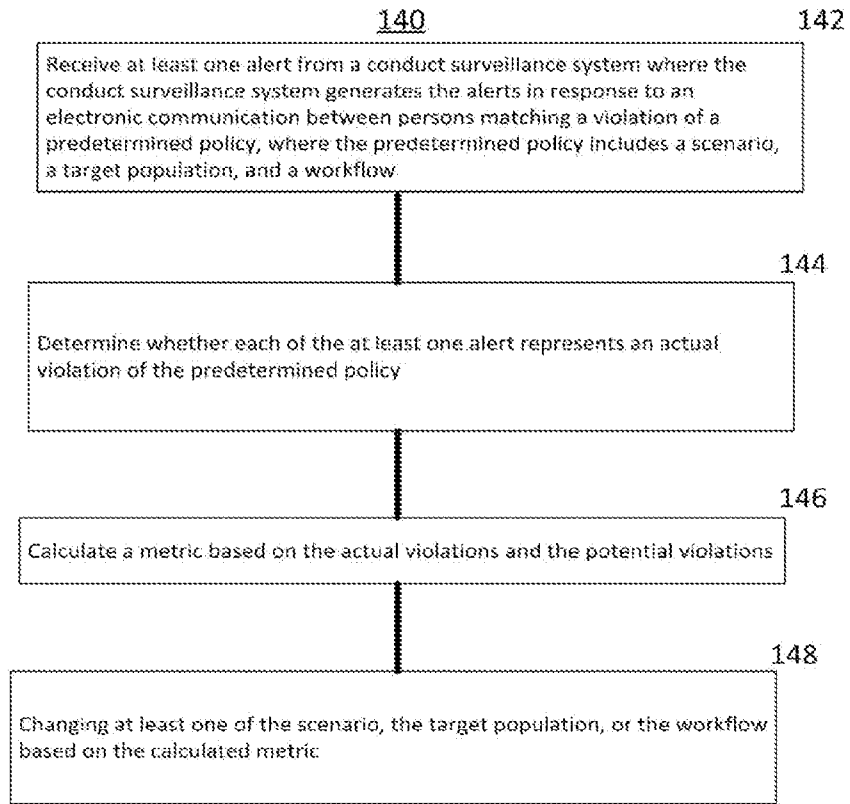


**FIG. 1A**

2/34



**FIG. 1B**



**FIG. 1C**

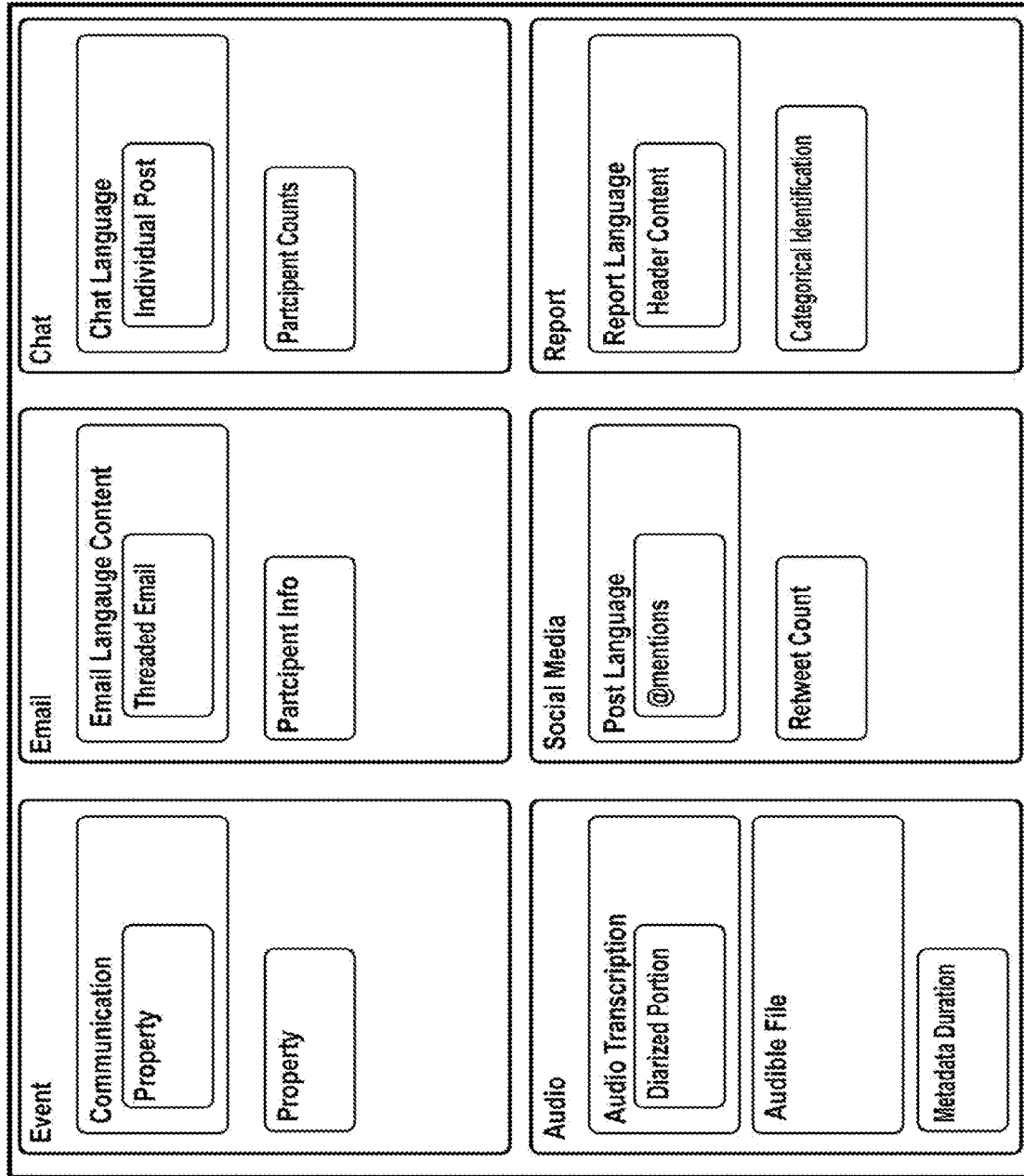
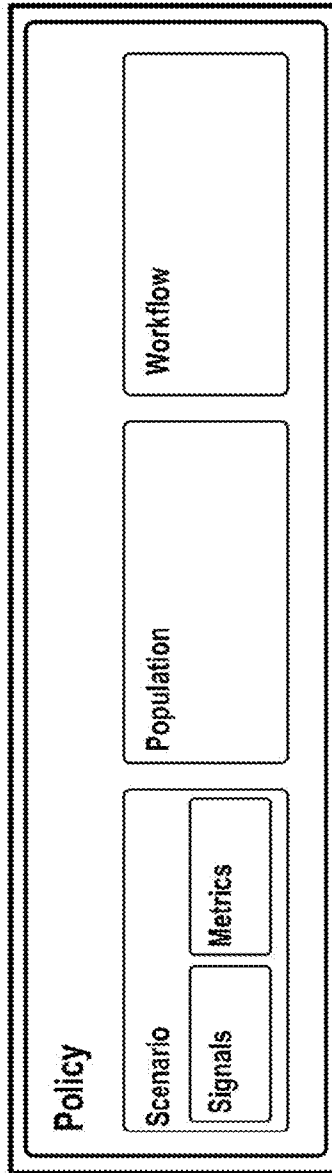
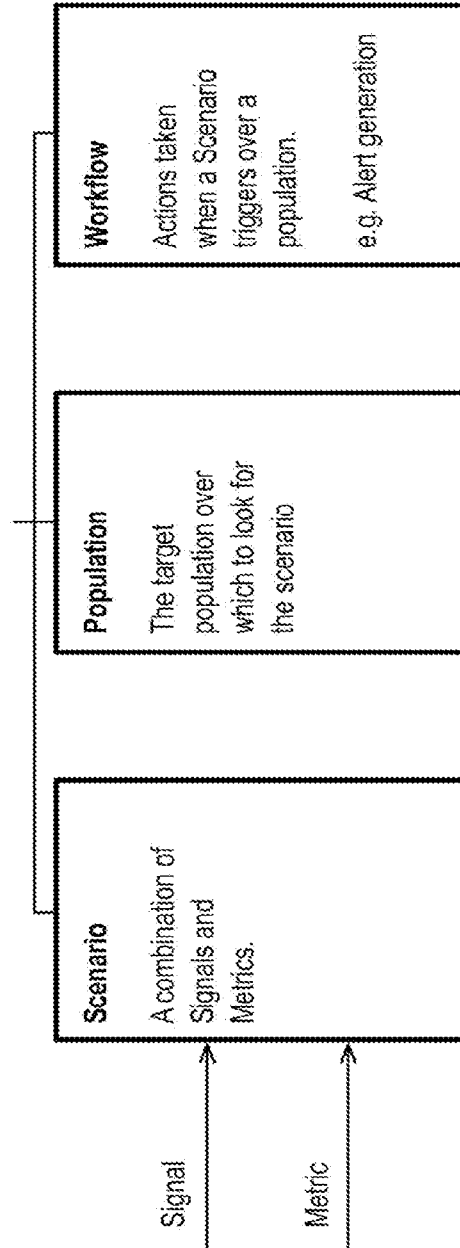


