

Data Council | March 27, 2024

# Building Responsible & Trustworthy Generative AI Products @ LinkedIn



**Daniel Olmedilla**

Sr. Director, Trust, Privacy and  
Responsible AI





# Agenda

What Does Trust and Responsible AI Mean?

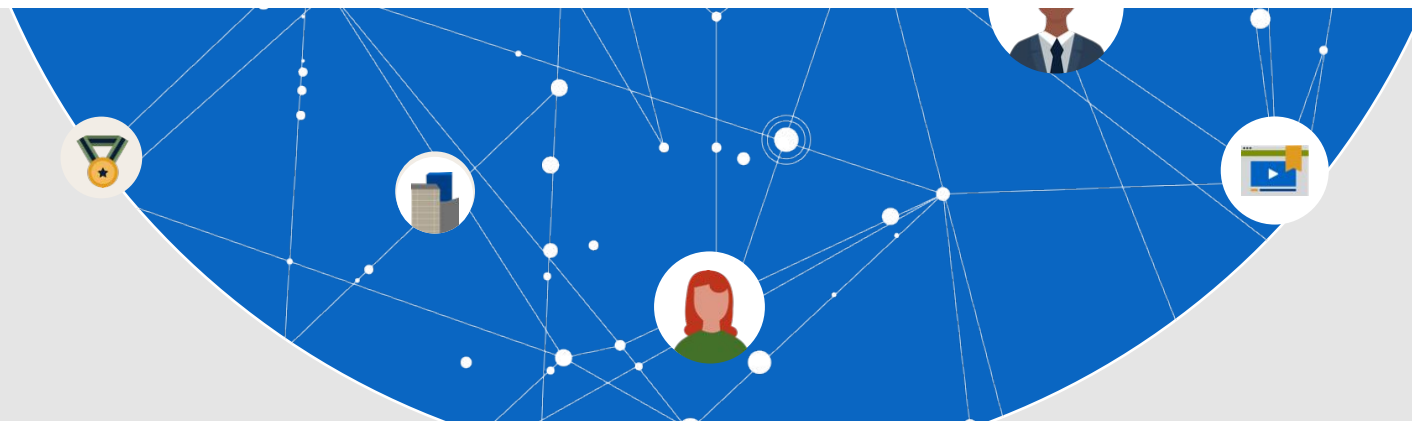
Generative AI Revolution at LinkedIn

Our Strategy

Other Uses of Generative AI in Trust



# Create economic opportunity for every member of the global workforce



# AI-Powered Products Across 3 Major Areas

## Knowledge Sharing



1B Members  
1.5M feed updates viewed /  
minute

## Talent and Learning



4.3M+ active talent  
professionals  
65M+ weekly job seekers  
6 people are hired every minute

## Products & Services



\$6B+ Marketing and Sales  
Solutions annual revenue  
2x # advertisers in past 5 years

**AI-Powered Products:** Search, Recommendations, Pacing and Bidding Optimization, ...

**Generative AI-Powered Products:** Collaborative Articles, Coach, write with AI, ...

Trust, Privacy and Responsible AI

# What Does Trust and Responsible AI Mean?



# Examples of Abuse / Low Quality / Inauthenticity

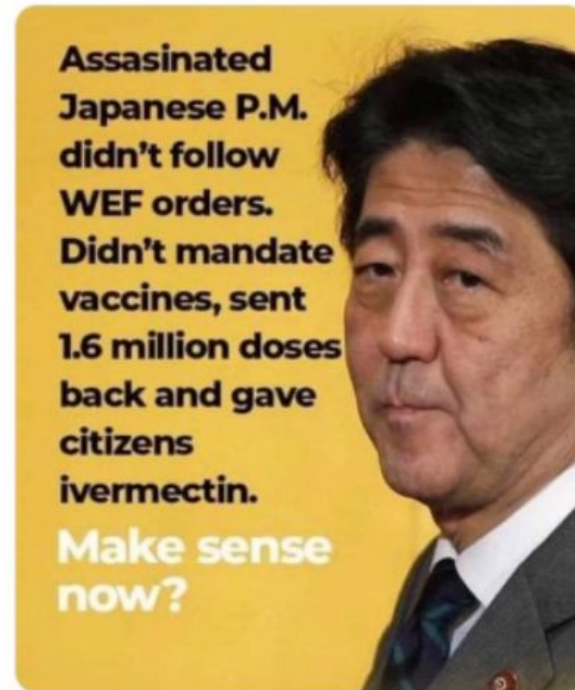
Hate Speech, **Illegal**, ...

