

## K2's Privacy Issues & Solutions -What You Need To Know!





#### **Course Description**

Generative Artificial Intelligence captured the imaginations of everyone when it was launched in 2022, but they also create opportunities for data leakage and inappropriate use of confidential data based on the terms of service. Unfortunately, everyone has at some point clicked "I agree" on a terms of service (TOS) or end user license agreement (EULA) document which they neither read nor understood. This session will review the common terms embedded in these documents with thousands of words and incomprehensible terminology. Attend this breakout and learn what you wish you knew about analyzing and evaluating license agreements and privacy policies.

🔊 K2 Enterprises





#### **Overview**

- Software licenses and privacy policies
- Artificial intelligence laws and regulations
- Key terms from NIST
- Language models
- AI Risk management

🔊 K2 Enterprises



Copyright 2024, K2 Enterprises, LLC

Copyright © 2024, K2 Enterprises, LLC. Reproduction or reuse for purposes other than a K2 Enterprises' training event is prohibited.









#### End User License Agreement (EULA) / Terms of Service (ToS)

- Document that states the terms which govern your use of an application
- Let software companies define the terms for the use of their software or service in a way which is favorable to their business
- Are often accepted without ever reading them at all
- May incorporate other documents, policies, and terms by reference

Signal K2 Enterprises

#### Word Count/Complexity of EULAs

- Intuit QuickBooks Online
  - 57,053 words
  - Flesch-Kincade grade 14.6
- QB Desktop US
  - 24,619 words
  - Flesch-Kincade grade 14.1
- Xero
  - 5,332 words
  - Flesch-Kincaid Grade 11.8

Many EULAs are long, complex, difficult to understand, and often disclose as few specifics as possible about any practices which might be a concern to end users

Copyright 2024, K2 Enterprises, LLC

#### Key Terms In EULA/ToS Documents

- EULAs are more commonly used for the licensing of software
- ToS are more commonly used for services hosted on websites
- Key terms/clauses include:
  - Terms and scope of license geography, exclusivity
  - Duration of the license
  - Permitted uses and prohibited uses of the application
  - Licensing fees, subscriptions, or royalties for use
  - Disclaimers of liability and warranties
  - Limitation of liability
  - How changes are made to the EULA/ToS

🔊 K2 Enterprises



## Key Terms In EULA/ToS Documents

- Key terms (continued)
  - How to get help with the application
  - Dispute resolution procedures
  - Required binding arbitration
  - Jurisdiction for dispute resolutions
  - Procedures for terminating relationship
  - Consent to send communicate via e-mail, SMS, and telephone
  - Any required indemnifications (guarantees)
  - How data will be shared with others, such as sub-processors of data (e.g., bank feed aggregators)

Copyright 2024, K2 Enterprises, LLC

Signal K2 Enterprises







#### Copyright © 2024, K2 Enterprises, LLC. Reproduction or reuse for purposes other than a K2 Enterprises' training event is prohibited.

#### **Privacy Policy Extracts - Intuit**



- Intuit Privacy Policy last updated 9/6/2023 (retrieved and was current as of 3/22/2024)
  - "<u>The Platform may also include information about or offers for third-party services or products or allow you to connect your account to or otherwise access third-party services or products.</u> Intuit does not warrant, and is not responsible for, such third-party services and products or claims made about them or the actions or inactions of any third party. You must review and comply with any Additional Terms.
    <u>Intuit may be compensated by those third parties</u>, which could impact whether, how and where the services and products are displayed."

Signal K2 Enterprises



#### **Privacy Policy Extracts - Xero**

• 29. Yodlee and bank feeds: Your use of automated bank account feeds enabled by Yodlee from within our services is subject to separate Yodlee terms. If your bank or credit union connects to Xero directly, you may use those feeds instead of Yodlee. Bank feeds are generally offered for free but may have associated charges that we will pass on to you. You may discontinue the use of a bank feed at any time - check out how to stop a bank feed on Xero Central. You can learn more about bank feeds on Xero Central. If you receive a bank feed and are based in the United Kingdom, Xero's additional terms for account information services apply.