FIG. 2A



**FIG. 2B**



**FIG. 2C**

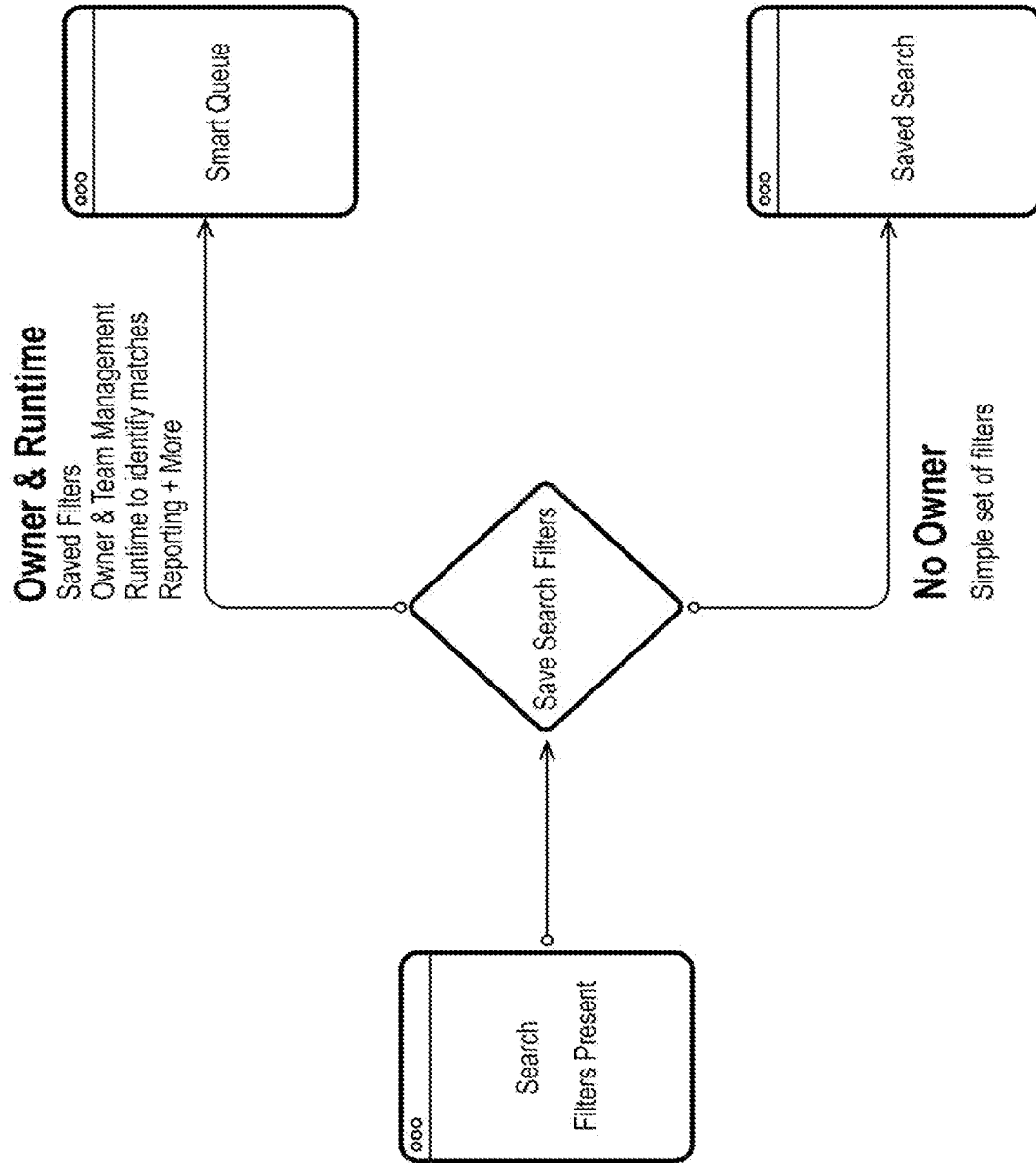


FIG. 3

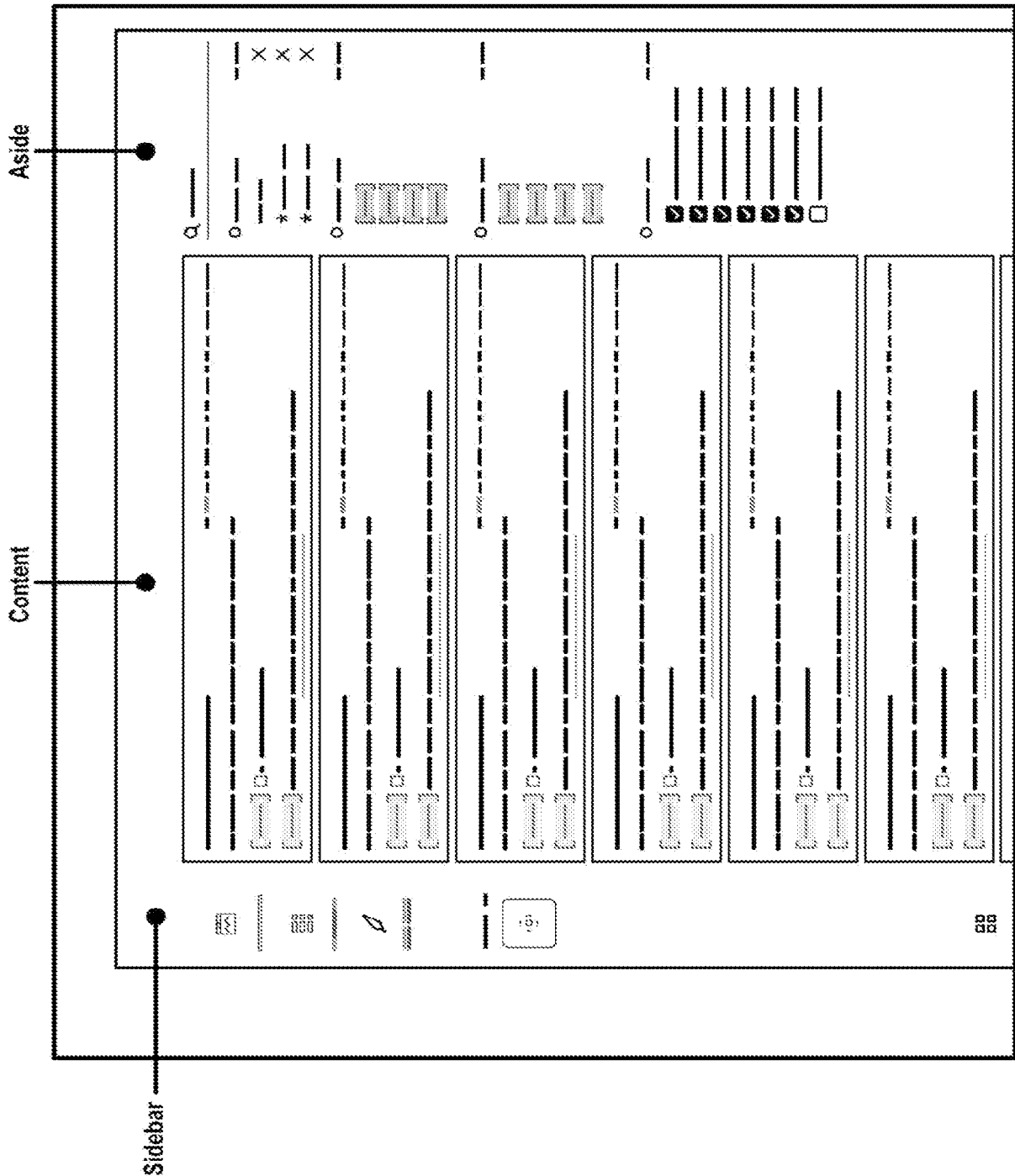


FIG. 4

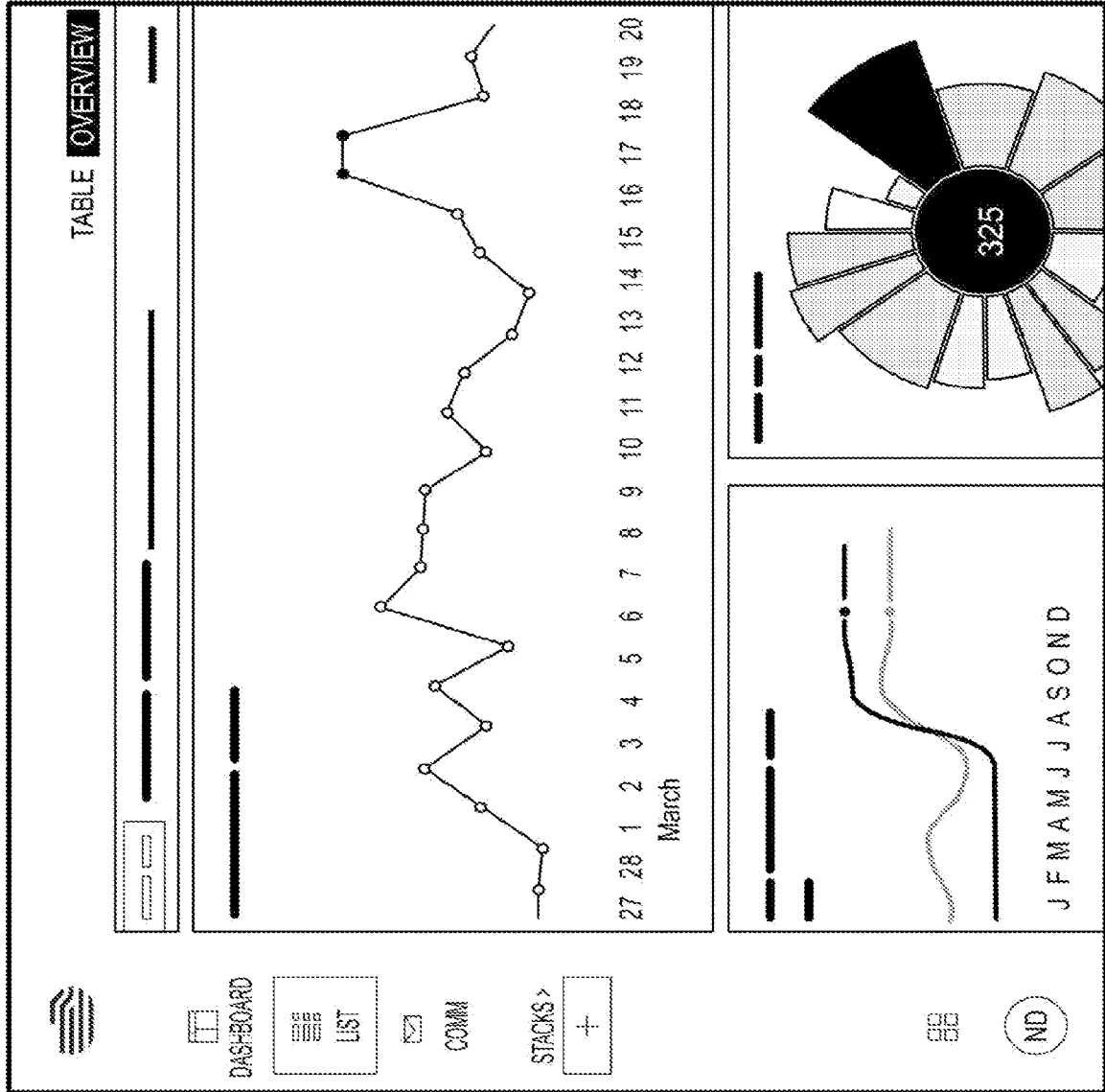


FIG. 5

Communication for Knowledge Task

Activity Metadata Message

# Conduct

To: richard.sanders@enron.com, wendy.davis@internalemail.com, salivar.dias@anotheremail.com, steve.rogers@marvel.com

From: frank.douglas@digitalreasoning.com

Sent: Tuesday, november 15 2018 - 12:24:34 UTC

0 : 4 attachments (total 12.6 mb)

**SUBJECT: RE: Important Wells Fargo Contact Information**

PLEASE. We have repairs that must be done ASAP, that's why we started the process in the beginning of October!!

Sent: Tuesday, November  
 TO: Herrera, Benjamin Subject: re: Important Wells Fargo Contact Information

I tried to call you today at the number in your email to me; however, your mailbox is full. You indicated in your email to me that if I had any questions, I was to contact you. I DO have a question. It has now been almost 3 months since we started the process.

Prev Next

Escalate 1

False Positive 2

News 3

Spam 4

Search

Focus View  On

CONCERNS -1 less

PEOPLE +23 more

Benjamin Herrera

Longtraine Jason Jones III

Time: 23:34

Large attachment size

Mispelling

High # of Attachments

Actions Required by a User

Communication for Knowledge Task

Conduct

Message

Activity Metadata

Search Print

To: richard.sanders@enron.com, wendy.davis@internalemail.com, salivar.dias@anotheremail.com, steve.rogers@marvel.com

From: frank.douglas@digitalreasoning.com

Sent: Tuesday, november 15 2018 - 12:24:34 UTC

0 : 4 attachments (total 12.6 mb)

SUBJECT: RE: Important Wells Fargo Contact Information

PLEASE. We have repairs that must be done ASAP, that's why we started the process in the beginning of October!!

Sent: Tuesday, November  
TO: Herrera, Benjamin Subject: re: Important Wells Fargo Contact Information

I tried to call you today at the number in your email to me; however, your mailbox is full. You indicated in your email to me that if I had any questions, I was to contact you. I DO have a question. It has now been almost 3 months since we started the process.