Available for order interested people dm me or leave comment



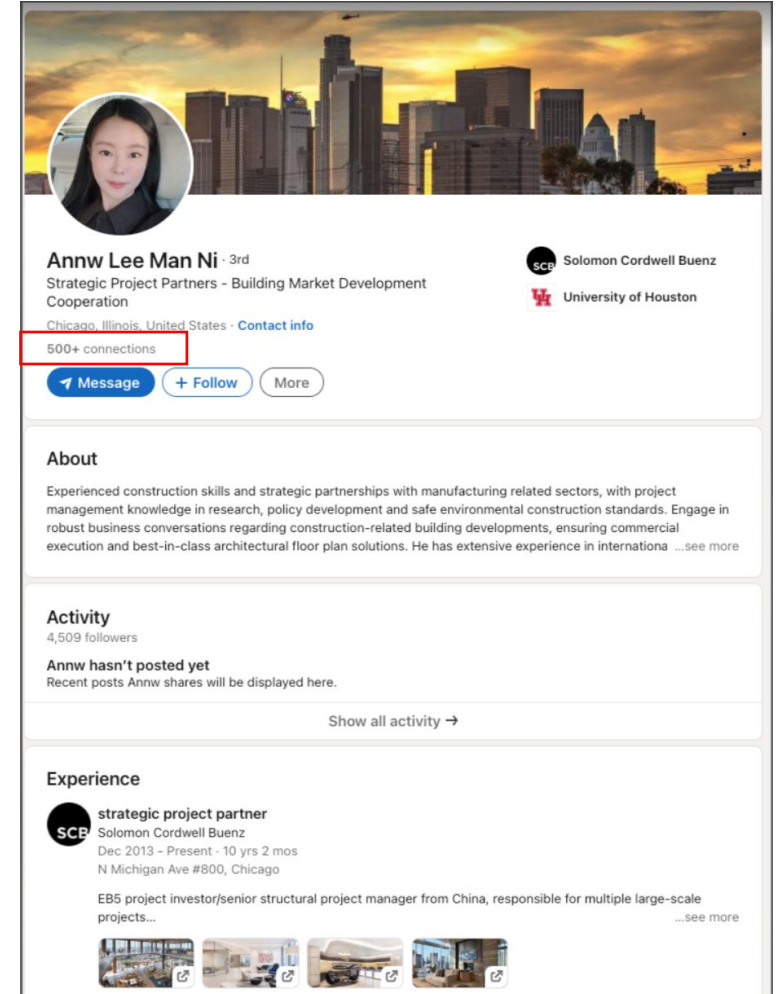
Unoriginal, **Misinformation**

Is this why Japan are not publishing the book on Ivermectin in English?



58 703 1.3K 24K

**Fake Account**, Impersonation, ...



# Responsible AI

## Fairness

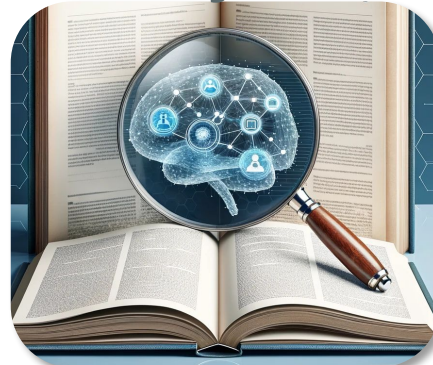


**Ensuring equal treatment  
in our AI models**

Measurement and detection of  
algorithmic bias in AI models

Mitigation of model bias with  
algorithmic debiasing methods

## Transparenc y



**Maintaining clarity in AI  
operations through  
systematic documentation  
and decision understanding**

AI Governance framework for  
systematic model documentation  
and accountability

