

# The Anecdotes AI Toolkit

A comprehensive list of risks, controls and policies to securely integrate Generative AI into your organization.



## Index

Contributors	02
Preamble	03
Security Tiers	04

#### The GenAl Risk Register

ANEC-AI-1 - Input of unsanitized PII/PHI	05
ANEC-AI-2 - Inaccurate information usage	06
ANEC-AI-3 - Copyrighted information usage	07
ANEC-AI-4 - Input of sensitive business information	80
ANEC-AI-5 - Vulnerable code usage	09
ANEC-AI-6 - Advanced social engineering attack	10
ANEC-AI-7 - LLM dependency	11
ANEC-AI-8 - Biased output data	12
ANEC-AI-9 - Prompt injection	13
ANEC-AI-10 - Direct LLM attack	14

#### The Anecdotes AI Framework - Controls Listing

Governance	15
Training	16
GenAl Vendors	17
Privacy	18
Customer Obligations	18
Secure Development	19
Technological Protection	20
Dependency Resilience	22
Generative AI Implementation Policy Template	23

Summary		27
---------	--	----





## Author



**Ethan Altmann,** Compliance Product Owner @ Anecdotes

### Contributors

The Anecdotes AI Toolkit was created in collaboration with this exceptional group of industry leading experts who have each lent their unique perspective and experience to this document. Their invaluable insight has enabled the vision of a practical toolkit for Security and GRC professionals to come to life.



Val Dobrushkin Director, Risk and Compliance @ NoName



**Prabhath Karanth** Global Head of Security & Trust @ Navan



Karl Mattson CISO @ NoName



Esther Pinto CISO @ Anecdotes



**Tamir Ronen** Global CISO @ HiBob



Olivia Rose CISO @ Rose CISO Group



Dineshwar Sahni Senior Cybersecurity Leader



Omer Singer Head of Cybersecurity Strategy @ Snowflake



Ronit Shlyfer Information Security & Compliance, Team Lead @ Riskified



Matt Szymanski Director of Security Engineering @ Yext



Tyler Young CISO @ BigID



## Preamble



New GenAI use cases are being discovered, tested and brought into the limelight on a daily basis. Striking the right balance between justified excitement and a healthy dose of skepticism is a challenge faced by countless tech leaders, but especially by Security and Compliance professionals. The core challenge they face lies in empowering their organization to reap the benefits of this new technology, whilst ensuring that the new processes do not exceed the organization's defined risk appetite.

This document aims to give organizations a clear framework through which Generative AI can be securely integrated into existing organizational operations, without impeding security and Compliance processes and obligations. By implementing the framework, security and Compliance leaders take an important step towards ensuring that their organization remains at the forefront of embracing modern technology, whilst refraining from excessive risk exposure and conforming to best practices.

Practically speaking, this toolkit takes a framework based methodology by which organizations can securely and effectively use and implement GenAI tools:

- By understanding the associated risks and their potential-impact on the organization
- By understanding and implementing appropriate mitigative controls to reduce these risks to an acceptable level
- By defining organizational policies that govern usage and outline control implementation expectations

Naturally, there is not a one-size-fits-all approach to security and Compliance, and as such, the risks, controls and policy outlined in this document serve as guidelines that should be modified and adapted to the precise needs of your organization.

## Security Tiers

By virtue of an organization's usage of GenAI, amongst a plethora of other parameters (such as industry, product/service offering, Compliance obligations, sensitivity of data etc.), the level and type of risk exposure will differ.

As such, this framework categorizes its recommendations based on the following three Security Tiers:

#### ORGANIZATION TYPE

• Employees use GenAI tools for general day-to-day tasks.

#### **Security Tier 1**

- Organization uses GenAl, such as a Large Language Model (LLM), to conduct business-critical processes and has dedicated resources to actively training the LLM.
- Organization has deployed GenAI within the production environment of the product/service offering.

#### ORGANIZATION TYPE

• Employees use GenAI tools for general day-to-day tasks.

#### **Security Tier 2**

• Organization uses GenAl, such as a LLM, to conduct business-critical processes and has dedicated resources to actively train the LLM.