DASHBOARD

QUEUE

MSG

STACKS >



Content Viewer For a Comm

FIG. 6

10/34

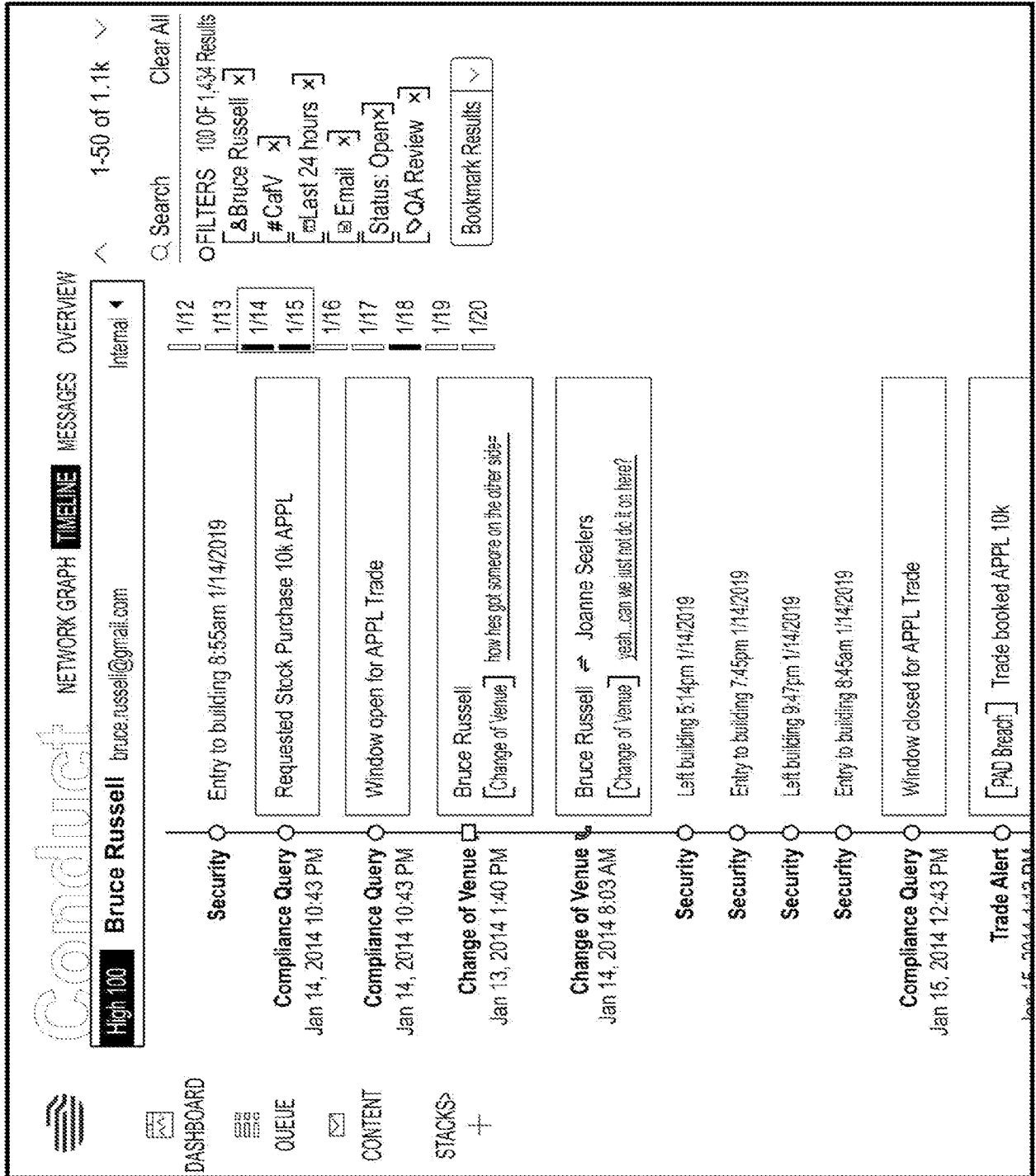


FIG. 7

11/34

Aside

Conduct

High 100 **Bruce Russell** bruce.russell@gmail.com Internal

MESSAGES OVERVIEW

1-50 of 1.1k Clear All

SEARCH 100 OF 1,434 RESULTS

Bruce Russell 
 #CaV 
 mLast 24 hours 
 Email 
 Status: Open 
 QA Review

Bookmark Results

---

DASHBOARD

QUEUE

CONTENT

STACKS

+

Thursday, January 01, 2014 4:17:20 AM	88 in   67 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] ...Duis mollis, est non commodo luctus, nisi erat porttitor ligula, eget lacinia odio... <input type="checkbox"/> Josh Goodwin <input type="checkbox"/> [Change of Venue] ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	
Wednesday, November 11, 2013 11:10:08 PM	24 in   78 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	
Tuesday, December 12, 2014 12:17:02 PM	37 in   52 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] & ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	
Saturday, January 01, 2014 6:52:17 PM	15 in   26 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] 1 ...how has got someone on the other side...	
Wednesday, December 12, 2016 6:10:45 PM	73 in   38 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] & ...Duis mollis, est non commodo luctus, nisi erat porttitor ligula, eget lacinia odio... <input type="checkbox"/> Josh Goodwin <input type="checkbox"/> [Change of Venue] ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	
Wednesday, February 02, 2016 11:21:10 AM	51 in   32 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] & ...Duis mollis, est non commodo luctus, nisi erat porttitor ligula, eget lacinia odio... <input type="checkbox"/> Josh Goodwin <input type="checkbox"/> [Change of Venue] ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	
Wednesday, July 07, 2017 4:28:49 AM	46 in   37 in
<input type="checkbox"/> Bruce Russell <input type="checkbox"/> [Change of Venue] & ...Duis mollis, est non commodo luctus, nisi erat porttitor ligula, eget lacinia odio... <input type="checkbox"/> Josh Goodwin <input type="checkbox"/> [Change of Venue] ...Donec sed odio dui. Praesent, commodo cursus magna, vel scelerisque nisl...	

FIG. 8

Alert: Hit Count

Back Next

ALERTS Add Manual Alerts

Subject 1

Corruption Class

Header 1

Rationalization 1

Quote Header 1

Corruption Regexp 1

Financial 1

MARK REMAINING AS

False Positive 1

Newsletter 2

Spam 3

Reviewed 4

Closed-Action Taken 5

Good hit-no issue 6

Escalate 7

Follow-Up 8

Other 9

FOCUS VIEW

From: jeb@jeb.com

Sent: Tuesday, July 25 2006 - 16:34:47 UTC

1 attachment (total 5.1 KB)

NCLB v2.0 oped.doc

Subject: FW: Proposed Op-ed

Alert

Communication METADATA HTML CONTENT

Message

Original Message----

From: Klein, Joel [mailto:JKlein@schools.nyc.gov]

Sent: Monday, July 24, 2006 3:01 PM

To: Jeb Bush

Subject: Proposed Op-Ed

Jeb, great taking today. I appreciate your understanding on all the sensitivities. You're a pro. The attached embodies the points that we discussed and your folks sent in draft taking points. If you approve, we should then have your communications people and the mayor's figure where, when etc. In the meantime, we should find a time for a joint appearance for you and me before a think tank in nyc (or anywhere else you'd like). If you have some dates in september/october, let me know and we'll get started. Again, thanks for your leadership and patience. Joel

Hit (Offset)

FIG. 9

List Card

Search Clear All

FILTERS 100 Of 1,434 Results

- Richard Sanders x
- #Bribery x
- Jul, 2019 x
- Email x
- Status: Open x
- QA Review x

LIST OPTIONS

- Common
- Alerts
- BlagielgrhdTime
- Participant Count
- Email
- To
- Subject

AUDIO

- Participants
- Run Time

Chat

- Active Contributors
- Inactive Attendees

Sort By

---

<p>Aubrey Henry</p> <p>RE: Subject line subject line subject line subject line</p> <p>[Bribery] @ : Camping2018.docx</p> <p>[Bribery] ...cras justo odio, dapibus ac facilisis in, egestas eget Nulla vitae eli aquam...</p>	<p>Bessis Nguyen</p> <p>RE: Subject line subject line subject line subject line</p> <p>[Bribery] @ : Camping2018.docx</p> <p>[Bribery] ...cras justo odio, dapibus ac facilisis in, egestas eget Nulla vitae eli aquam...</p>	<p>Regina Bell</p> <p>RE: Subject line subject line subject line subject line</p> <p>[Bribery] @ : Camping2018.docx</p> <p>[Bribery] ...cras justo odio, dapibus ac facilisis in, egestas eget Nulla vitae eli aquam...</p>	<p>Kathryn Warren</p> <p>RE: Subject line subject line subject line subject line</p> <p>[Bribery] @ : Camping2018.docx</p> <p>[Bribery] ...cras justo odio, dapibus ac facilisis in, egestas eget Nulla vitae eli aquam...</p>	<p>Theresa Lane</p> <p>RE: Subject line subject line subject line subject line</p>
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------

Alert Hit Preview

FIG. 10

**Conduct**

Sort By: Internal/External

**Activity** letziata **Message**

Search Download Print

**Metric** [100] vandy.davis@gmail.com & [89] jason.jones@intermail.com

ext [34] salivar.dias@anothermail.com, [31] Steve Rogers &

**From:** sam.davis@intermail.com

**Sent:** Tuesday, November 15 2018 - 12:23:34 UTC

**0** : 4 attachments (total 12.6 mb)

**Subject:** RE: Important Wells Fargo Contact Information

**Communication** METADATA HTML **CONTENT**

PLEASE, We have repairs that must be done ASAP, that's why we started the process in the beginning of October!

Sent: Tuesday, November

To: Herrera, Benjamin Subject: re: Important Wells Fargo Contact Information

Bullet List

Bullet 1

Back Next

ALERTS [Corruption] 15 ^ ^ [Inappropriate L.] 99 ^ ^

News 1

Spam 2

Escalate 3

False Positive 4

FOCUS VIEW  ON

CONCERNS -1 less [Time: 23:34] [Large attachment size] [Misspelling] [High # of attachments]

PEOPLE +23 more [Benjamin Herrera] [Longname Jason Jones III] [Benjamin Herrera] [Jason]