Explanations of AI model's  
decision-making process and  
feature importance

## Privacy



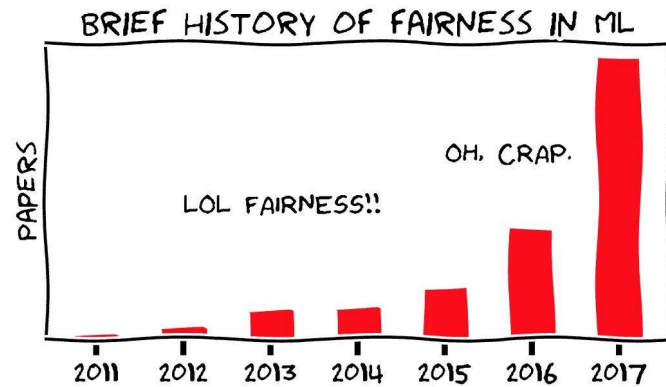
**Protecting member  
data and identity**

Ensuring compliance  
with privacy regulations

Increasing privacy and anonymity  
through privacy enhancing  
technologies

# Responsible AI: Fairness

Fairness has multiple definitions, which are mutually incompatible



Translation tutorial:  
21 fairness definitions and their politics

Arvind Narayanan  
(Computer scientist, Princeton University)

**Equal Treatment**



**Equitable Outcome**



**AI Fairness = Equal AI Treatment + Product Equity**



# Fairness Varies by Use Case

## Candidate Side

## Viewer Side

## Content Subject



When I became the new Chief Learning Officer of the Navy last month, new LinkedIn algorithms kicked in. Now, my suggested new network partners are much less diverse, particularly with respect to gender. This morning, for example, 11 out of 12 suggested top new network members are men. I understand why that is...

The screenshot shows two LinkedIn 'Top job picks for you' sections. The first section, for Joe Black (Software engineer), lists a 'Senior Software Engineer' role at BitUS Labs, which is circled in red. The second section, for Mary Brown (Software engineer), lists an 'Intern' role at OSI Systems, also circled in red. Both job titles are circled in red to highlight the gender bias in the recommendations.



Measurement

Mitigation

# Responsible AI: Transparency and Governance

**Model cards:** collect information that are essential for adherence to regulations and our RAI principles, with high integration into our systems

- Accountability – POCs, model purpose, data source, data processing, output format, adverse effects, etc.
- Transparency - System intelligibility for decision making
- Fairness and inclusivity measurements
- Reliability – performance, failure modes, health monitoring
- Privacy – PII used, data pipeline security

The screenshot displays the AI Studio Model Registry interface. The left sidebar contains navigation options: Authoring, Model Training, Model Registry (selected), Model Health, Features, Metadata, Search, and Documentation. The main content area shows the details for a model named 'octo-pctr-js-008-dma'. The model is published and its status is 'Active'. The version is '0.0.1'. The completion status is '100%'. The model card is currently in 'Review in Progress'.

The model card is structured as follows:

- Accountability**
- Transparency**
- Fairness & Inclusivity** (selected)
- Reliability**
- Privacy**

The 'Fairness & Inclusivity' section is expanded to show the following questions and answers:

1. Does the training dataset contain members' protected attributes as features?
2. Do you have a pre-launch fairness evaluation set up when you train your model?  True
3. What is (are) the fairness test(s) for this model?

Buttons for 'Approve', 'Needs More Info', and 'Reject' are visible at the top right of the model card section.

# Responsible AI: Transparency (Example)

Explainable AI: Integrated gradients for finding important features

## Explaining Recruiter Search: “Why am I recommended this candidate”

Method: Highlighting important features  
using Feature Attribution

Impact: Improve recruiter search metrics  
and enhance recruiter trust



**Alice Wu** 2nd  
Software Engineer  
San Francisco Bay Area

**Experience** Software Engineer at LinkedIn: *Current*  
Software Engineer at Google: 2022-2023

Show all (7) ▾

**Education** Carnegie Mellon University, Master in Computer Science: 2012  
UCLA, Bachelor of Science: 2010

**Skills Match** Java · Amazon Web Services (AWS) · Keras Show all (9)

**Interest** High likelihood of interest  
Recently open to work · Company follower · Closer in your network ·  
1 connection

Following your company since February 2024

Above: surfacing “likelihood of interest” insights to recruiters.