#### **ORGANIZATION TYPE**

#### **Security Tier 3**

• Employees use GenAl tools for general day-to-day tasks.



## The GenAI Risk Register

## ANEC-AI-1 Input of unsanitized PII/PHI

GenAl tools are commonly used for performing functions/analysis on large data sets. As such, it must be ensured that inputs are void of data that is subject to regulatory or contractual limitations, as to avoid infringement. Employees may knowingly input such data as the process of masking or anonymizing the data may be excessively strenuous. Alternatively, employees may unwittingly input such data due to a lack of awareness of what constitutes PII/PHI, or simply due to a lack of awareness of the data content. Regardless, formalized processes should be implemented to ensure employees understand their individual responsibilities and the broad implications of PII/PHI misuse. Employee access to PII/PHI should be restricted, based on 'least privilege'. Furthermore, employees accessing PII/PHI for training LLMs provided with an obscured data set (such as data that has been subjected to sanitization mechanisms).





### ANEC-AI-2 Inaccurate information usage

The output generated by GenAI tools is delivered in a manner that is both highly efficient and highly convincing, and therefore employees (ST2-3) or end-users (ST1) are likely to take outputs at face value. However, in doing so, the organization is at risk of not only consuming inaccurate information, but also of perpetuating it. If the organization were to stand accused of spreading misinformation, this could have legal implications, cause significant reputational damage, and result in a loss of customer trust. Therefore, the organization should define a manual QA methodology for reviewing output data (such as a formal legal review, peer review, technical review or managerial review). Where GenAl is used in the production environment, a disclaimer should be presented to end-users regarding 3rd party responsibility for the accuracy of outputted information.

0



### ANEC-AI-3 Copyrighted information usage

Many LLMs draw data from an abundance of sources, some of which have been shown to contain copyrighted/proprietary information. Usage of this information or models trained/fine-tuned on copyrighted information would therefore constitute copyright infringement and potential legal action. As such, formal protocols should be implemented to validate the source of all data within prompt responses, datasets used for initial model training as well as datasets used to fine-tune trained models prior to usage (this may include a formal legal review).





## ANEC-AI-4 Input of sensitive business information

GenAl tools' ability to provide instant, valuable feedback on inputted data such as source code, financial data, architecture diagrams etc. may lead to an unintentional intellectual property leakage, in the event that the sensitive data is then presented to other users. As such, technical controls should be in place to ensure that GenAl usage is conducted in an environment that is segmented, where the tool does not utilize the input prompts for training any other models.



### ANEC-AI-5 Vulnerable code usage

GenAl tools are capable of instantaneously generating code that may take developers significant time to write, and as such there are circumstances under which specific prompt types may provide crucial shortcuts. However, these 'shortcuts' may prove costly in the long run if secure software development best practices are not followed, such as performing SAST, peer review processes, and testing/QA prior to production deployment. Moreover, the aforementioned best practices should be monitored and enforced by leadership, creating a secure development culture.

0





## ANEC-AI-6 Advanced social engineering attack

Prior to GenAl tools, adversarial phishing/smishing attacks most often had tell-tale signs, allowing riskaware recipients to raise a flag and not fall victim. In turn, awareness training provided to employees tends to focus on typos/unconvincing language/out-of-the-ordinary context. However, with the help of GenAl, adversaries can now overcome these shortcomings, and produce high-quality phishing/smishing attacks that are much more difficult to detect. As such, organizations must ensure that employee awareness training curricula shift focus to teaching more technical approaches to recognizing these attack vectors, as well as ensuring that there are sufficient technical safeguards in place to filter out these attacks before they reach employees.

