



(19) **United States**

(12) **Patent Application Publication**

Burns et al.

(10) **Pub. No.: US 2025/0258907 A1**

(43) **Pub. Date: Aug. 14, 2025**

(54) **SCANNING AND DETECTING THREATS IN MACHINE LEARNING MODELS**

(71) Applicant: **HiddenLayer, Inc.**, Austin, TX (US)

(72) Inventors: **Tanner Burns**, Austin, TX (US); **Chris Sestito**, Austin, TX (US); **James Ballard**, Fredericksburg, VA (US); **Thomas Bonner**, Steeple Claydon (GB); **Marta Janus**, Twickenham (GB); **Eoin Wickens**, Skibbereen (IE)

(21) Appl. No.: **19/196,620**

(22) Filed: **May 1, 2025**

Related U.S. Application Data

(63) Continuation of application No. 18/113,444, filed on Feb. 23, 2023, now Pat. No. 12,314,380.

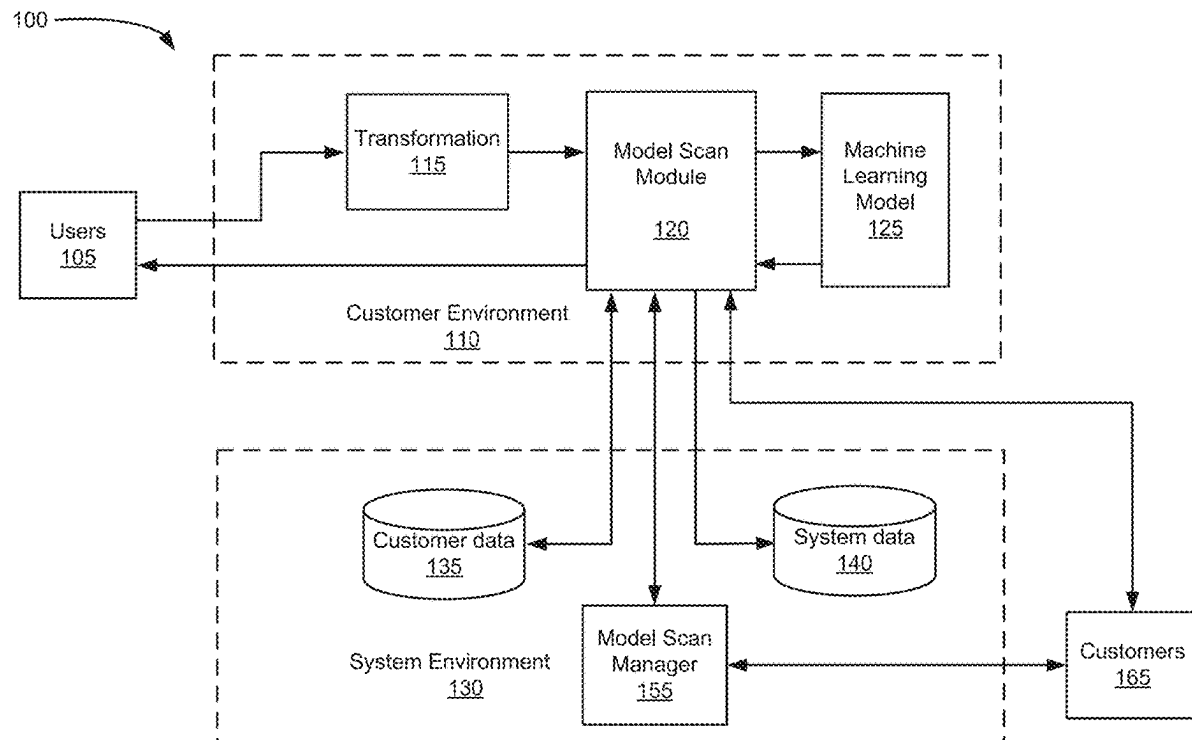
Publication Classification

(51) **Int. Cl.**
G06F 21/53 (2013.01)
G06F 21/56 (2013.01)

(52) **U.S. Cl.**
CPC **G06F 21/53** (2013.01); **G06F 21/56** (2013.01)

(57) **ABSTRACT**

A machine learning model is scanned to detect actual or potential threats. The threats can be detected before execution of the machine learning model or during an isolated execution environment. The threat detection may include performing a machine learning file format check, vulnerability check, tamper check, and stenography check. The machine learning model may also be monitored in an isolated environment during an execution or runtime session. After performing a scan, the system can generate a signature based on actual, potential, or absence of detected threats.