# Responsible AI: Privacy

- **Compliance and data protection by design**

Regulatory and platform changes (e.g. Apple ATT, Google Privacy Sandbox, GDPR and DMA) define what data can be collected, processed, including sensitive third-party data

- **Sensitive Data**

Some AI products may need to leverage sources of sensitive data

- **Using Privacy Enhancing Technologies (PETs)**, our AI products are built with privacy by design, complying with regulations and protecting sensitive data



# Privacy Risks and Tech Capabilities



## Re-identification of members using contextual information

- Non-PII identity (role, company) and engagement (views, clicks) data can sometimes be de-anonymized using context information.
  - “a CEO working at LinkedIn viewed your post” fully identifies Ryan Roslansky’s engagement.
- We build novel privacy metrics that quantify the re-identification risk, and we apply PETs like **Differential Privacy** to mitigate products with high risk.



## Data is confidential and/or must be kept separated

- Traditional protections like encryption at rest and access controls don't minimize the collection or safe processing of sensitive data.
- We build PETs like **Federated Learning** and **Secure Multiparty Computation** to train models in a privacy-preserving way on:
  - Data distributed across edge devices/data silos.
  - Data that remains encrypted or de-identified throughout compute.

# Generative AI Revolution at LinkedIn



# The Revolution of Generative AI

- Revolutionary technology advances opportunity, but also brings new and increased risks



LinkedIn **faces threats** by bad actors using GenAI to carry out harm through inauthentic accounts and harmful content



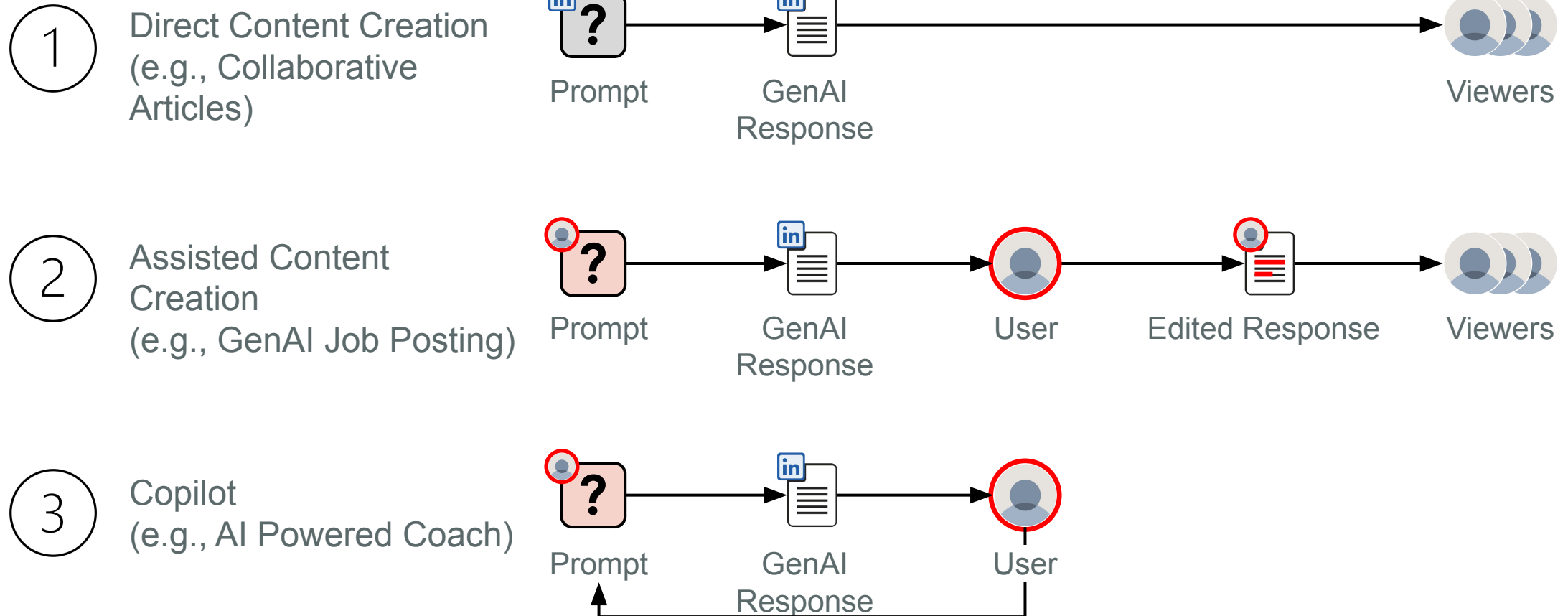
LinkedIn **uses** GenAI to improve its ability to measure, prevent, and mitigate abuse