Tab

FIG. 11

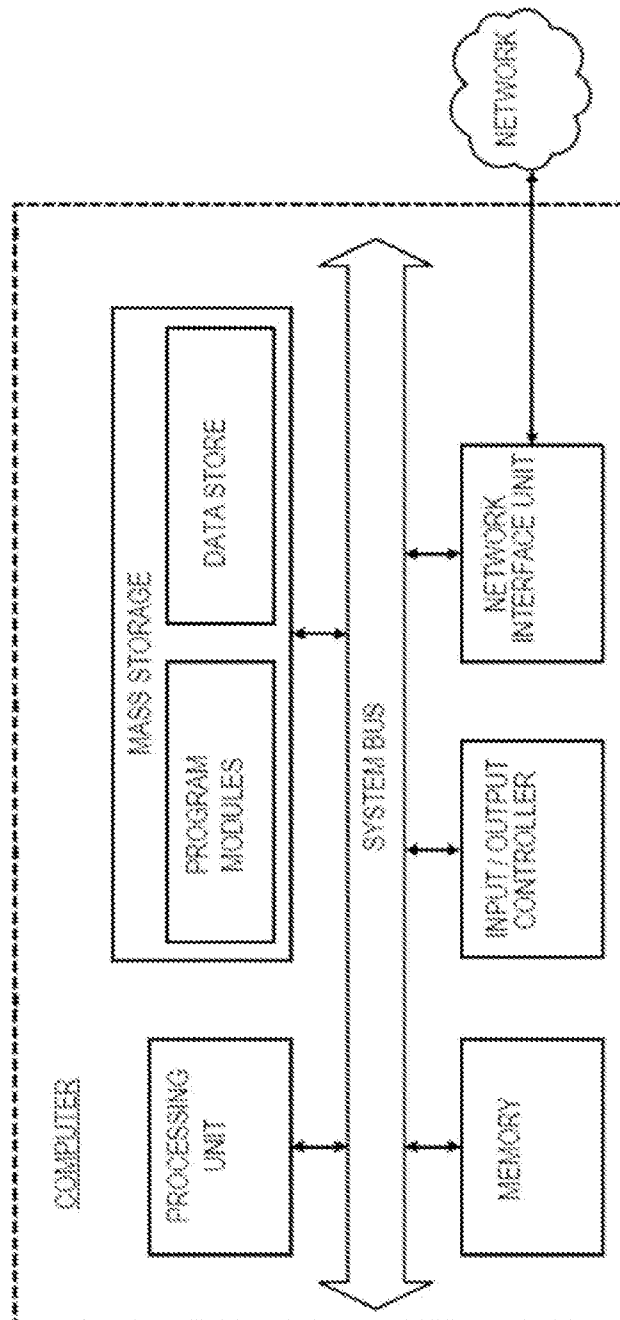


FIG. 12

16/34

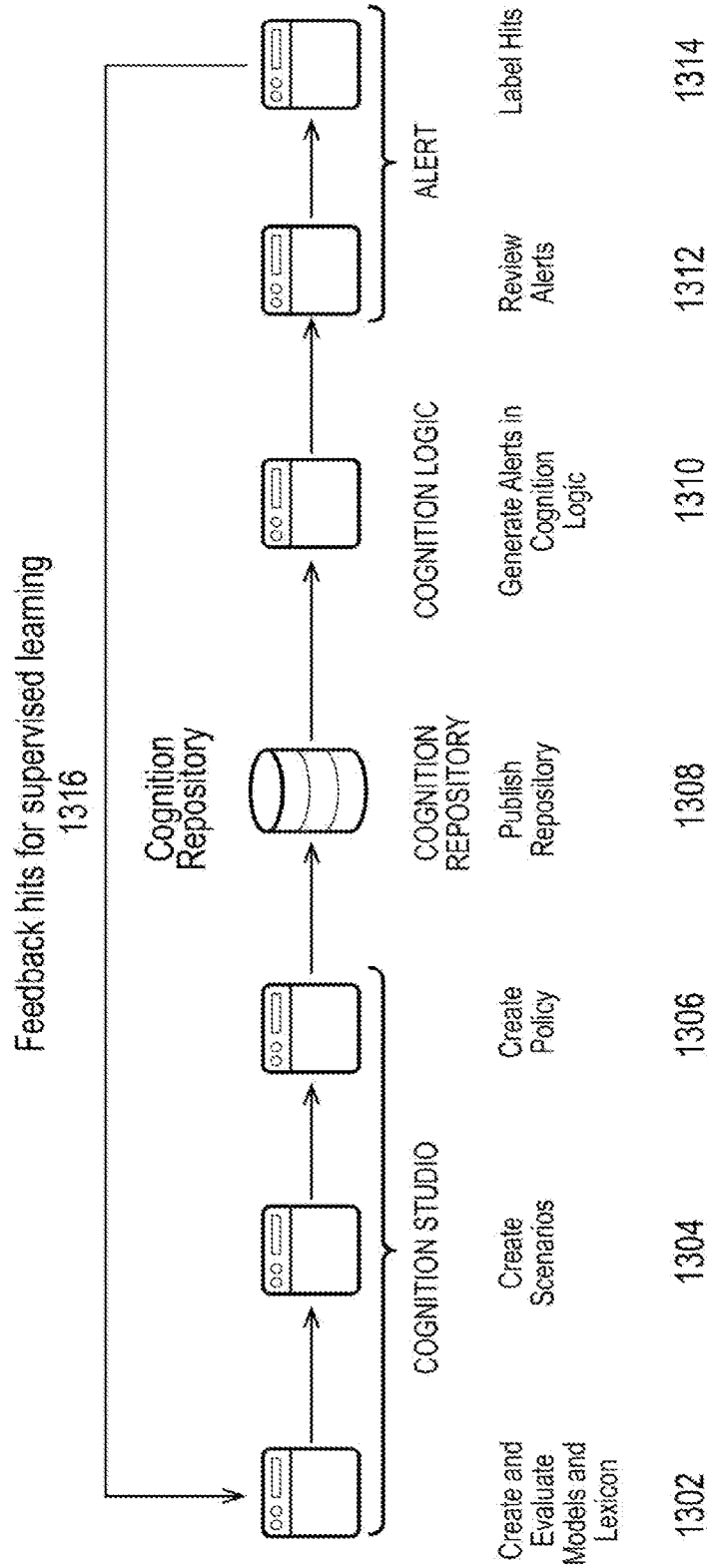


FIG. 13

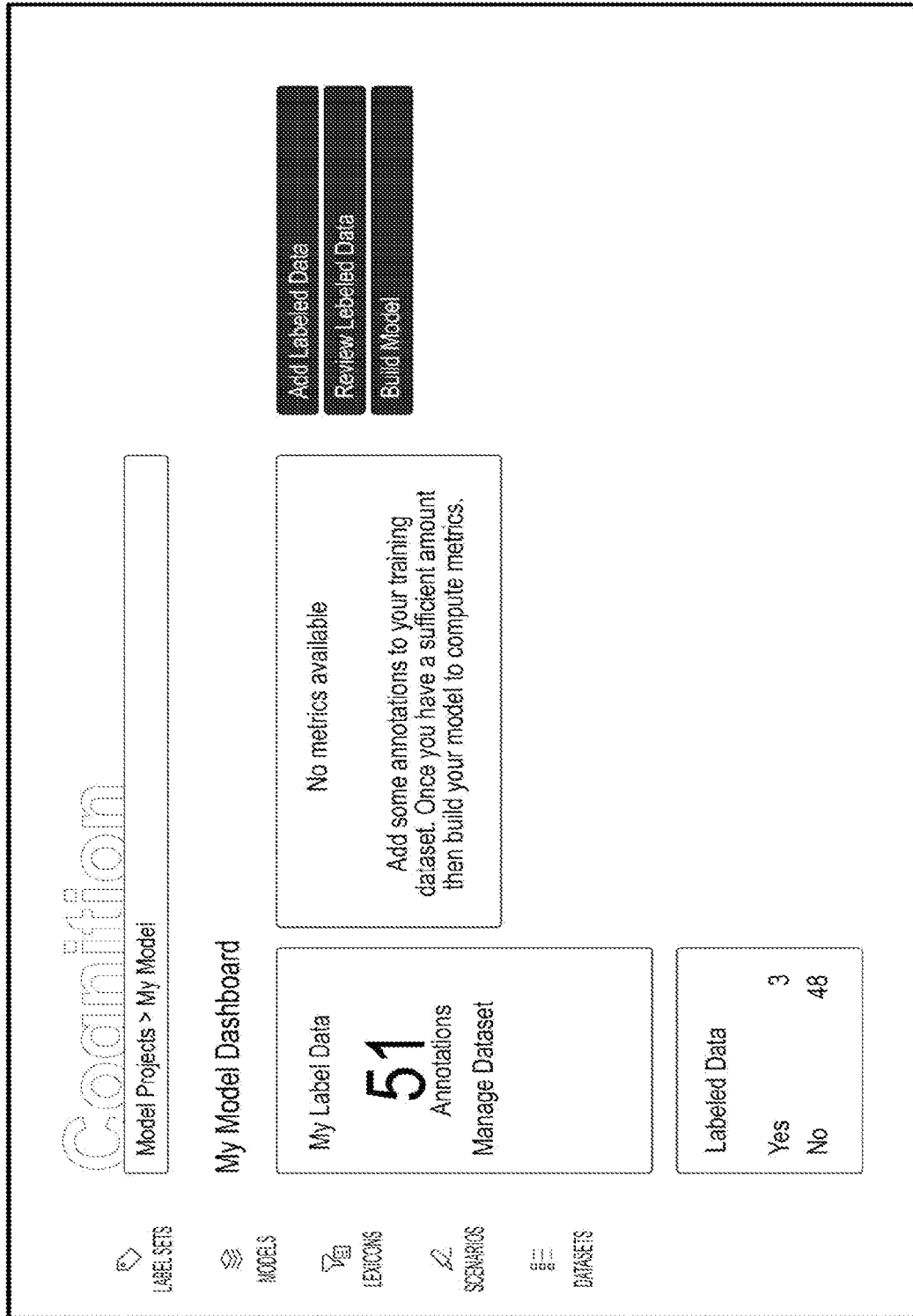


FIG. 14

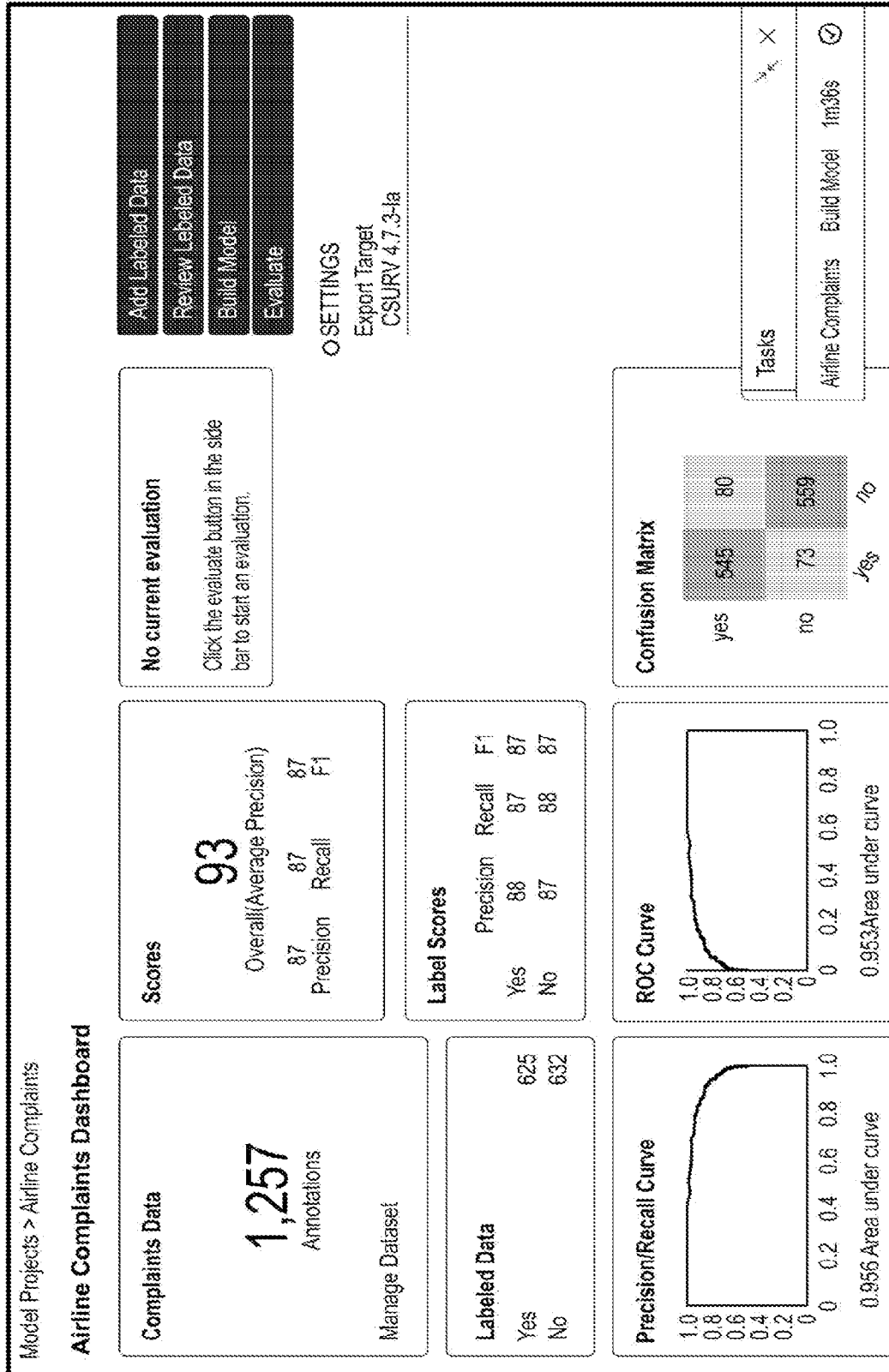


FIG. 15

Lexicons > Secrecy > Evaluate

**Evaluate Lexicon**

Sample 1-50 of 1,627

Sample	Matches	Hit
Fwd: A secret to tell	10	Yes
----- Inline attachment follows ----- +Read more		
Fwd: A secret to tell	10	Yes
----- Inline attachment follows ----- +Read more		
CMTA Legislative Weekly - 10/26/01	8	Yes
The CMTA Fall Conference schedule for October 31 - November 2 has been cancelled. ----- +Read more		
Best Kept Secret in the Financial Market	7	Yes

oMETRICS Hits Dataset 1627 (100%) enron\_services\_el

^ Search

Labels  
Models  
Lexicons  
Datasets

FIG. 16

Lexicons > Secrecy > Evaluate

**Evaluate Lexicon**

1-50 of 3,437

Search

o METRICS  
Hits 117 (3%)  
Dataset secrecy\_train

Sample	Matches	Hit	Status
this is not to be shared with the others.	1	Yes	True Positive
ask your employees to keep silent.	1	Yes	True Positive
ask your boy to keep quiet.	1	Yes	True Positive
ask your friends to keep hush-hush.	1	Yes	True Positive
can we plan to keep this hush-hush.	1	Yes	True Positive
ask your employees to keep hush-hush.	1	Yes	True Positive
bury this audit.	1	Yes	True Positive
please don't tell shit.	1	Yes	True Positive
tell your friends to keep quiet.	1	Yes	True Positive
bury this paper trail.	1	Yes	True Positive

Confusion Matrix

Hit	117	1342
Miss	0	1979
	Hi	Miss
	Predicted	

FIG. 17

Lexicons > Secrecy > Evaluate > Sample

**Evaluate Lexicon Sample**

RE: Welch hunt

8 of 1,627

---

There is a direct analogy to the Civil War here. Soon after the Union began to lose with much embarrassment in the field, a Committee For Overseeing the Conduct of the War comprised of big angry guys from Congress was formed. I'm not sure how. I think they just kind of announced that they were formed. They began interviewing people and were quite vicious about their confessions. I recall correctly, they gave themselves the power to jail people. Lincoln endured them and eventually they atrophied, but this kind of savage second guessing and vigilante vengeance by elected legislators is definitely a theme of politics. The notion that Dunn is seeking to keep things quiet to prevent rumors is total bullshit. If I were one of those poor saps being interviewed, I would get a lawyer and scream bloody murder about due process. Dunn can only do harm with this tactic, and again, it shows that the pathology of the body politic in California has advanced to its most fatal degree. These people are sick. ...cg

California state senator Joe Dunn is conducting a secret McCarthyesque inquisition of ISO Board members in an apparent effort to find a scapegoat to soak up the political blame for the energy crisis and divert negative criticism from Dunn. Dunn better pick up the pace, the March primary election is only a few months away.

Energy prices kept under wraps. In secret, a Capitol panel is deposing key players in the lifting of a power price cap.

By John Hill  
See Capitol Bureau

SEARCH

o METRICS  
Matches 6  
Hit Yes

o FOCUS VIEW

ON

o TERMS  
keep FOLLOWEDBY() (quote OR silent)

secret?

---

LABEL SETS

MODELS

LEXICONS

DATA SETS

FIG. 18