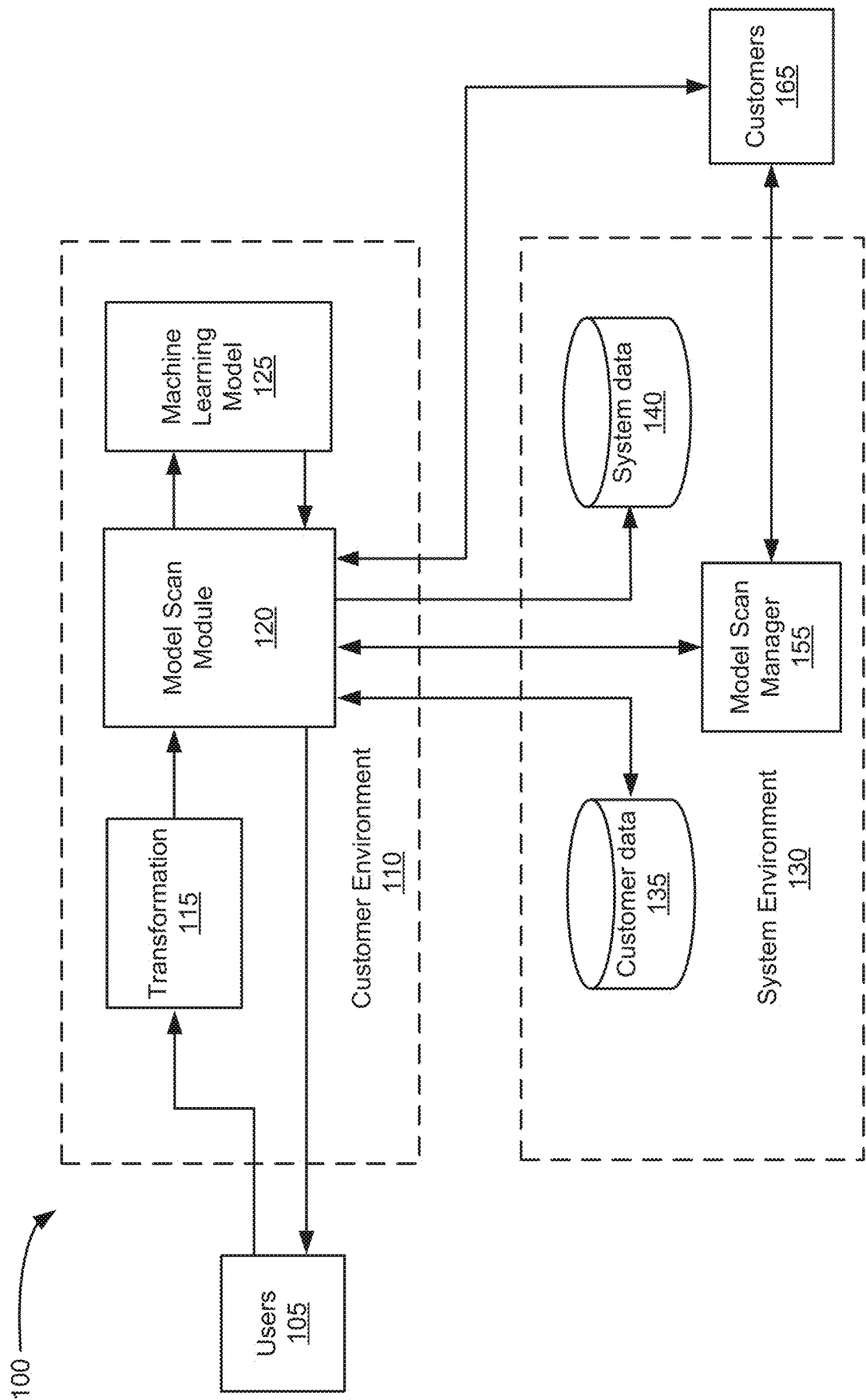


FIGURE 1

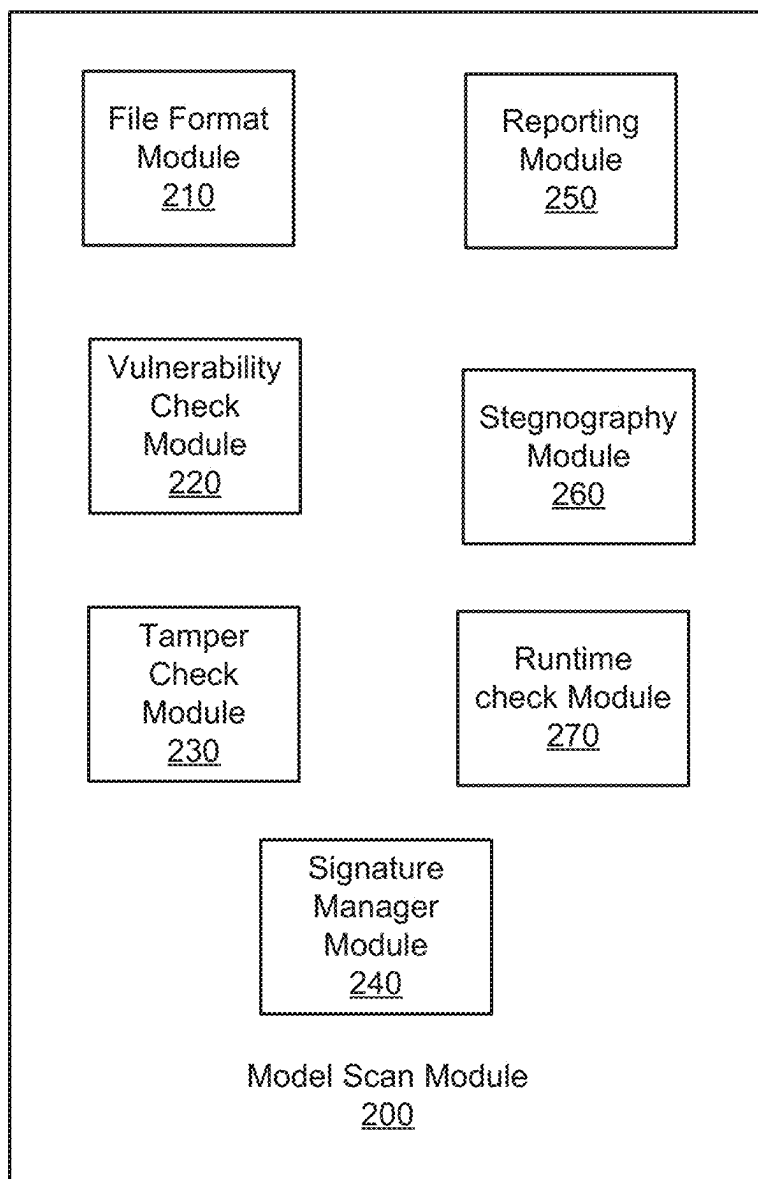
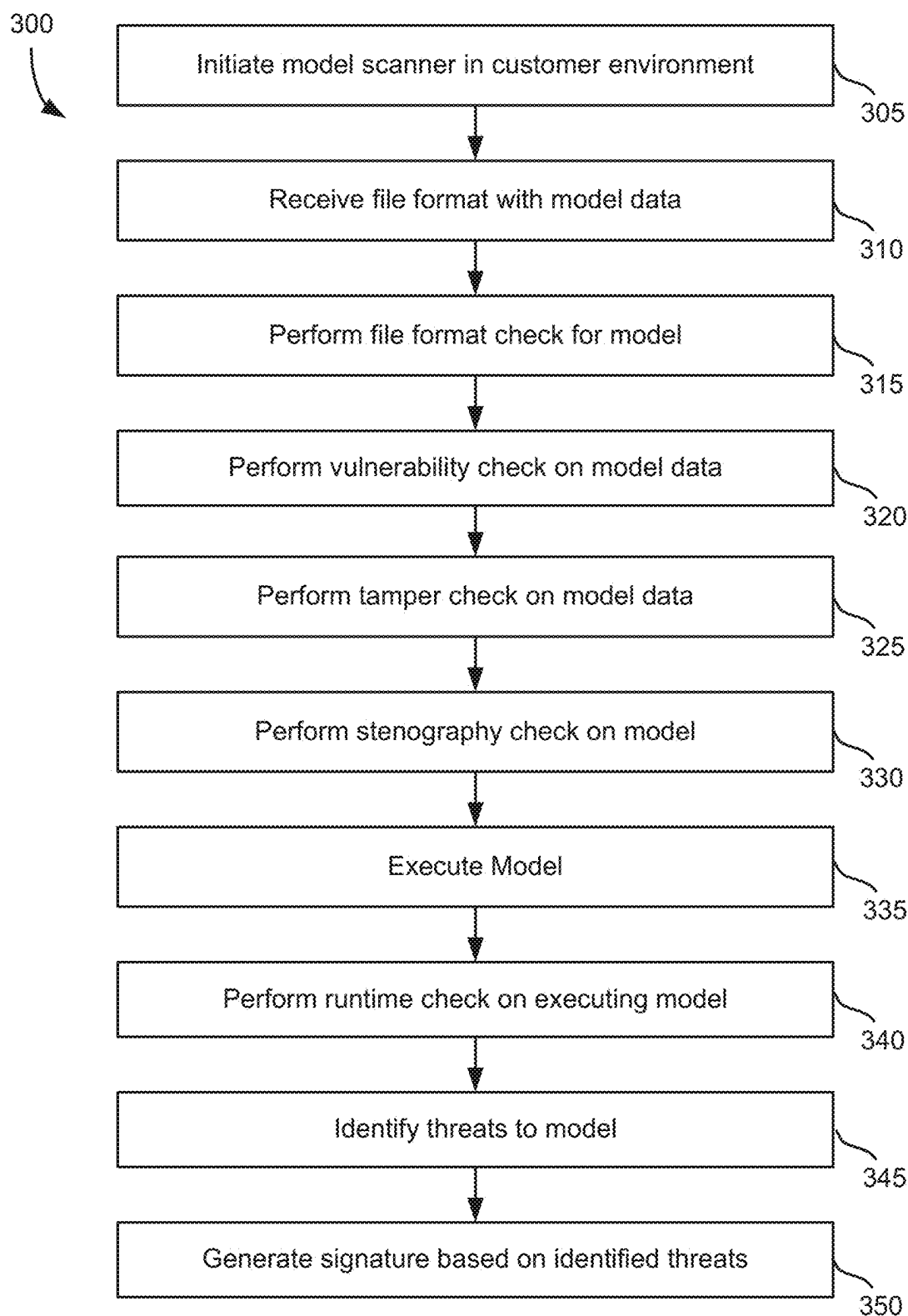