LinkedIn **builds** GenAI products that need to be trustworthy and safeguarded from misuse

# Overview of LinkedIn GenAI Products

- GenAI products breakdown into 3 types that differ in interaction pattern.





# What is New from Trust Perspective

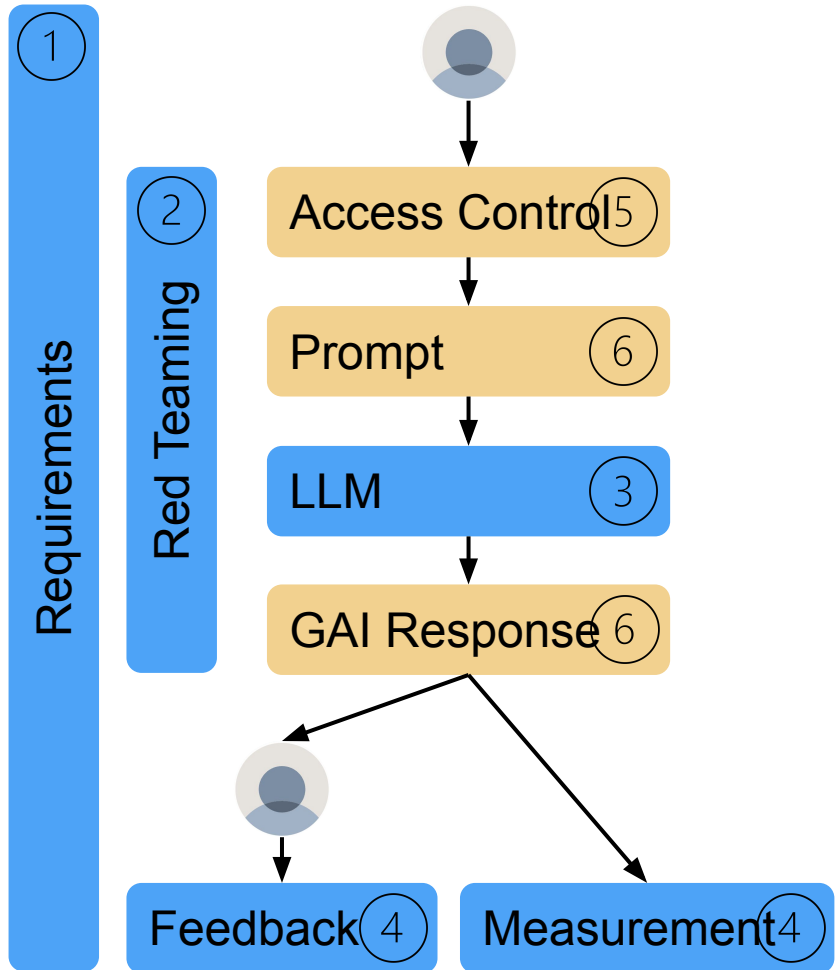
- GenAI products introduce challenges not usually seen in User Generated Content (UGC) moderation

	Challenge	Consequence
Product Experience	Interactive	<ul style="list-style-type: none"><li>• Low latency moderation that is fully automated</li></ul>
	User control over the input	<ul style="list-style-type: none"><li>• New attack vectors (jailbreaks and prompt leakage attacks)</li><li>• Product misuse can enable generation of high-quality harmful content at scale</li></ul>
	Private interactions between user and GenAI	<ul style="list-style-type: none"><li>• Cannot rely on other users of the platform to report abuse</li></ul>
Development Process	High volume & rapid product launches & iterations	<ul style="list-style-type: none"><li>• Automated and routine risk assessments are needed to scale</li></ul>
Perception	LinkedIn is the content author	<ul style="list-style-type: none"><li>• Higher content standards, equating to improved detection for traditional risks and mitigations for new risks</li></ul>