Most software providers use sub processors like Yodlee or Fiserv to handle tasks like establishing and maintaining bank feeds. The use of those services may result in information being shared in accordance with additional terms and services

Yodlee's <u>Privacy Notice</u> is an additional 3,812 words, and the organization's <u>Privacy</u> <u>Policy</u> is an additional 2,158 words.

Copyright 2024, K2 Enterprises, LLC

#### ADP's AI Ethics Statement

ADP's privacy statements include an <u>Al</u> <u>ethics statement</u>, which explains things like:

- The organization's approach to AI,
- where AI is used in its products,
- data quality processes,
- principles which are followed when using AI, including "privacy-by-design",
- How AI models are evaluated,
- How AI is governed, and

more – see it <u>on ADP's website</u>

#### ADP's AI Ethics Principles

- 1. Human oversight
- 2. Governance
- 3. Privacy-by-design
- 4. Explainability & transparency
- 5. Data quality
- 6. Culture of responsible AI
- 7. Inclusion and training

🔊 K2 Enterprises



### Acceptable Use Policy

Areas addressed in an AUP commonly include:

- Ownership and general use of technology
- Security and handling confidential information
- Types of unacceptable use of systems
- E-mail and communication guidelines
- Blogging and social media
- Consequences for violation of these policies



#### Signal K2 Enterprises













#### **Executive Order 14110**



Copyright 2024, K2 Enterprises, LLC

- Create cybersecurity program to develop AI tools to find and fix vulnerabilities in critical infrastructure
- Increases government investment in AI and government use of AI
- Actions to provide guidance to landlords, federal benefit programs, and federal contractors to avoid algorithmic discrimination, and provide best practices in the criminal justice system
- Safety program to remedy harms and unsafe healthcare practices using AI
- Transform education by creating resources to support educators with Alenabled education tools
- Create best practices to maximize benefits of AI for workers
- Increase global collaboration surrounding AI

Signal K2 Enterprises

**NIST Trustworthy & Responsible AI Resource Center**  Available online at https://airc.nist.gov Resources include NIST AI Risk Management Framework (AIRMF) • AIRMF Playbook and glossary AIRMF Roadmap Crosswalks to various standards and frameworks Get an introduction to the RMF by watching the video at NIST 🔊 K2 Enterprises Copyright 2024, K2 Enterprises, LLC

> Copyright © 2024, K2 Enterprises, LLC. Reproduction or reuse for purposes other than a K2 Enterprises' training event is prohibited.

#### NIST AI Risk Management Framework

- Issued in January 2023 by National Institutes of Standards and Technology (NIST), the US Federal Government Agency which creates technology, privacy, and security standards for all government departments and agencies
- Currently voluntary, not mandatory, but we expect this to change before the end of 2024



Artificial Intelligence Risk Management Framework (AI RMF 1.0)

> This publication is available free of charge from: https://doi.org/10.6028/NIST.AI.100-1

Copyright 2024, K2 Enterprises, LLC

#### NIST AI Risk Management Framework

- It seeks to incorporate trustworthiness considerations into the design, development, use, and evaluation of AI products, services, and systems
- Four major components
  - Govern: Establishing governance structures and policies for AI risk management
  - Map: Identifying and assessing AI-related risks
  - Measure: Quantifying and evaluating risks
  - Manage: Implementing risk mitigation strategies



🖎 K2 Enterprises

🔊 K2 Enterprises



## NIST AI Risk Management Framework Playbook

- A set of practice aids called the "AI Risk Management Playbook" which provides voluntary suggestions for documenting, identifying, and mitigating risk
  - IMPORTANT This document contains voluntary suggestions and agencies are NOT currently required by NIST to use this tool (as of 3/2024)
- The free tools offer assistance with governance, mapping of risks, risk measurement, and risk management
- Like risk management frameworks from COSO, ISACA (COBIT), and others