FIGURE 2

**FIGURE 3**

325

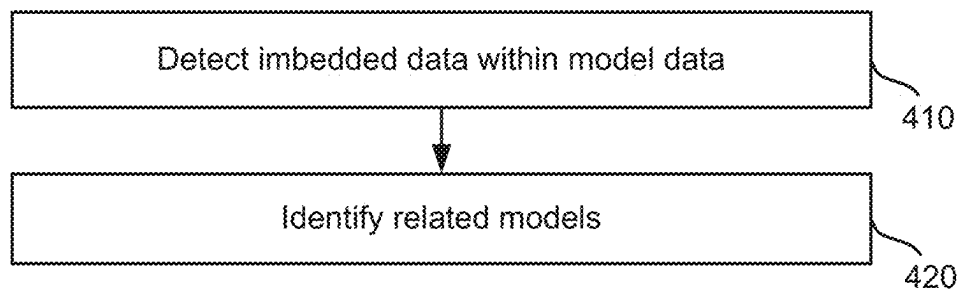


FIGURE 4

410
↓

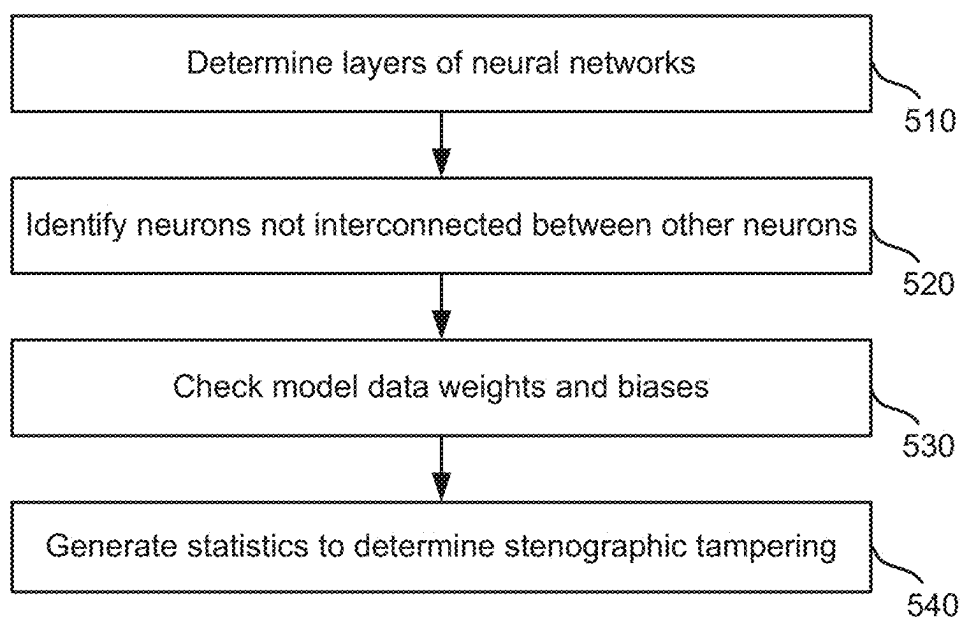
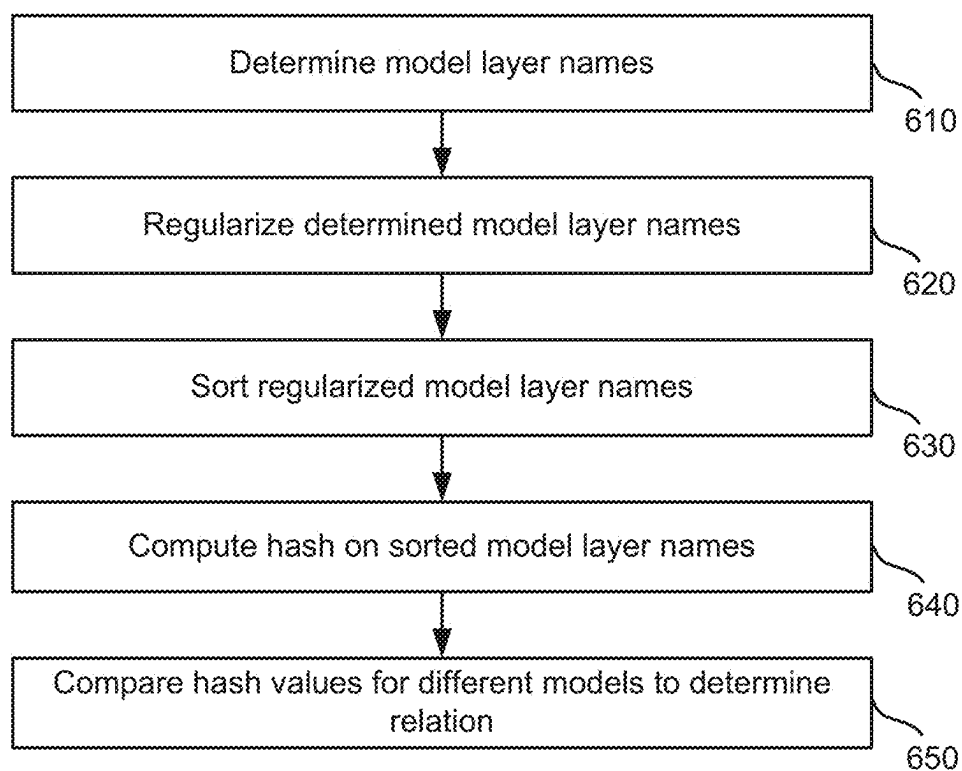