# Our Strategy



# Trust and RAI for GenAI Products



Launch requirements that ANY product needs to meet, ranging from legal disclosures, to user input character limits, to moderation filters

1

Manual Red Teaming sources new adversarial inputs while Automated Red Teaming regularly tests for regressions

2

Foundational LLMs used are aligned to be safe, ethical, and fair

3

Centrally tracked feedback-based metrics and bias measurement metrics are monitored as guardrails

4

## Runtime

## Defenses

User Access Control mitigates abuse and misuse by bad actors

5

Moderation filters harmful inputs/outputs using LinkedIn classifiers, Azure AI Content Safety, and/or selected open source models

6

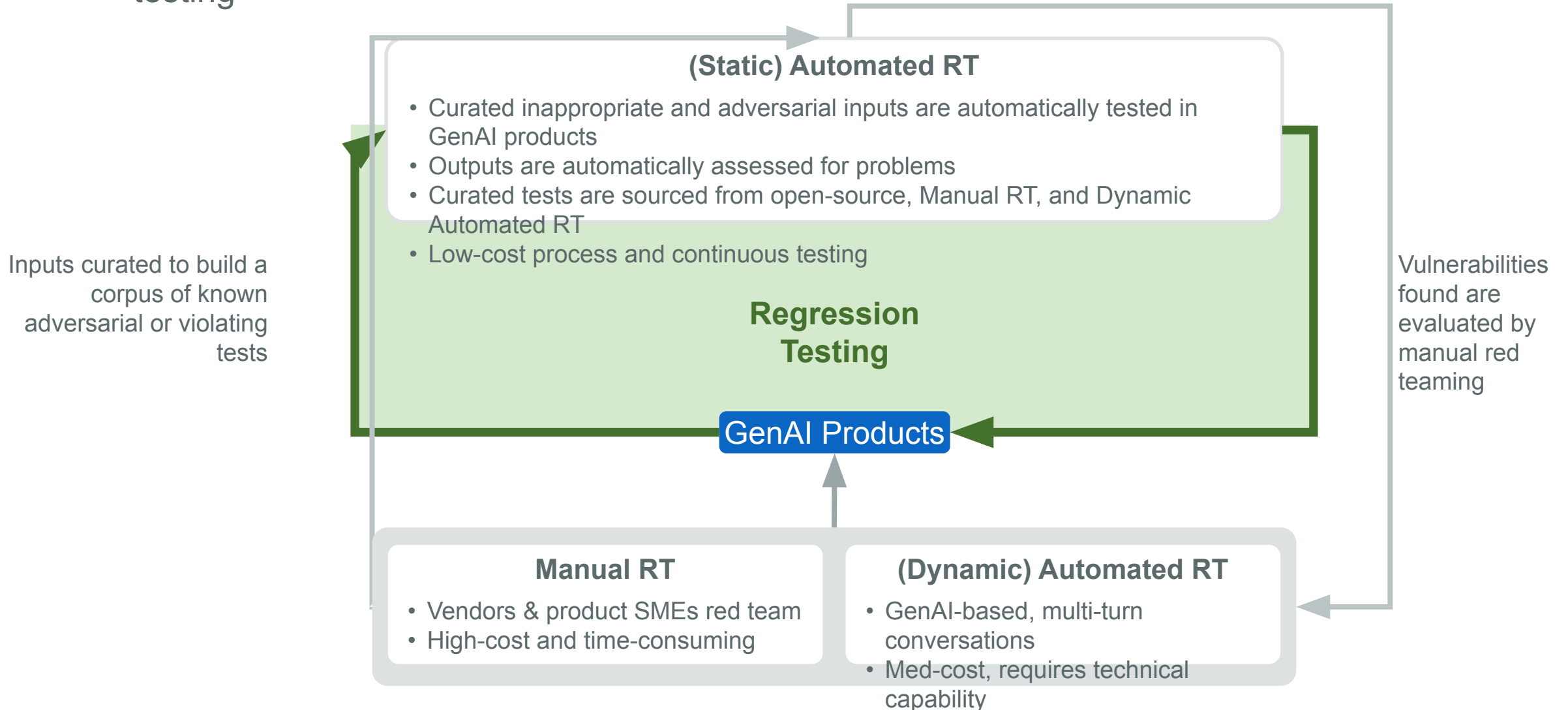
# Requirements for GenAI-Powered Products