#### 🔊 K2 Enterprises



Source: "Artificial Intelligence Risk Management Framework" (AI 100.1) by US National Institutes for Standards and Technology (NIST)

🔊 K2 Enterprises



## <section-header><section-header><list-item><list-item><list-item><list-item><list-item>

## **NIST Definitions**

- AI Model
  - A function that takes features as input and predicts labels as output
  - Machine Learning algorithms and data processing designed, developed, trained and implemented to achieve set outputs, inclusive of datasets used for said purposes unless otherwise stated
- Al System an engineered or machine-based system that can, for a given set of objectives, generate outputs such as predictions, recommendations, or decisions influencing real or virtual environments. Al systems are designed to operate with varying levels of autonomy

#### 🔊 K2 Enterprises



- Adaptive Learning Updating predictive models online during their operation to react to concept drifts
- Artificial Intelligence Learning The ingestion of a corpus, application of semantic mapping, and relevant ontology of structured and/or unstructured data that yields inference and correlation leading to the creation of useful conclusive or predictive capabilities in a given knowledge domain. Strong AI learning also includes the capability of creating unique hypotheses, attributing data relevance, processing data relationships, and updating its own lines of inquiry to further the usefulness of its purpose

Signal K2 Enterprises

#### **NIST Definitions**

• Artificial Narrow Intelligence (ANI) - Artificial Narrow Intelligence, also known as weak or applied intelligence, represents most of the current artificial intelligent systems that usually focus on a specific task. Narrow AIs are mostly much better than humans at the task they were made for: for example, look at face recognition, chess computers, calculus, and translation. The definition of artificial narrow intelligence is in contrast to that of strong AI or artificial general intelligence, which aims at providing a system with consciousness or the ability to solve any problems. Virtual assistants and AlphaGo are examples of artificial narrow intelligence systems.

🔊 K2 Enterprises

Copyright 2024, K2 Enterprises, LLC



Copyright 2024, K2 Enterprises, LLC

- Artificial Neural Networks A computing system inspired by how the human brain processes information. A neural network is made up of a number of simple, highly interconnected processing elements, which processes information by its dynamic state response to external inputs.
- Breach The loss of control, compromise, unauthorized disclosure, unauthorized acquisition, or any similar occurrence where: a person other than an authorized user accesses or potentially accesses personally identifiable information; or an authorized user accesses personally identifiable information for another than authorized purpose.

🔊 K2 Enterprises

# <section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item>



- Data Governance A set of processes that ensures that data assets are formally managed throughout the enterprise. A data governance model establishes authority and management and decision-making parameters related to the data produced or managed by the enterprise
- Data Poisoning Machine learning systems trained on userprovided data are susceptible to data poisoning attacks, whereby malicious users inject false training data with the aim of corrupting the learned model

Signature Karante Kara

### **NIST Definitions**

• **Deep Learning** – An approach to AI that allows computers to learn from experience and understand the world in terms of a hierarchy of concepts, with each concept defined through its relation to simpler concepts. By gathering knowledge from experience, this approach avoids the need for human operators to formally specify all the knowledge that the computer needs. The hierarchy of concepts enables the computer to learn complicated concepts by building them out of simpler ones. If we draw a graph showing how these concepts are built on top of each other, the graph is deep, with many layers.

Signation K2 Enterprises

Copyright 2024, K2 Enterprises, LLC



- **Privacy by Design** Embedding privacy measures and privacy enhancing technologies directly into the design of information technologies and systems
- **Responsible AI** An AI system that aligns development and behavior to goals and values. This includes developing and fielding AI technology in a manner that is consistent with democratic values.
- **Robust AI** An AI system that is resilient in real-world settings, such as an object-recognition application that is robust to significant changes in lighting. The phrase also refers to resilience when it comes to adversarial attacks on AI components.