FIGURE 5

420

FIGURE 6

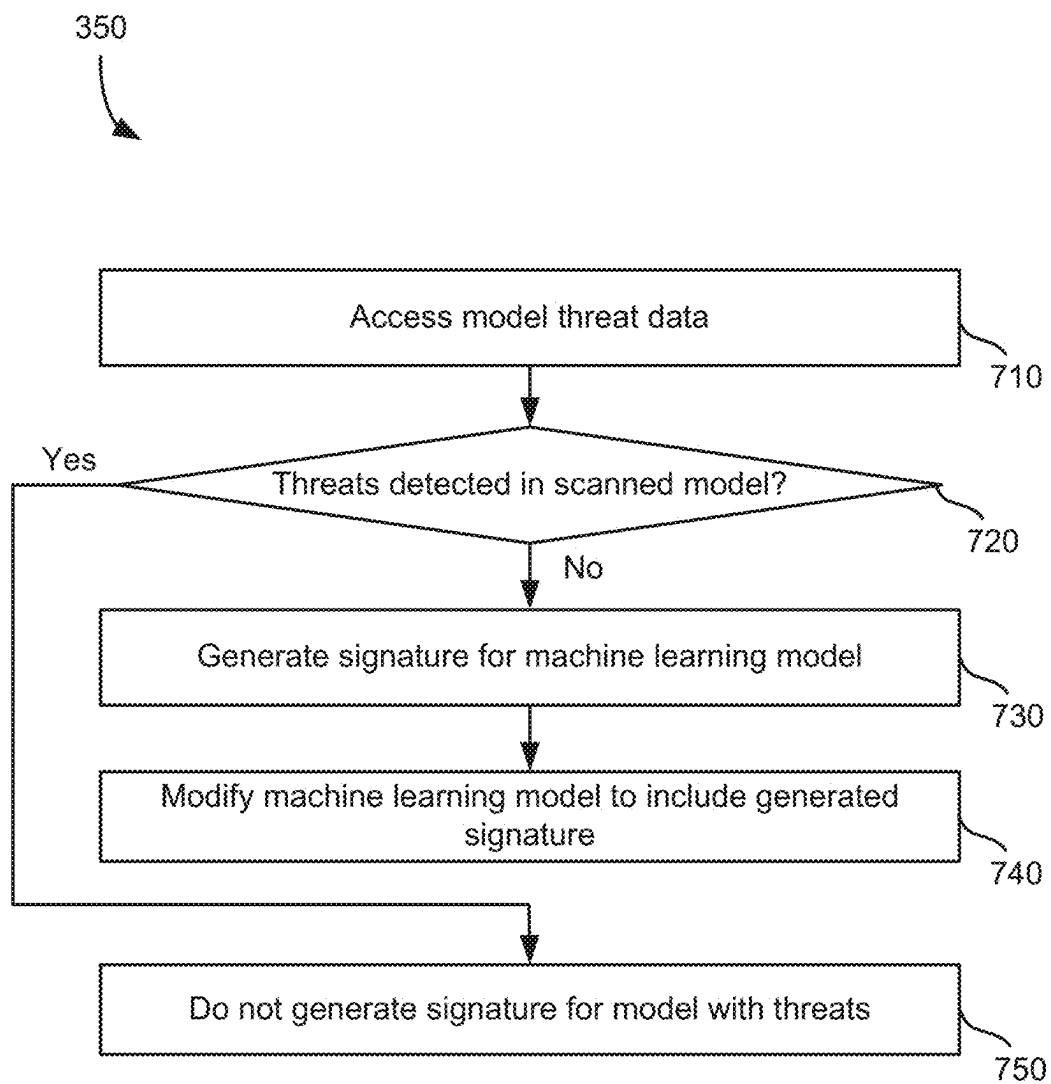


FIGURE 7

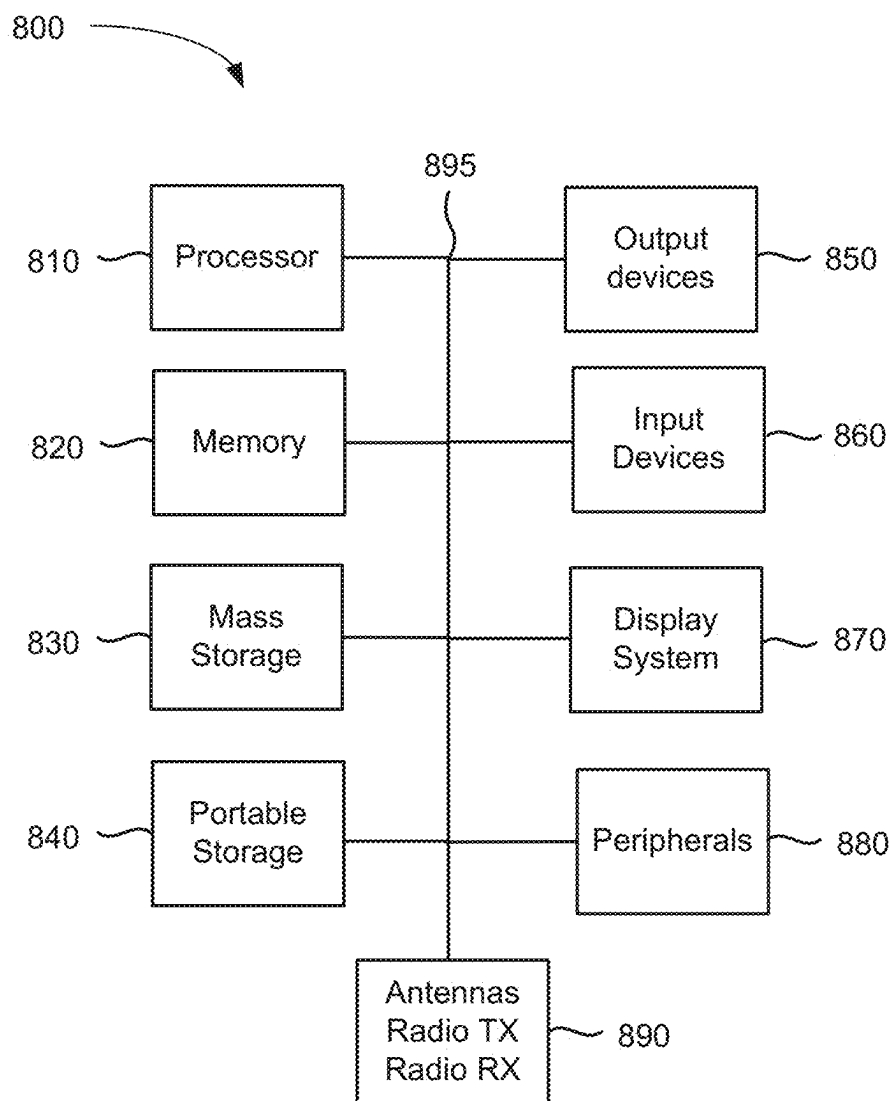


FIGURE 8

SCANNING AND DETECTING THREATS IN MACHINE LEARNING MODELS

RELATED APPLICATION

[0001] The current application claims priority to U.S. patent application Ser. No. 18/113,444 filed on Feb. 23, 2023, the contents of which are hereby fully incorporated by reference.

BACKGROUND

[0002] Machine learning computing resources are becoming more popular in products and computing systems. With the increased presence of machine learning model resources, so have the attacks perpetrated on machine learning-based systems by bad actors. Usually, threats to machine learning models are not discovered until the threat materializes and damage is done to the model or the system in which it operates. As such, what is needed is an improved method for detecting threats on machine learning systems.

SUMMARY

[0003] The present technology, roughly described, scans a machine learning model for threats. The threats can be detected before execution of the machine learning model or during execution in an isolated environment, such as for example a sandbox. The threat detection may include performing a machine learning file format check, vulnerability checks, tamper checks, and/or stenography checks. The machine learning model may also be monitored in an isolated environment during an execution or runtime session.

[0004] After performing a scan, the scanning system generates a set of results based on the threat detection operations and identifies any actual or potential threats that exist in the machine learning model. The system can generate a signature based on actual, potential, or lack of detected threats. For example, if no threats are detected, the system may generate a signature and embed the signature within the model data or file format of the machine learning model.

[0005] In some instances, the present technology provides a method for scanning a machine learning model for threats. The method begins by receiving data for a machine learning model by a scanning module, where the scanning module is stored on a first server, and the data is associated with model parameters and received before execution of the machine learning model. The method continues with performing a plurality of checks by the scanning module based on the received machine learning model data, wherein the checks are performed while the machine learning model is not executing. Next, the scanning module identifies whether the machine learning model includes a threat within the machine learning model based on one or more results of the plurality of checks. The one or more results are reported to a remote computing device by the scanning module.

[0006] In some instances, a non-transitory computer readable storage medium includes embodied thereon a program, the program being executable by a processor to perform a method for scanning a machine learning model for threats. The method begins with scanning a receiving data for a machine learning model by a scanning module, where the scanning module is stored on a first server, and the data is associated with model parameters and received before execution of the machine learning model. The method