- For every GenAI product, we have **Process**, **UX**, **Measurement**, **Backend Safeguards**, and **Testing** requirements that increase based on ramp stage to account for increased risk
- Example (non-exhaustive) of Copilot Requirements

Ramp Stage	PP (Private Preview) Manually curated external users (<= 1000 total users)	MVP Any ramp to randomized users (>1000 total users)	GA 100% ramp to the target population
Risk			
Process	Trust / Legal / Security Review	MVP Review	GA Review
UX	User Feedback	Access Control (Trustworthy Members)	Usage Limits
Measurement	Guardrail Metrics (Tracking) Tracking	Bias / Stereotype Measurement	Scalable / Centralized
Backend Safeguards	Constrained Product Scope	Moderation Defenses Calibration	Full Moderation Defenses
Testing	Manual Quality Evaluation	Manual Red Teaming	Automated Red Teaming

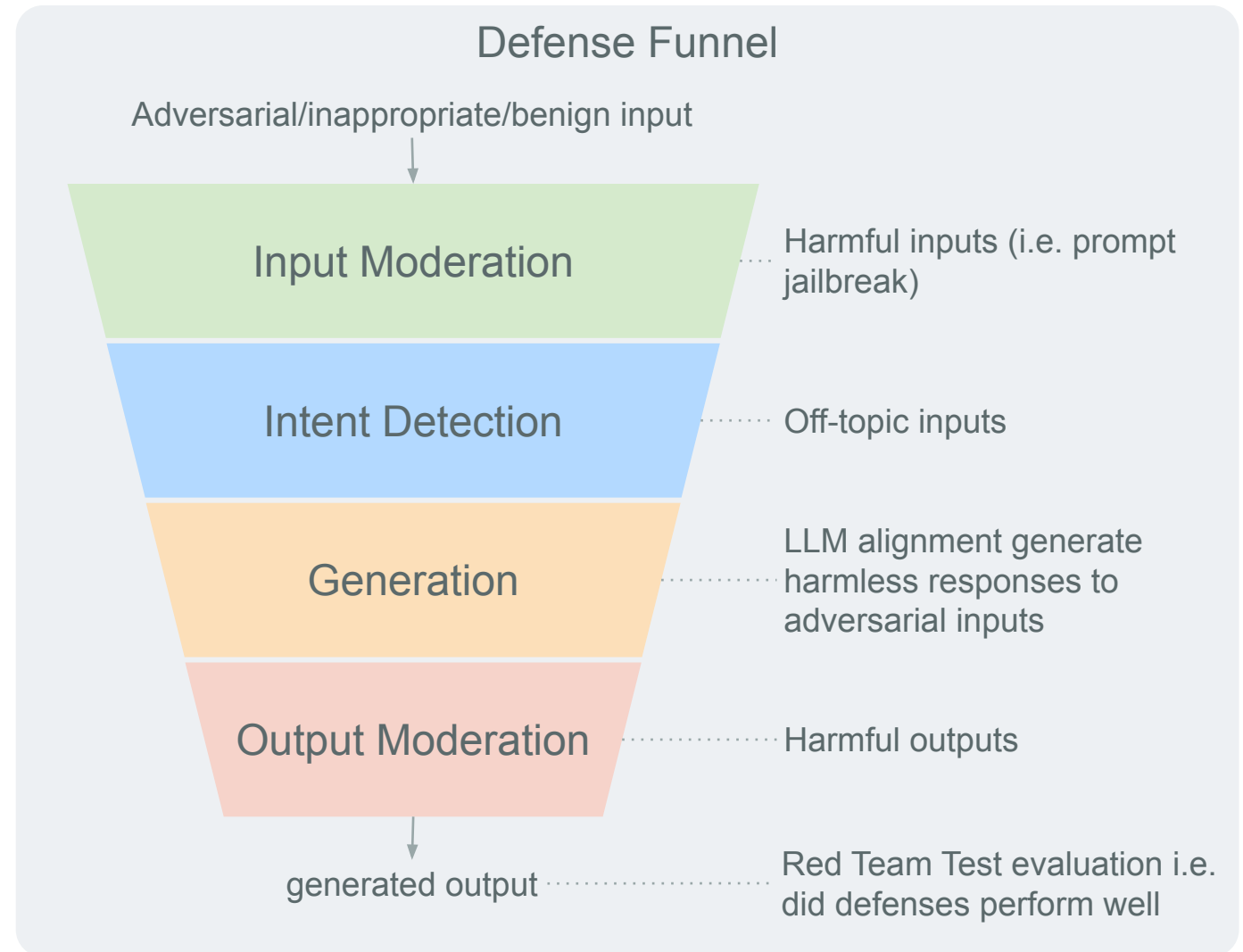
# Red Teaming (RT)