🔊 K2 Enterprises



#### • Trustworthy AI

- "Characteristics of trustworthy AI systems include: valid and reliable, safe, secure and resilient, accountable and transparent, explainable and interpretable, privacy-enhanced, and fair with harmful bias managed." (NIST definition)
- "AI that can be trusted by humans. Conditions for such trust can refer to (other) ethical principles such as human dignity, respect for human rights, and so on, and/or to social and technical factors that influence whether people will want to use the technology. The use of the term 'trust' with regard to technologies is controversial." (from AI Ethics by Mark Coeckelbergh)
- "Trustworthy AI has three components: (1) it should be lawful, ensuring compliance with all applicable laws and regulations (2) it should be ethical, demonstrating respect for, and ensure adherence to, ethical principles and values and (3) it should be robust, both from a technical and social perspective, since, even with good intentions, AI systems can cause unintentional harm. Trustworthy AI concerns not only the trustworthiness of the AI system itself but also comprises the trustworthiness of all processes and actors that are part of the system's life cycle." (EU draft report on Ethics Guidelines for Trustworthy AI)

Signature Karente Kare

#### **NIST Definitions**

Copyright 2024, K2 Enterprises, LLC

- Generative AI describes algorithms (such as ChatGPT) that can be used to create new content, including audio, code, images, text, simulations, and videos
- Hallucination generated content that is nonsensical or unfaithful to the provided source content; when a bot confidently says something that is not true
  - Intrinsic Hallucination a generated output that contradicts the source content
  - Extrinsic Hallucination a generated output that cannot be verified from the source content (i.e., output can neither be supported nor contradicted by the source)

🔊 K2 Enterprises



- Interpretable Model An interpretable machine learning model obeys a domain-specific set of constraints to allow it (or its predictions, or the data) to be more easily understood by humans. These constraints can differ dramatically depending on the domain
- Language Model A language model is an approximative description that captures patterns and regularities present in natural language and is used for making assumptions on previously unseen language fragments

Signal K2 Enterprises

### **NIST Definitions**

- Machine Learning (ML) A general approach for determining models from data; machine learning is based on algorithms that can learn from data without relying on rules-based programming
- Minimization When using personal information or other high-risk data, the data available is limited as much as possible to those items which are relevant and those things which are absolutely necessary to perform the task
  - As an example, many payroll departments minimize their payroll processes to not require access to social security numbers/SINs - usually by using an employee number instead of an SSN/SIN as the unique identifying number for an employee

Signation K2 Enterprises

Copyright 2024, K2 Enterprises, LLC



- Robust AI An AI system that is resilient in real-world settings, such as an object-recognition application that is robust to significant changes in lighting; the phrase also refers to resilience when it comes to adversarial attacks on AI components
- Supervised Learning A type of machine learning in which the algorithm compares its outputs with the correct outputs during training. In unsupervised learning, the algorithm merely looks for patterns in a set of data
- **Transfer Learning** A technique in machine learning in which an algorithm learns to perform one task, such as recognizing cars, and builds on that knowledge when learning a different but related task, such as recognizing cats

🖎 K2 Enterprises





#### There Is A Need For Multiple Models

- Narrow models
- Small Language Models (SLMs)
- Medium Language Models (MLMs)
- Large Language Models (LLMs)

🔊 K2 Enterprises

## <section-header><section-header><section-header><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item><list-item>











#### One Model Doesn't Fit All

- Narrow models
- Small Language Models (SLMs)
- Medium Language Models (MLMs)
- Large Language Models (LLMs)

🔊 K2 Enterprises





#### AI Leakage Research – November 2023

- "Scalable Extraction of Training Data from (Production) Language Models", published 11/28/2023
- Discusses memorization in large language models two kinds
  - **Discoverable memorization** all model training data which can be recovered with other pieces of the training data set
  - Extractable memorization the subset of discoverable memorization which can be efficiently recovered by an adversary