continues with performing a plurality of checks by the scanning module based on the received machine learning model data, wherein the checks are performed while the machine learning model is not executing. Next, the scanning module identifies whether the machine learning model includes a threat within the machine learning model based on one or more results of the plurality of checks. The one or more results are reported to a remote computing device by the scanning module.

[0007] In some instances, a system for scanning a machine learning model for threats includes a server having a memory and a processor. One or more modules can be stored in the memory and executed by the processor to receive data for a machine learning model by a scanning module, the scanning module stored on a first server, the data associated with model parameters and received before execution of the machine learning model, perform a plurality of checks by the scanning module based on the received machine learning model data, the checks performed while the machine learning model is not executing, identify, by the scanning module, whether the machine learning model includes a threat within the machine learning model based on results of the plurality of checks, and report the results to a remote computing device by the scanning module.

BRIEF DESCRIPTION OF FIGURES

[0008] FIG. 1 is a block diagram of a system for scanning a machine learning model.

[0009] FIG. 2 is a block time of a machine learning model scan module.

[0010] FIG. 3 is a flowchart of a method for scanning a machine learning model.

[0011] FIG. 4 is a method for performing a tamper check on machine learning model data.

[0012] FIG. 5 is a method for detecting embedded data within machine learning model data.

[0013] FIG. 6 is a method for identifying related machine learning models.

[0014] FIG. 7 is a method for generating a signature based on identified threats to a machine learning model.

[0015] FIG. 8 is a computing environment for implementing the present technology.

DETAILED DESCRIPTION

[0016] The present technology, roughly described, scans a machine learning model for actual or potential threats. The threats can be detected before execution of the machine learning model or during execution in an isolated execution environment. The threat detection may include performing a machine learning file format check, vulnerability check, tamper check, and/or stenography check. The machine learning model may also be monitored in an isolated environment during an execution or runtime session.

[0017] After performing a scan, the scanning system generates a set of results based on the threat detection operations and identifies any actual or potential threats that exist in the machine learning model. The system can generate a signature based on actual, potential, or lack of detected threats. For example, if no threats are detected, the system may generate a signature and embed the signature within the model data or file format of the machine learning model.