<p>Dunn is seeking to <u>keep things quiet</u> to prevent rumors is total bullshit. If I were one of those poop saps being interviewed, I would get a lawyer and scream bloody murder about due process. Dunn can only do harm with this tactic, and again, it shows that the pathology of the body politic in California has deviated to its most florid degree. There people are sick. ---CGY</p> <p>California state senator Joe Dunn is conducting a <u>secret</u> McCarthyesque inquisition of ISO Board members in an apparent effort to find a scapegoat to soak up the political blame for the energy crisis and divert negative criticism from Davis. Dunn better pick up the pace, the March primary election is only a few months away.</p> <p>Energy probe kept under wraps: In <u>secret</u>, a Capitol panel is deposing key players in the filing of a power price cap.</p> <p>By John Hill Bee Capitol Bureau (Published Oct. 30, 2011)</p>	<p>Hit</p> <p><input type="radio"/> FOCUS VIEW</p> <p><input checked="" type="checkbox"/> On</p> <p><input type="radio"/> TERMS</p> <p>Keep FOLLOWED BY (quiet OR silent) <input checked="" type="radio"/></p> <hr/> <p>secret? <input checked="" type="radio"/></p> <p><input type="radio"/></p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

FIG. 19

**Coanition**

Key Indicators > Secret Spam > Definition

**Secret Spam**    HISTORY    VERSIONS    DASH    **DEFINITION**

Model: Secrety v2.1.1

- ANY message
- EVERY section
- ONEOF
- Lexicon: Spam v4.7
- Model: Spam v4.7
- Model: Secrety v2.1.1
- ANY message
- Model: Newsletter v2.1.1

Search: name: |

- Standard Features
- Participant Count v2.9
- Models
- Secrecy Model v2.1.1
- Newsletter Model v0.1
- Spam Model v4.7
- Lexicons
- Operators

Navigation: LABEL SETS, MODELS, LEXICONS, KIs, DATASETS

FIG. 20

**POLICY ADMIN**

**Create/Edit New Policy**

**Save**

Cancel

Status

Active

Inactive

**Policy Name**

Name

Gifts and Entertainment (UK)

**Policy Description**

Label

G&E English and Spanish UK comp only

**Active Key Indicator**

Search Key Indicators

Bribery	V2 dd/mm/yyyy
Customer Complaints	V3 dd/mm/yyyy
Gifts and Entertainment	V2 dd/mm/yyyy
Harassment	V2 dd/mm/yyyy
Secrecy	

**Policy Group**

Search Policy Group

UK\_Compliance

**Knowledge Graph**

Search Knowledge Graph

- english\_20200426
- english\_1283726
- spanish\_20397191
- spanish\_4902836
- cantonese\_1039826
- french\_09372909
- japanese\_9382792
- english\_20200426

FIG. 21

Repository

Repository

All  
 Labeled Dataset  
 Lexicon  
 Model  
 Scenario

Name	Path	Artifact Type
quidproquo	test.notebook.17f060004308a0cc3ad1244e1234foac	Model
guaranteessurances-test	smash-guaranteessurances	Labeled Dataset
Complaints	dir.notebook.27f06252e308a7ab3ac3688e5362bnda	Model
guaranteessurances-1r	smash-guaranteessurances	Model
tes2	sm.11111a000ccc1a0ccc11ababab000.2	Lexicon
nrm-1234-testing	sm.11111a000ccc1a0ccc11ababab000.2	Scenario
kyleTest2	sm.11111a000ccc1a0ccc11ababab000.2	Scenario
name1		Lexicon
name2		Lexicon

Filter

Import Artifacts

LABEL SETS  
 MODELS  
 LEXICONS  
 SCENARIOS  
 REPOSITORY

FIG. 22

**€**

LABEL SETS

MODELS

LEADONS

SCENARIOS

FB

REPOSITORY

DATASETS

NOTECARDS

INQUESTIONS

Datasets	Documents	Other	
Datasets imported	0		<input type="checkbox"/>
Demo-dataset	0		<input type="checkbox"/>
Test	9	sentence 3,539	<input type="checkbox"/>
ds-chats	0		<input type="checkbox"/>
non-english emails	6	sentence 11,182,145	<input type="checkbox"/>
presales-chats	33	sentence 306	<input type="checkbox"/>
reply-chains-emails	2	sentence 542	<input type="checkbox"/>
		sentence 36	<input type="checkbox"/>

**Create Dataset**

**Create Model Project**

**INSTRUCTIONS**

Copilot imports data from your Syntheys knowledge graphs

Please make sure you've imported at least one to get started.