Milad Nax<sup>+1</sup> Nicholas Carlint<sup>+1</sup> Jonathan Hayase<sup>1,2</sup> Matthew Jagielski<sup>1</sup> A. Feder Cooper<sup>3</sup> Daphne Ippoliuo<sup>1,4</sup> Christopher A. Choquette-Choo<sup>3</sup> Eric Wallace<sup>5</sup> Florian Trans<sup>1,6</sup> Katherine Lee<sup>1,1,2</sup> <sup>1</sup>Google DeepMind <sup>2</sup>Uivers Washington <sup>-5</sup>Cornell <sup>+</sup>CMU <sup>3</sup>U Derkeley <sup>6</sup>ETH Zurich <sup>+</sup>Equal contribution <sup>+</sup>Senior author

#### Abstract

This paper studies *extractable memorization*: training data that an adversary can efficiently extract by querying a machine learning model without prior knowledge of the training dataset. We show an adversary can extract gigabytes of training data from open-source language models like Pythia or GPT-Neo, semi-open models like LLaMA or Falcon, and closed models like ChatGPT. Existing techniques from the literature suffice to attack unaligned models; in order to attack the aligned ChatGPT, we develop a new *divergence* attack that causes the model to diverge from its chatbot-style generations and emit training data at a rate 150× higher than when behaving properly. Our methods show practical attacks can recover far more data than previously thought, and reveal that current alignment techniques do not eliminate memorization.

🔊 K2 Enterprises

Copyright 2024, K2 Enterprises, LLC

#### AI Leakage Research – November 2023



- A new divergence attack was created for ChatGPT which caused it to emit memorized training data more frequently
- "Our methods show <u>that practical</u> <u>attacks can recover more than</u> <u>previously thought</u>, and reveal that <u>current alignment techniques do not</u> <u>eliminate memorization."</u>

Scalable Extraction of Training Data from (Production) Language Model

Milad Nas<sup>\*1</sup> Nicholas Carlini<sup>\*1</sup> Jonahan Hayas<sup>1,2</sup> Mathew Jagletski<sup>1</sup> A. Feder Cooper<sup>3</sup> Daphne Ippolitol.<sup>4</sup> Christopher A. Chaquette-Choo<sup>1</sup> Eric Wallace<sup>6</sup> Florian Tramèr<sup>6</sup> Katherine Lee<sup>+1,3</sup> <sup>1</sup>Google DeepMind <sup>2</sup>University of Washington <sup>7</sup>Cornell <sup>4</sup>CMU <sup>5</sup>UC Berkeley <sup>6</sup>ETH Zurich <sup>\*</sup>Equal contribution <sup>+</sup>Senior author

#### Abstract

This paper studies *extractable memorization*: training data that an adversary can efficiently extract by querying a machine learning model without prior knowledge of the training dataset. We show an adversary can extract gigabytes of training data from open-source language models like Pythia or GPT-Neo, semi-open models like LLaMA or Falcon, and closed models like ChatGPT. Existing techniques from the literature suffice to attack unaligned models; in order to attack the aligned ChatGPT, we develop a new *divergence* attack that causes the model to diverge from its chatbot-style generations and emit training data at a rate 150× higher than when behaving properly. Our methods show practical attacks can recover far more data than previously thought, and reveal that current alignment techniques do not eliminate memorization.