[0018] FIG. 1 is a block diagram of a system for performing a model scan. The block diagram of FIG. 1 includes

customer environment 110 and system environment 130. Customer environment 110 includes transformation module 115, model scan module 120, and machine learning model 125. Customer environment 110 may receive and communicate data with one or more users 105, system environment 130, and customers 165. System environment 130 may include customer data 135, system data 140, and model scan manager 155.

[0019] The model scan module 120 can scan the machine learning model 125 before the module is used by users 105. When used by users, the users may provide a stream of data, such as timeseries data, generalized input, or some other data type, to transformation module 115. The transformation module may convert the received timeseries into a series of vectorized data. In some instances, the vectorized data may include an array of float numbers. The vectorization of the received data is then provided to machine learning model 125 for processing. After processing the vectorized data, machine learning model will provide an output, such as a prediction, intended for the requesting user 105.

[0020] Model scan module 120 may collect the vectorized data provided by transformation 115 as well as the output provided by machine learning model 125. The model scan module 120 may then scan the model input, model output, and aspects of the machine learning model itself. In some instances, the model scan module may perform a file format check, vulnerability check, stenography, detect tampering, and perform runtime scanning of the machine learning model.

[0021] The model scan module 120 may be implemented in several ways. In some instances, a model scan module may be implemented as an API placed between the requesting user and the machine learning model. The API may intercept the request, and then send the request to the machine learning model.

[0022] In some instances, the model scan module may be implemented by an API gateway as well as a proxy application. The API gateway may receive the request, provide the request to the proxy application, which may then forward the request to the machine learning model 125.

[0023] Returning to FIG. 1, system environment 130 includes customer data store 135, system data store 140, and model scan manager 155, and customers 165. Each of customer environment 110 and system environment 130 may be implemented as one or more servers implementing the physical or logical modules 115-125 and 135-155 illustrated in FIG. 1. In some instances, each environment is located in one or more cloud computing environments.

[0024] Environments 110 and 130 may communicate over a network. In some instances, one or more modules may be implemented on separate machines in separate environments which may also communicate over a network. The network may be implemented by one or more networks suitable for communication between electronic devices, including but not limited to a local area network, wide-area networks, private networks, public network, wired network, a wireless network, a Wi-Fi network, an intranet, the Internet, a cellular network, a plain old telephone service, and any combination of these networks.

[0025] The customer data store 135 of FIG. 1 includes data associated with one or more customers. The customer data stored includes data for customers of the present system, and may be accessed by any module within system environment 130.

[0026] System data 140 includes data related to system environment 130. System data may include event data, traffic data, timestamp data, and other data. The data may be accessed by any of modules 145-160, and may be used to generate one or more dashboards for use by customers 165.

[0027] Model scan manager 155 may communicate with model scan module 120 and provide output regarding one or more scanned machine learning modules to customers 165. In some instances, some or all of the functionality described with respect to model scan module 120 may be implemented at least in part by model scan manager 155. For example, model scan manager may provide output in the form of a dashboard regarding the status of one or more scans of one or more machine learning models 125.

[0028] FIG. 2 is a block diagram of a model scan module. Model scan module 200 provides more detail for model scan module 125 of FIG. 1. Model scan module 200 includes file format module 210, vulnerability check module 220, temperature module 230, signature manager module 240, reporting module 250, stenography module 260, and runtime check module 270.

[0029] File format module 210 may detect and process file formats of a machine learning model. A file format may include model data, such as for example numeric weights and other model data. In some instances, the file format data is a wrapper around model data to convey data, for example in .zip, pickle, protobuf, HDF, or JSON format.

[0030] Vulnerability check module 220 may analyze and determine vulnerability of a machine learning module. A vulnerability check may be based on the machine learning model type, layers, changes in relation to a base model topology, or level of dormant activity.

[0031] Tamper check module 230 may determine whether a machine learning model has been tampered with. In some instances, detecting tampering of a machine learning model may be based on an analysis of model data such as layers or weights, or on identifying related models.

[0032] Reporting module 250 may report scan status and any actual or potential threats identified for a machine learning model. Signature manager module 240 may determine whether a machine learning model includes a previously generated signature. Module 240 may also generate and embed a signature within a machine learning model.

[0033] FIG. 3 is a flowchart of a method for performing a model scan. A model scanner is initiated in a customer environment at step 305. The Model scanner may be installed in a customer environment such that it may access the machine learning model code, inputs and outputs, within the environment. The machine learning model file format may be received by the model scanner at step 310. The file format may include model data, such as machine learning model layer data, weights, and other data. The file format check may be performed for the model at step 315.

[0034] A vulnerability check may be performed on the model data at step 320. The vulnerability check may be applied to neural net, random forest, and other machine learning model types. For example, with a neural net machine learning model, there may be parameters that specify its policy, number of players, units per layer, the nature of an interconnect player, and other parameters. The vulnerability check may determine how closely related the scanned model is to a typical base model. For example, most machine learning models are fine-tuned versions of a particular base model. The scanner, in some instances, may

determine if the fine-tuned machine learning model has the same topology architecture of the base model. Determining a vulnerability check may include detecting changes to the based apology and parameters, detecting backdoors, and dormant portions and activated upon a certain input. Performing a tamper check may include checking model data layers and specific weights to determine if they may suggest a threat.

[0035] A tamper check may also involve identifying related models, such as determining if a current machine learning model is similar to a previous model in which a threat was detected. Performing a tamper check is discussed in more detail below with respect to FIGS. 4 through 6.

[0036] A stenography check may be performed on the machine learning model at step 330. A stenography check may analyze the lower order bits and other aspects of the machine learning model data. For example, after processing the file format and compression, the model scanner can obtain actual weights and biases used in the model. The lower order bits in each individual number have little impact in resulting calculations and can be zeroed out through compression. Lower significant bits are typically random. If the scan model extracts the lower bits and, after studying the distribution, determines they are not random, it may be identified as a potential or actual threat based on tampering, depending on the degree of randomness or lack of randomness.

[0037] In some instances, performing a stenography check may include analyzing entropy associated with a machine learning model. A baseline may be determined for certain types of models. Clustering can be used to understand model types, baselines, identify predicted entropy, detect the actual entropy system to compare the predicted statistic to the actual entropy statistic to determine if malware is within the model at static runtime.