RISK EVENT DESCRIPTION	Adversary uses high-quality GenAl outputs to create advanced phishing campaigns, successfully achieving a foothold in organizational systems.
SECURITY TIER APPLICABILITY	ST1 ST2 ST3
RISK EFFECT	Confidentiality Integrity Availability
RISK SOURCE	Malicious actor.
THREAT	GenAl tool outputs phishing content without expected tell-tale signs.
VULNERABILITY	Lack of employee ability/technical ability to detect malicious emails.
MITIGATING CONTROLS	2.2 Phishing, 7.3 Email filtering.



```

### ANEC-AI-7 LLM dependency

GenAl LLMs possess data processing capabilities that are largely unparalleled, which may lead to excessive dependence, especially if they come in place of human expertise that is capable of performing the function in the event of LLM outage/downtime. If this business process directly ties to a contractual obligation (such as an SLA), this could pose significant risk to the organization. As such, effective disaster recovery and business continuity processes should be developed, and periodically tested.

0





### ANEC-AI-8 Biased output data

GenAl LLMs output data is based on training data that may stem from an unreliable or unverified source. As such, output data may possess the same biases or discriminations as the source data. This can include racial or gender discrimination, or any other protected characteristics, and can therefore, in turn, result in legal action against the organization, as well as reputational damage. This can be partly mitigated by the implementation of technical controls that ensure the diversity and integrity of training data, as well as an end-user facing disclaimer.





### ANEC-AI-9 Prompt injection

GenAI LLMs can be configured to restrict output to specific parameters, using guardrails. However, adversaries are constantly seeking to design prompts that are capable of circumventing these guardrails, allowing them to utilize LLMs for potentially malicious activity. Many GenAI tools are currently releasing frequent updates to attempt to improve prompt injection resistance, and as such it should be ensured the most up-to-date version of the tool is used. Furthermore, guardrails should be put in place and continuously improved.

0





### ANEC-AI-10 Direct LLM attack

Where GenAI LLMs are used to provide production environment functionality to end-users, adversaries may seek to conduct attacks that result in events such as a data breach, poisoned output data, denial of service and more. This may be achieved by directly attacking the LLMs back-end and gaining a foothold into the LLM management layer and data sets (e.g. through the source code or by accessing the GenAI tool itself). As such, maintaining good credential hygiene as well as ensuring that GenAI vendors adhere to industry-standard security practices should be ensured.





## The Anecdotes AI Framework

The following controls listing aims to provide proactive objectives for effective organizational measures in reducing Al risk exposure.

#### 1. Governance

| CONTROL NAME<br>1.1<br>Al Policy              | CONTROL DESCRIPTION<br>The organization should develop<br>an AI-specific policy, outlining<br>usage limitations and<br>expectations, and require all<br>employees to read and formally<br>acknowledge the policy prior to<br>GenAI tool usage and<br>development. | APPLICABILITY<br>ST1 ST2 ST3<br>EQUIVALENT CONTROLS<br>iso/iec 27001 - isms 5.2<br>SA STAR A&A - 1<br>ST<br>NIST 800-53 REV 5 - AU-1       |
|-----------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------|
| CONTROL NAME<br>1.2<br>Policy<br>Augmentation | CONTROL DESCRIPTION<br>The organization should augment<br>existing policies (such as the<br>Acceptable Use policy, the SDLC<br>policy etc.) to explicitly reference<br>boundaries for AI usage and<br>provide effective implementation<br>guidance.               | APPLICABILITY<br>ST1 ST2 ST3<br>EQUIVALENT CONTROLS<br>ISO/IEC 27001 - ISMS 5.2<br>SA STAR A&A - 1<br>STA STAR A&A - 1<br>STA STAR A&A - 1 |

AAIF

 $\bigcirc$ 

#### 2. Training



#### 3. GenAl Vendors

| CONTROL NAME       | CONTROL DESCRIPTION                                                                                                                                                                                              | APPLICABILITY                                                                                        |
|--------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------|
| 3.1<br>Segregation | The organization should ensure<br>that where possible, a version of<br>the GenAl tool is used that does<br>not use inputted data to train or<br>improve non-organization owned<br>models (instance segregation). | ST1 ST2 ST3                                                                                          |
| CONTROL NAME       | CONTROL DESCRIPTION<br>The organization should ensure                                                                                                                                                            | APPLICABILITY<br>ST1 ST2 ST3                                                                         |