Signature Karente Kare

#### Nine Common Types Of AI Attacks From AI Time Journal



- 1. Adversarial Attacks
- 2. Data Poisoning Attacks
- 3. Model Inversion Attacks
- 4. Membership Inference Attacks
- 5. Evasion Attacks
- 6. Transfer Attacks
- 7. Distributed Denial of Service (DDoS) Attacks
- 8. Data Manipulation Attacks
- 9. Misuse of AI Assistants

🔊 K2 Enterprises

## AI Risks Explored By OpenAI - GPT-4

- Hallucinations
- Harmful content
- Harms of representation, allocation, and quality of service
- Disinformation and influence operations
- Proliferation of conventional and unconventional weapons
- Privacy

- Cybersecurity
- Potential for risky emergent behaviors
- Interactions with other systems
- Economic impacts
- Acceleration
- Overreliance

K2 Enterprises

Copyright 2024, K2 Enterprises, LLC

### NIST's Four Major Types Of Attacks On Al Systems

- Evasion Attacks
  - Occur after an AI system is deployed, attempt to alter an input to change how the system responds to it
  - Examples include adding markings to speed limit signs to make an autonomous vehicle misinterpret them as stop signs or creating confusing lane markings to make the vehicle veer off the road
- Poisoning Attacks
  - Occur in the training phase by introducing corrupted data
  - One example would be slipping numerous instances of inappropriate language into conversation records, so that a chatbot interprets these instances as common enough parlance to use in its own customer interactions



Copyright 2024, K2 Enterprises, LLC

### NIST's Four Major Types Of Attacks On Al Systems



- Privacy Attacks
  - Occur during deployment, are attempts to learn sensitive information about the AI or the data it was trained on in order to misuse it also called "jailbreak attacks"
  - An adversary can ask a chatbot numerous legitimate questions, and then use the answers to reverse engineer the model so as to find its weak spots — or guess at its sources
  - Adding undesired examples to those online sources could **make the AI behave inappropriately**, and making the AI unlearn those specific undesired examples after the fact can be difficult
- Abuse Attacks
  - Involve the insertion of incorrect information into a source, such as a webpage or online document, that an AI then absorbs
  - Unlike the aforementioned poisoning attacks, abuse attacks attempt to give the AI incorrect pieces of information from a legitimate but compromised source to **repurpose the AI system's intended use**

Signal K2 Enterprises

## **Types Of Jailbreak Attacks**

Category	Description	
Attempt to change system rules	This category comprises, but is not limited to, requests to use a new unrestricted system/AI assistant without rules, principles, or limitations, or requests instructing the AI to ignore, forget and disregard its rules, instructions, and previous turns	
Embedding a conversation mockup to confuse the model	This attack uses user-crafted conversational turns embedded in a single user query to instruct the system/AI assistant to disregard rules and limitations	
Role-Play	This attack instructs the system/AI assistant to act as another "system persona" that does not have existing system limitations, or it assigns anthropomorphic human qualities to the system, such as emotions, thoughts, and opinions	
Encoding Attacks	This attack attempts to use encoding, such as a character transformation method, generation styles, ciphers, or other natural language variations, to circumvent the system rules	











## Microsoft's AI Copilots



Name	Products	Monthly Cost	Commercial Data Protection Included?
Copilot for MS 365 (Business/Enterprise)	Microsoft 365 apps (Word, Excel, PowerPoint, Outlook, Teams)	\$30 per user	Yes
Copilot in Windows (Bing Chat)	Windows OS	Free	Not for home users, included with most business/enterprise O365/M365 plans
Copilot Pro for Individuals	Advanced features on top of standard Copilot, plus integration with home Microsoft 365 apps	\$20/user/mo.	Not specified
Copilot for Security	Microsoft's cybersecurity products	Consumption-based fee - \$4/hour	Not specified
Copilot for Finance, Sales, and Service	Financial operations, sales optimizations, service enhancements	\$50/user/mo., \$20/user/mo. if already have MS 365	Not available to other customers, runs on Microsoft cloud in separate instance of ChatGPT, not used by MS to train models by default
Designer for Copilot	Image creation and editing	Not available	No
Copilot GPTs and Azure AI Studio	Custom generative AI assistants and solutions	Not specified	Not specified

#### 🔊 K2 Enterprises