[0038] The machine learning model being scanned may be executed at step 335. During execution, a runtime check may be performed on the executing model at step 340. In some instances, the model may be executed in a sandbox or other isolated environment or virtual machine. While executing, the scanning tool may identify executing shell scripts, and other processes in the isolated environment, where no processes should start by themselves. In some instances, loading the machine learning model in memory could cause malicious code to execute and be identified as a threat.

[0039] After performing the various checks and analyses on the machine learning model, model and information and threats to the model are identified at step 345. Information from the scan may include the time of the scan, scan failures and successes, the model file scanned, as was the threats detected. A threat may be a condition or other aspect of a machine learning model that allows the model to be manipulated, controlled, or otherwise affected in an unintended and/or undesired way.

[0040] A signature may be generated based on the identified threats, if any, at step 350. Generally, a machine learning model does not include a signature. The scanner tool may generate a signature if no threats are detected and embed the signature within the machine learning model data. In some instances, manipulation of the machine learning model file may invalidate the signature. In some instances, manipulation of any material portion of the machine learning model will invalidate the signature. Hence,

if a signature is generated, and the model subsequently modified, it will be detectable using the generated signature.

[0041] In some instances, any of the checks in steps 410-430 and 440 may be performed using machine learning models. The machine learning model may receive input in the form of file format data, layer data, or other data discussed with respect to steps 410-430 and 440 for the model being scanned. The machine learning model receiving the data may generate a prediction or other indication as to whether the scanned model likely includes a threat. The output may then be used as part of step 445 to determine whether the scanned model includes a threat. For example, for a prediction value ranging from 0 to 1, a value above 0.6 may indicate a threat, while a value of 0.3 to 0.6 may indicate a potential threat.

[0042] FIG. 4 is a method for performing a tamper check on model data. The method of FIG. 4 provides more data for step 325 of the method of FIG. 3. Embedded data within the model data may be detected at step 410. Detecting embedded data may be based on analysis of neural network layers, neuron interconnections, and model data weights and biases. More detail for detecting embedded data is discussed with respect to the method of FIG. 5.

[0043] Related machine learning models may be identified at step 420. Identifying related models may be useful to compare a new model with a previously known infected model or a model detected to have one or more threats. More detail for identifying related models is discussed with respect to the method of FIG. 6.

[0044] FIG. 5 is a method for detecting embedded data within model data. The method of FIG. 5 provides more detail for step 410 of the method of FIG. 4. First, layers of neural networks may be determined at step 510. Neurons not interconnected with other neurons are identified at step 520. For example, a neuron that is connected to a point completely outside the machine learning model may be identified as a potential threat.

[0045] Model data weights and biases may be checked at step 530. Model data weights and biases may be used to generate statistics at step 540. The statistics are generated to determine stenographic tampering. For example, the entropy of a machine learning model may be used to determine whether a model has a threat. To identify a threat based on entropy, a baseline entropy is determined for the machine learning model. The baseline entropy can be determined for certain types of machine learning models. Clustering can then be performed to understand the machine learning model type being used, identify predicted entropy for the particular model, and to detect the actual entropy for the particular machine learning model being scanned. The Scanning module can compare the actual entropy to the predicted entropy to identify if malware is within the model during static runtime.

[0046] FIG. 6 is a method for identifying related models. The method of FIG. 6 provides more detail for step 420 of the method of FIG. 4. First, model layer names are determined at step 610. The model layer names may be regularized at step 620. The regularization of the model names may include removing spaces, periods, and making the names in lowercase letters. The regularized model layer names may be sorted at step 630. In some instance, the regularized names are alphabetically sorted. A hash may then be performed on the sorted model layer names at step 640. The hash generated using any hashing technique suitable to

generate a comparable value. Hash values for different models may be compared to determine a relationship at step 650. For example, if a present model hash value is similar to the hash value of a machine learning model known to have threats, the present model may be treated as having a similar threat

[0047] FIG. 7 is a method for generating a signature based on identified threats. The method of FIG. 7 provides more detail for step 350 of the method of FIG. 3. First, model threat data is accessed at step 710. A determination is then made as to whether the data generated based on the machine learning model scan includes any threats indicating that the model has viable threats.

[0048] If the threat data indicates that threats do exist in the machine learning model, a signature is not generated for the model at step 750. If the threat data indicates that no threat is associated with the machine learning model, a signature is generated for the machine learning model at step 730. Machine learning models and modified to include the generated teacher at step 740.

[0049] FIG. 8 is a block diagram of a computing environment for implementing the present technology. System 800 of FIG. 8 may be implemented in the contexts of the likes of machines that implement model scan module 120, customer data store 135, system data store 140, and model scan manager 155. The computing system 800 of FIG. 8 includes one or more processors 810 and memory 820. Main memory 820 stores, in part, instructions and data for execution by processor 810. Main memory 820 can store the executable code when in operation. The system 800 of FIG. 8 further includes a mass storage device 830, portable storage medium drive(s) 840, output devices 850, user input devices 860, a graphics display 870, and peripheral devices 880.

[0050] The components shown in FIG. 8 are depicted as being connected via a single bus 890. However, the components may be connected through one or more data transport means. For example, processor unit 810 and main memory 820 may be connected via a local microprocessor bus, and the mass storage device 830, peripheral device(s) 880, portable storage device 840, and display system 870 may be connected via one or more input/output (I/O) buses.

[0051] Mass storage device 830, which may be implemented with a magnetic disk drive, an optical disk drive, a flash drive, or other device, is a non-volatile storage device for storing data and instructions for use by processor unit 810. Mass storage device 830 can store the system software for implementing embodiments of the present invention for purposes of loading that software into main memory 820.

[0052] Portable storage device 840 operates in conjunction with a portable non-volatile storage medium, such as a floppy disk, compact disk or Digital video disc, USB drive, memory card or stick, or other portable or removable memory, to input and output data and code to and from the computer system 800 of FIG. 8. The system software for implementing embodiments of the present invention may be stored on such a portable medium and input to the computer system 800 via the portable storage device 840.

[0053] Input devices 860 provide a portion of a user interface. Input devices 860 may include an alpha-numeric keypad, such as a keyboard, for inputting alpha-numeric and other information, a pointing device such as a mouse, a trackball, stylus, cursor direction keys, microphone, touchscreen, accelerometer, and other input devices. Additionally, the system 800 as shown in FIG. 8 includes output devices

850. Examples of suitable output devices include speakers, printers, network interfaces, and monitors.

[0054] Display system 870 may include a liquid crystal display (LCD) or other suitable display device. Display system 870 receives textual and graphical information and processes the information for output to the display device. Display system 870 may also receive input as a touchscreen.