| Agreement          | that the GenAl vendor formally<br>bears responsibility for the<br>protection of inputted data, as well<br>as liability for potential data<br>inaccuracies.                                                       | EQUIVALENT CONTROLS<br>(50) ISO/IEC 27001 - A.15.1.3<br>(05) CIS V8 - 15.4<br>(05) CSA STAR STA - 09 |
|                    |                                                                                                                                                                                                                  |                                                                                                      |
| CONTROL NAME       | CONTROL DESCRIPTION                                                                                                                                                                                              | APPLICABILITY                                                                                        |
| 3.3<br>Assessment  | The organization should ensurethat the GenAl vendor hasundergone a security and privacy                                                                                                                          | ST1 ST2                                                                                              |
|                    | the organizational Vendor                                                                                                                                                                                        | ISO/IEC 27001 - A.15.2.1                                                                             |
|                    | Management policy) to ensure that the vendor adheres to industry-                                                                                                                                                | CIS V8 - 15.5                                                                                        |
|                    | standard security practices.                                                                                                                                                                                     | CSA STAR STA - 13                                                                                    |
|                    |                                                                                                                                                                                                                  | NIST 800-53 REV 5 - SR-3                                                                             |
|                    |                                                                                                                                                                                                                  |                                                                                                      |

📙 anecdotes

--

#### 4. Privacy



#### 5. Customer Obligations

| CO | NΤ | RO | LN | IAI | ME |
|----|----|----|----|-----|----|
|    |    |    |    |     |    |

#### 5.1 Disclaimer

#### CONTROL DESCRIPTION

The organization should ensure that customers are presented with a disclaimer prior to usage of GenAI elements within the product offering, which recommends input data practices and warns of potential output data inaccuracies, and defers responsibility for outputted data to the GenAI vendor.

#### APPLICABILITY



📙 anecdotes

The Anecdotes AI Toolkit The Anecdotes AI Framework

-----

| CONTROL NAME<br>5.2<br>Opt-out      | CONTROL DESCRIPTION<br>The organization should ensure that<br>customers are able to opt-out of<br>GenAI elements within their product<br>offering, and are informed of GenAI<br>usage/functionality when it is<br>enabled by default. | APPLICABILITY<br>ST1 |
|-------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------|
| CONTROL NAME<br>5.3<br>Consent      | CONTROL DESCRIPTION<br>The organization should ensure that<br>customers formally consent to GenAl<br>usage within the product offering.                                                                                               | APPLICABILITY<br>ST1 |
| CONTROL NAME<br>5.4<br>Terms of Use | CONTROL DESCRIPTION<br>The organization should update the<br>Terms of Use to reflect GenAI usage<br>within the product offering and make<br>existing customers aware of the<br>changes retroactively.                                 | APPLICABILITY<br>ST1 |

#### 6. Secure Development

#### CONTROL NAME

Training

6.1

#### CONTROL DESCRIPTION

The organization should ensure that all relevant employees undergo training in AI secure development practices (as per the augmented SDLC policy) prior to using GenAI tools for development purposes.

#### APPLICABILITY





| CONTROL NAME              | CONTROL DESCRIPTION                                                                                                         |
|---------------------------|-----------------------------------------------------------------------------------------------------------------------------|