You will need to use the CLI method to import a KG; documentation for this can be found at [Copilot documentation under 'Advanced Dataset Imports'](#)

FIG. 23A

Labeled Datasets > Change of Venue

**Change of Venue Labeled Dataset**

**Examples**  
1  
Total Examples

**Manage Examples**

**Keywords**  
0  
Total Keywords

**Manage Keywords**

**Labelled Samples**  
0  
Total Labelled Samples

**Review**

**Datasets**

Test (0 samples)

reply-chains-emails (2 samples)

Add a dataset...

**Language**  
English

**Labels**

Yes 1

No 2

Add Label

**SETTINGS**

Sampling window: document

Pre-trained model: None

Labelled sample word embeddings: Default English Binary

Related word embeddings: Default English Binary

Ignore samples when advancing

Inject whitespace on token boundaries

**SAMPLING SETTINGS**

Include negative samples

Threshold: 0.5

Value must be between 0 and 1

Sampling Ratio: 1

Value must be between 1 and 10,000

Label New Samples

Label New Samples

Review

Upload

FIG. 23B

Labeled Datasets > Demo dataset > Annotate

### Demo dataset Annotation

Random samples you have not yet encountered

1 of 6 < >

LABELS	
Yes	1 2
No	2 2
Cross Set	C
Total 4	

Keyword search

RELATED WORDS

Search for a keyword to see related words

**Train Now**

FIG. 23C

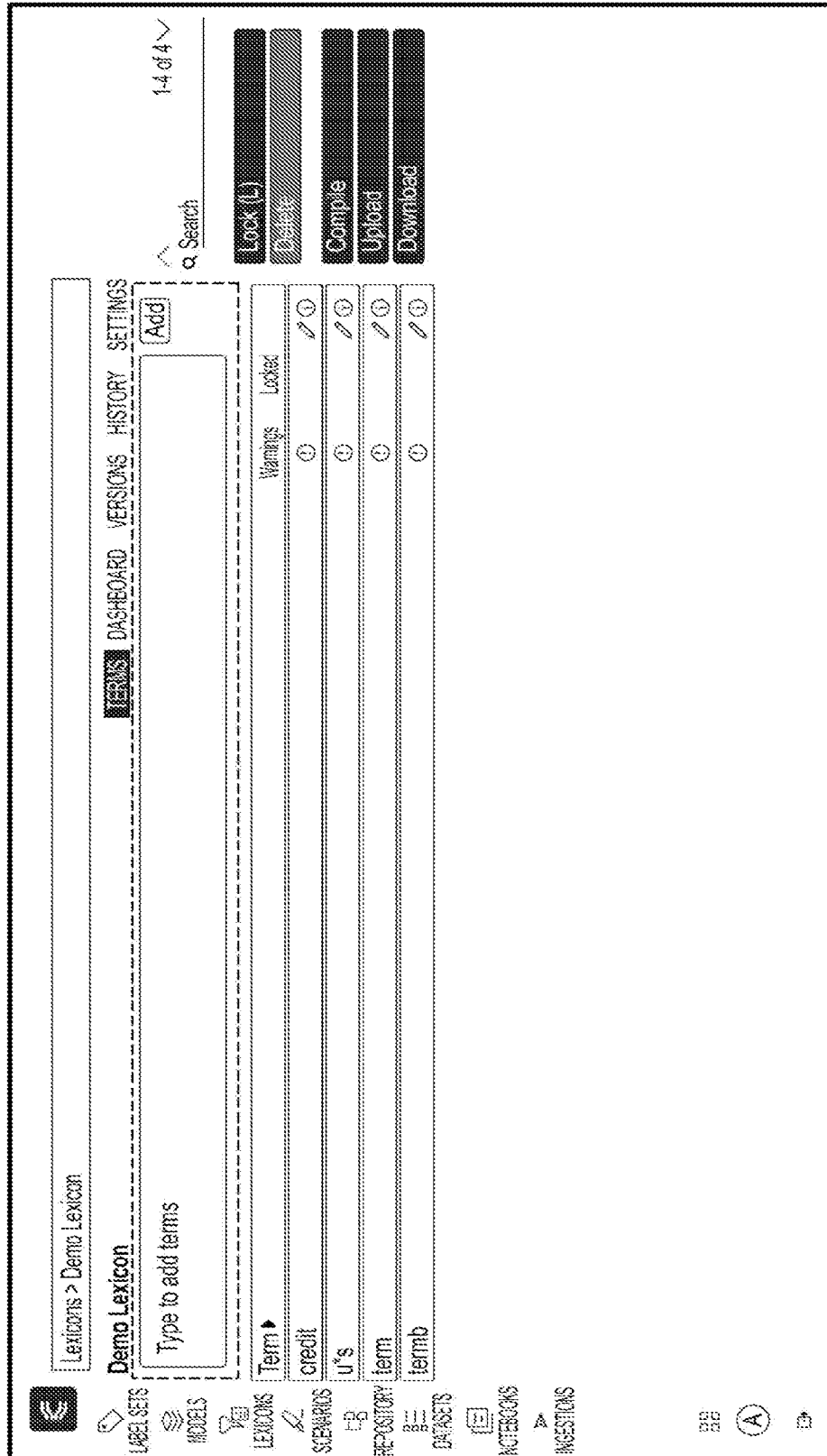


FIG. 23D

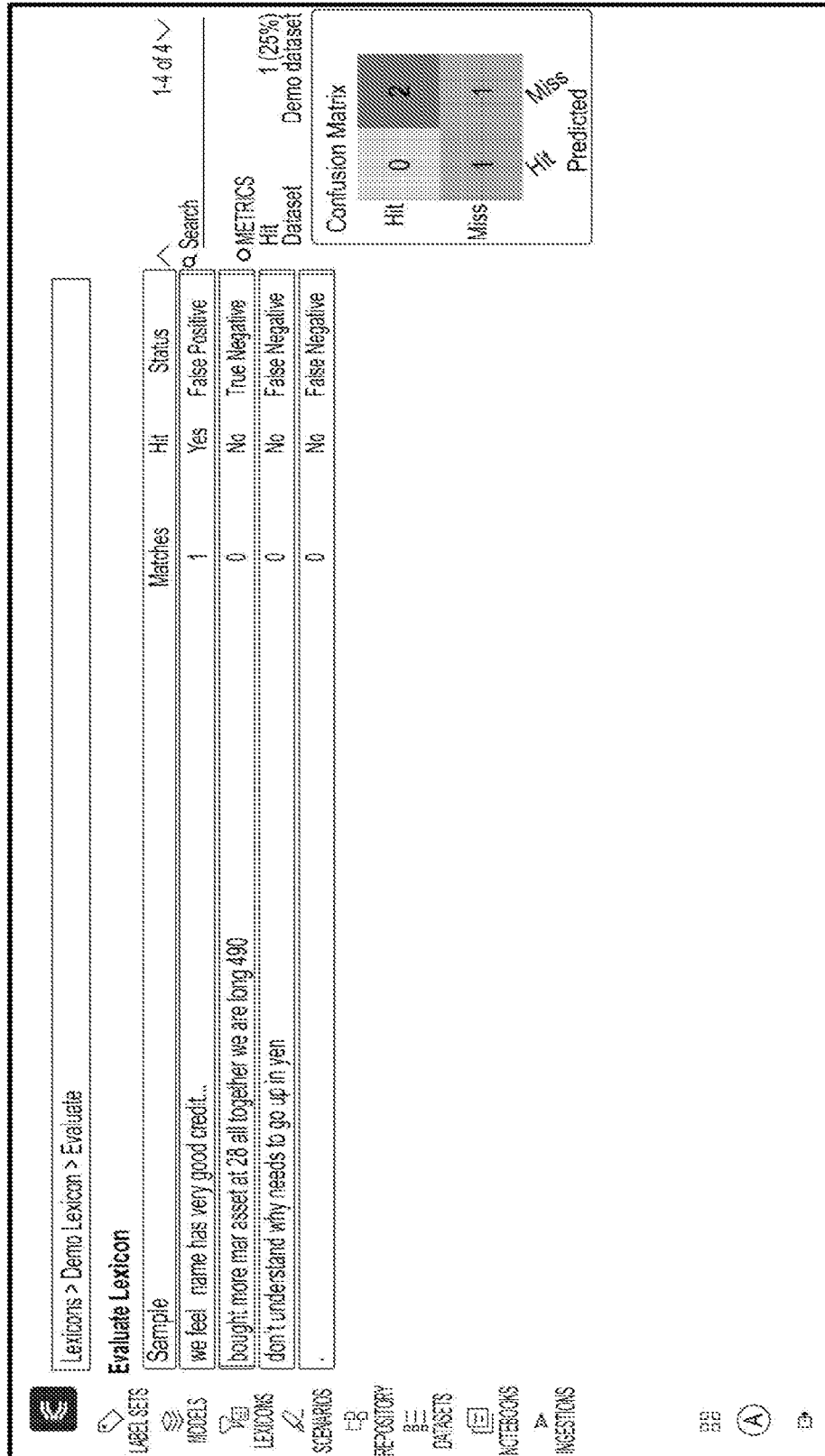


FIG. 23E

Scenarios > Demo Scenario

**DEFINITION** DASHBOARD VERSIONS HISTORY

Demo Scenario

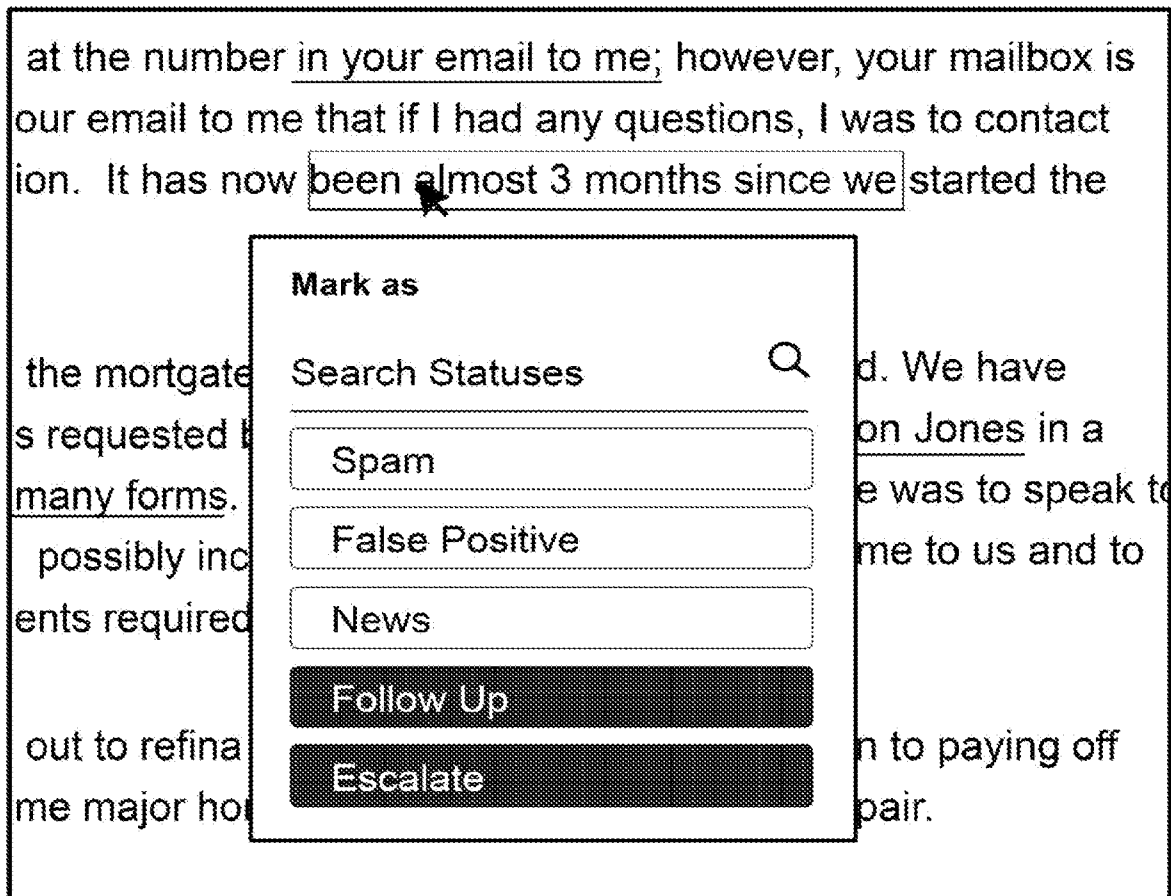
[ALL OF]

- = > [ANY: span] ⊗
- = [ALL OF] ⊗
- = [NOT] ⊗
- = > [ANY: context] ⊗
- = > Lexicon Classifier: Boasting v1.4.0 v1.2.1 ⊗
- = > Number of Recipients: >5 v1.0.0 ⊗

Models	
Lexicons	
<b>Matches</b>	
Attachment File Extension Matches	v1.0.0 ⊗
Attachment Name Matches	v1.0.0 ⊗
800 Not Empty	v1.0.0 ⊗
Communications Type	v1.1.0 ⊗
Contains Encrypted Attachment	v1.0.0 ⊗
Contains Image	v1.0.0 ⊗
Other Skills	
Boolean	
Score	

LABEL SETS   
 MODELS   
 LEXICONS   
 SCENARIOS   
 REPORTING   
 DATASETS   
 NOTECARDS   
 INGESTIONS

FIG. 23F



**FIG. 24**

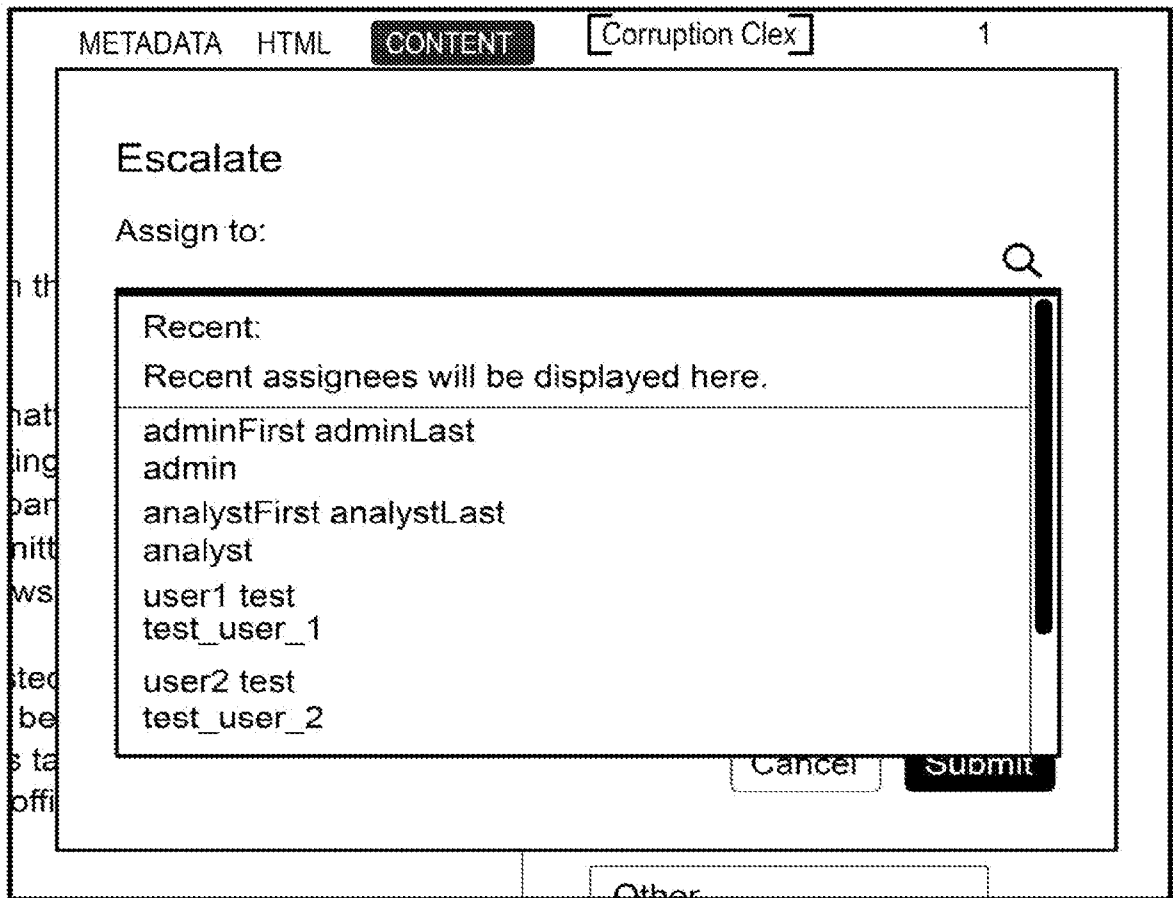


FIG. 25

<b>Assignments</b>	Total Remaining
Assigned to you.	6
Assigned to you only.	

**FIG. 26**