[0055] Peripherals 880 may include any type of computer support device to add additional functionality to the computer system. For example, peripheral device(s) 880 may include a modem or a router, printer, and other device.

[0056] The system of 800 may also include, in some implementations, antennas, radio transmitters and radio receivers 890. The antennas and radios may be implemented in devices such as smart phones, tablets, and other devices that may communicate wirelessly. The one or more antennas may operate at one or more radio frequencies suitable to send and receive data over cellular networks, Wi-Fi networks, commercial device networks such as a Bluetooth device, and other radio frequency networks. The devices may include one or more radio transmitters and receivers for processing signals sent and received using the antennas.

[0057] The components contained in the computer system 800 of FIG. 8 are those typically found in computer systems that may be suitable for use with embodiments of the present invention and are intended to represent a broad category of such computer components that are well known in the art. Thus, the computer system 800 of FIG. 8 can be a personal computer, handheld computing device, smart phone, mobile computing device, workstation, server, minicomputer, mainframe computer, or any other computing device. The computer can also include different bus configurations, networked platforms, multi-processor platforms, etc. Various operating systems can be used including Unix, Linux, Windows, Macintosh OS, Android, as well as languages including Java, .NET, C, C++, Node.JS, and other suitable languages.

[0058] The foregoing detailed description of the technology herein has been presented for purposes of illustration and description. It is not intended to be exhaustive or to limit the technology to the precise form disclosed. Many modifications and variations are possible in light of the above teaching. The described embodiments were chosen to best explain the principles of the technology and its practical application to thereby enable others skilled in the art to best utilize the technology in various embodiments and with various modifications as are suited to the particular use contemplated. It is intended that the scope of the technology be defined by the claims appended hereto.

1. A method for implementation by one or more computing devices for scanning a machine learning model for threats comprising

receiving, by a scanning module, data for a machine learning model which is associated with model parameters and is received before execution of the machine learning model;

performing, by the scanning module, a plurality of checks based on the received machine learning model data, the checks performed while the machine learning model is not executing, at least one of the checks identifying differences between the machine learning model and a corresponding base machine learning model;

identifying, by the scanning module, whether the machine learning model includes a threat within the machine learning model based on results of the plurality of checks; and

reporting, by the scanning module, the results to a remote computing device.

2. The method of claim 1 further comprising:

generating a certificate based on determining that the machine learning model does not include a threat; and embedding the generated certificate within the data for the machine learning model.

3. The method of claim 1, wherein the plurality of checks includes two or more of a file format check, a vulnerability check, a tampering check, and a stenography check.

4. The method of claim 1 further comprising:

executing the machine learning model in an isolated environment;

monitoring the execution of the machine learning in the isolated environment to detect suspicious activity.

5. The method of claim 1, wherein performing a plurality of checks comprises:

determining an expected entropy for the machine learning model;

determining the actual entropy for the machine learning model;

calculating the difference between the expected entropy and the actual entropy; and

identifying a threat if the difference in expected entropy and actual entropy exceeds a threshold.

6. The method of claim 1, wherein an identified threat may be a potential threat or an actual threat.

7. The method of claim 1, wherein the scanning module is stored on a server forming part of a computing environment that includes the machine learning model.

8. A method for implementation by one or more computing devices for scanning a machine learning model for threats comprising:

receiving, by a scanning module, data for a machine learning model which is associated with model parameters and is received before execution of the machine learning model;

performing, by the scanning module, a plurality of checks based on the received machine learning model data, the checks performed while the machine learning model is not executing, at least one of the checks identifying differences between the machine learning model and a previously known infected machine learning model;

identifying, by the scanning module, whether the machine learning model includes a threat within the machine learning model based on results of the plurality of checks; and

reporting, by the scanning module, the results to a remote computing device.

9. The method of claim 8 further comprising:

generating a certificate based on determining that the machine learning model does not include a threat, and embedding the generated certificate within the data for the machine learning model.

10. The method of claim 8, wherein the plurality of checks includes two or more of a file format check, a vulnerability check, a tampering check, and a stenography check.

11. The method of claim 8, wherein the operations further comprise: executing the machine learning model in an isolated environment, and monitoring the execution of the machine learning in the isolated environment to detect suspicious activity.

12. The method of claim 8 further comprising:

determining an expected entropy for the machine learning model;

determining the actual entropy for the machine learning model;

calculating the difference between the expected entropy and the actual entropy; and

identifying a threat if the difference in expected entropy and actual entropy exceeds a threshold.

13. The method of claim 8, wherein an identified threat may be a potential threat or an actual threat.

14. The method of claim 8, wherein the first server which stores the scanning module is within a computing environment that includes the machine learning model.

15. A system for scanning a machine learning model for threats comprising:

at least one data processor; and

memory storing instructions which, when executed by the at least one data processor, results in operations comprising:

receiving, by a scanning module, data for a machine learning model comprising model parameters, a vectorized input to the machine learning model, and an output generated by the machine learning model in response to ingesting the vectorized input;

performing, by the scanning module, a plurality of checks based on the received machine learning model data including an input to the machine learning model and a resulting output;

identifying, by the scanning module, whether the machine learning model includes a threat within the machine learning model based on results of the plurality of checks; and

reporting, by the scanning module, the results to a remote computing device.

16. The system of claim 15, wherein the operations further comprise: generating a certificate based on determining that the machine learning model does not include a threat, and embedding the generated certificate within the data for the machine learning model.

17. The system of claim 15, wherein the plurality of checks includes two or more of a file format check, a vulnerability check, a tampering check, and a stenography check.

18. The system of claim 15, wherein the operations further comprise: executing the machine learning model in an isolated environment, and monitoring the execution of the machine learning in the isolated environment to detect suspicious activity.

19. The system of claim 15, wherein performing a plurality of checks comprises:

determining an expected entropy for the machine learning model;

determining the actual entropy for the machine learning model;

calculating the difference between the expected entropy and the actual entropy; and
identifying a threat if the difference in expected entropy and actual entropy exceeds a threshold.

20. The system of claim **15**, wherein an identified threat may be a potential threat or an actual threat.

* * * * *