| 6.2<br>Threat<br>Modeling | The organization should conduct<br>threat modeling as part of the<br>development process of GenAl<br>product functionality. |

#### APPLICABILITY



#### EQUIVALENT CONTROLS

CIS V8 - 16.14

NIST 800-53 REV 5 - SA-11(2) אוכד

### 7. Technological Protection

| CONTROL NAME<br>7.1<br>Input<br>Validation  | CONTROL DESCRIPTION<br>The organization should ensure<br>that acceptable input parameters<br>are defined and tested and are<br>enforced by the GenAl tool.                                       | APPLICABILITY<br>ST1 ST2<br>EQUIVALENT CONTROLS<br>INIST 800-53 REV 5 - SI-10                     |
|---------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------|
| CONTROL NAME<br>7.2<br>Output<br>Validation | CONTROL DESCRIPTION<br>The organization should<br>implement an output validation<br>process for LLMs in order to ensure<br>that they continue to meet defined<br>accuracy and non-bias criteria. | APPLICABILITY<br>ST1 ST2<br>EQUIVALENT CONTROLS<br>CSA STAR CCC - 02<br>NIST 800-53 REV 5 - SI-15 |



The Anecdotes AI Toolkit The Anecdotes AI Framework

| CONTROL NAME<br>7.3<br>Email<br>Filtering             | CONTROL DESCRIPTION<br>The organization should<br>implement email filtering<br>technologies (such as DMARC) to<br>prevent potentially malicious<br>emails from reaching employee<br>inboxes.                   | APPLICABILITY<br>ST1 ST2 ST3<br>EQUIVALENT CONTROLS<br>CIS V8 - 9.5, 9.7                                  |
|-------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------|
| CONTROL NAME<br>7.4<br>Sensitive<br>Data<br>Discovery | CONTROL DESCRIPTION<br>The organization should<br>implement a sensitive data<br>discovery mechanism to ensure<br>data repository sanitization.                                                                 | APPLICABILITY<br>ST1 ST2<br>EQUIVALENT CONTROLS<br>CSA STAR DS - 03<br>NIST 800-53 REV 5 - AC-4(25)       |
| CONTROL NAME<br>7.5<br>Zero<br>Trust<br>Architecture  | CONTROL DESCRIPTION<br>The organization should<br>implement a zero trust architecture<br>to prevent unauthorized access to<br>sensitive organizational assets<br>(such as DBs containing PII/PHI, or<br>LLMs). | APPLICABILITY<br>ST1 ST2<br>EQUIVALENT CONTROLS<br>CSA STAR IAM - 05<br>NIST 800-53 REV 5 - AC-6          |
| CONTROL NAME<br>7.6<br>Updates                        | <b>CONTROL DESCRIPTION</b><br>The organization should ensure<br>that the GenAl tool used is the<br>most up-to-date and secure<br>version.                                                                      | APPLICABILITY<br>ST1 ST2 ST3<br>EQUIVALENT CONTROLS<br>CIS V8 - 12.1<br>ST<br>NIST 800-53 REV 5 - SI-2(4) |



#### CONTROL NAME

#### 7.7 Credentials

#### CONTROL DESCRIPTION

The organization should enforce user credentials (such as passwords and access keys) that are periodically rotated and meet industry-standard strength. Multifactor authentication should be enforced where possible.



#### **EQUIVALENT CONTROLS**



#### 8. Dependency Resilience

| CONTROL NAME                                      | CONTROL DESCRIPTION                                                                                                                                               |                                                                                                    |
|---------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------|
| 8.1<br>Impact                                     | The organization should perform<br>business impact analysis (BIA) for<br>all GenAI tool usage.                                                                    | ST1   ST2     EQUIVALENT CONTROLS     ISO/IEC 27001 - 17.1.1, 17.1.2, 17.1.3     CSA STAR BCR - 02 |
| CONTROL NAME<br>8.2<br>Recovery<br>and Continuity | CONTROL DESCRIPTION<br>The organization should develop<br>clearly defined disaster recovery<br>and business continuity plans for<br>any GenAl tool usage that has | APPLICABILITY<br>ST1 ST2<br>EQUIVALENT CONTROLS                                                    |



been deemed to be business critical. These plans should be tested periodically.

## Generative AI Implementation Policy Template

## Purpose

[Org. name] is committed to preserving organizational data confidentiality, integrity, and availability and in turn adhering to contractual and regulatory obligations. As such, [Org. name] defines and outlines this organizational Generative Artificial Intelligence (henceforth referred to as "GenAI") policy, in order to uphold these obligations by establishing a governance baseline, upon which technological controls can be built.



This policy applies to [Org. name] employees in all offices, whether working in the office or remotely, and whether employed in a full or part-time capacity. This policy applies to GenAl usage across two different areas of application:

 Employee usage of GenAl tools for day-to-day activities (for example, inputting prompts to a service such as ChatGPT).
[Employee usage of GenAl tools, specifically Large Language Models (henceforth referred to as "LLMs") for business-critical processes.]

[Certain requirements outlined within this policy that are relevant to the organization as a whole and additional GenAI requirements not outlined within this policy that are more department-specific are further detailed in [Org. name] policies, which have been augmented in order to accommodate GenAI tool usage. These policies are:

- [Acceptable Use policy]
- [Secure Development Lifecycle policy]
- [Security Awareness and Training policy]
- [Vendor Management policy]
- [Disaster Recovery and Business Continuity policy]

For organization-wide requirements, this policy provides a brief overview and then makes reference to the topic-specific policy.]



•

The Anecdotes Al Toolkit Generative Al Implementation Policy Template

## Policy life-cycle

This policy is made readily available to all [Org. name] employees via the [internal portal/shared drive/HR tool] and is actively communicated to all employees upon first release, as part of the onboarding process and annually thereafter.

This policy is made available to any interested parties upon request, as deemed appropriate by the [CISO] [and the DPO].

This policy is reviewed and formally approved by the [CISO] [and the DPO] on an annual basis or following any changes.

## **Background**

[Org. name] acknowledges that GenAl tools represent significant opportunities for increased efficiency and productivity, and as such does not wish to prevent or discourage employees from their usage. However, [Org. name] also acknowledges that GenAl tools present new, significant risks to the organization, that without the implementation of effective mitigating controls, will exceed the organizational risk appetite.

As such, [Org. name] has designed a control set that aims to reduce the risks associated with GenAl usage to an acceptable level. These control areas include:

- Defining acceptable use of GenAl tools
- [Implementing additional awareness training content]
- [GenAl vendor security and privacy requirements]
- [Customer facing obligations]
- [Data input/output validation]
- [Dependency resilience]



# Acceptable

Employees using GenAI tools shall conform to the behavioral expectations laid out in the organizational [Code of Conduct], in particular pertaining to refraining from the input of profanity or of discriminatory content or views.

Furthermore, employees shall refrain from:

- Creating prompts that contain organizational intellectual property (such as [source code, financial information, trade secrets etc.]).
- Creating prompts that contain any data that is subject to regulatory or compliance limitations, such as personal identifiable information (PII) or protected health information (PHI).
- Using GenAl tool outputs without an effective accuracy validation process.
- Publishing GenAl tool outputs without a [Legal dept.] approved disclaimer.

In the event that exceptions to the aforementioned are required in order to support a business process, an exception request shall be made in writing and formally approved by the [CISO] [and the DPO].

The [Org.name] [Acceptable use policy] has been appended to include the aforementioned requirements.

## Employee training

- Employees using GenAI tools shall undergo additional information security awareness training that specifically includes:
- Acceptable use of GenAl tools as outlined above.
- Recognizing advanced social engineering attacks.

Engineering teams using GenAl tools shall undergo additional training that specifically includes:

- Risks associated with usage of code outputted by GenAl tools.
- Static Application Security Testing (SAST) process requirements.
- Peer review requirements.

The [Org.name] [Security awareness and training policy] and [Secure development lifecycle policy] have been appended to include and further detail the aforementioned requirements.



----



The [Org.name] [Disaster recovery and Business Continuity plans] and [Professional development program] have been appended to include and further detail the aforementioned requirements.



. - -

## Summary

As we conclude the creation of this document, the continued innovation in Generative AI remains at the front and center of not only the Security and Compliance ecosystem, but of global news as well. And as this is unlikely to change in the near future, new use cases will surely raise questions around new, unintended risks and threats.

It is imperative that organizations stay abreast of new developments, and implement effective mitigation strategies. Anecdotes remains committed to providing the community with practical solutions and tools and will continue to evolve this framework to address new challenges (so stay tuned for the next iteration!).



Anecdotes enables GRC teams to strengthen and scale their GRC program. By giving teams actionable GRC data, which is instantly mapped to any use case, coupled with Al-enhanced analysis tools and configurable automations, Anecdotes empowers GRC teams to quickly identify gaps—and attest to their organization's state of compliance with confidence. Learn more by visiting www.anecdotes.ai