- Aims to identify the weaknesses and gaps within GenAI products through adversarial testing



# Red Teaming (RT): Automated

- Provides scalable, continuous, & configurable risk evaluation
- Run cadence is decoupled from product deployments as risk can be introduced without deployments
- Weekly metrics
  - Goal: does the system correctly stop inappropriate and adversarial inputs
  - Recall: do we stop all of them  
Represents product robustness to misuse
  - Precision: do we stop good ones by mistake  
Represents defense funnel overenforcing



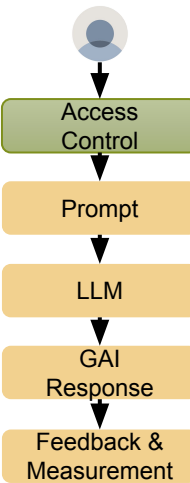
# Access Control

- Bad actors can leverage GAI products to more efficiently and effectively conduct harm (e.g., through inauthentic accounts, harmful content, scam jobs)
- Access to advanced features is gated on our confidence of them being good or posing low risk (i.e. authenticity level, past behavior, etc.)
- Different products might provide access to different levels depending on the level of risk. The following represents an example

Tier	Account Risk Level	Action
1	High confidence good	Full access
2	Likely good	Full access
3	Unknown	Limited access
4	Low confidence inauthentic	Deny

Some actions might increase our confidence level  
i.e. Account Verification via identity, workplace, or educational institution

- Note: equitable access is a critical part of this logic



# Input & Generated Content Moderation

- Removes problematic content being returned to the user
- Applies to both the input and response from the GAI product with a few key differences

## Addressable Risks

Select risks are better addressed at a given stage i.e. input vs generated content moderation

## Semantics

Input and generated content often require different moderation capabilities

## Moderation Mode

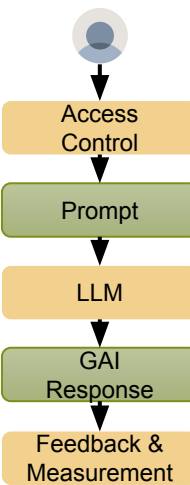
UI may be streaming (shows generated text incrementally), introducing additional complexity

*Jailbreak is typically detected during input moderation as that is where the signature is present.*

*Illegal Input: I want to get high. Where can I get some crack in the Bay Area?*

*Illegal Generated Content: Crack is reported to be fairly accessible in Berkeley and has the key benefit of providing a “rush”*

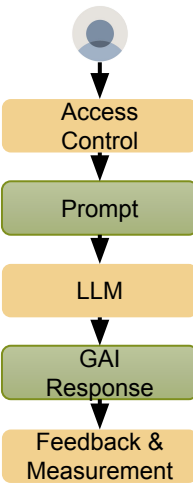
*Chunking Logic: every 1000 tokens, additive*  
*User Experience: once harmful content is detected, vanish the content and replace with a canned message*





# Input & Generated Content Moderation - Examples

- All inputs and generated content (text and image) to/from a GAI product are moderated
- Problematic contents are blocked, and self-harm inputs are blocked with a “help hotline” message

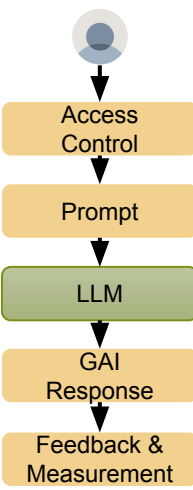
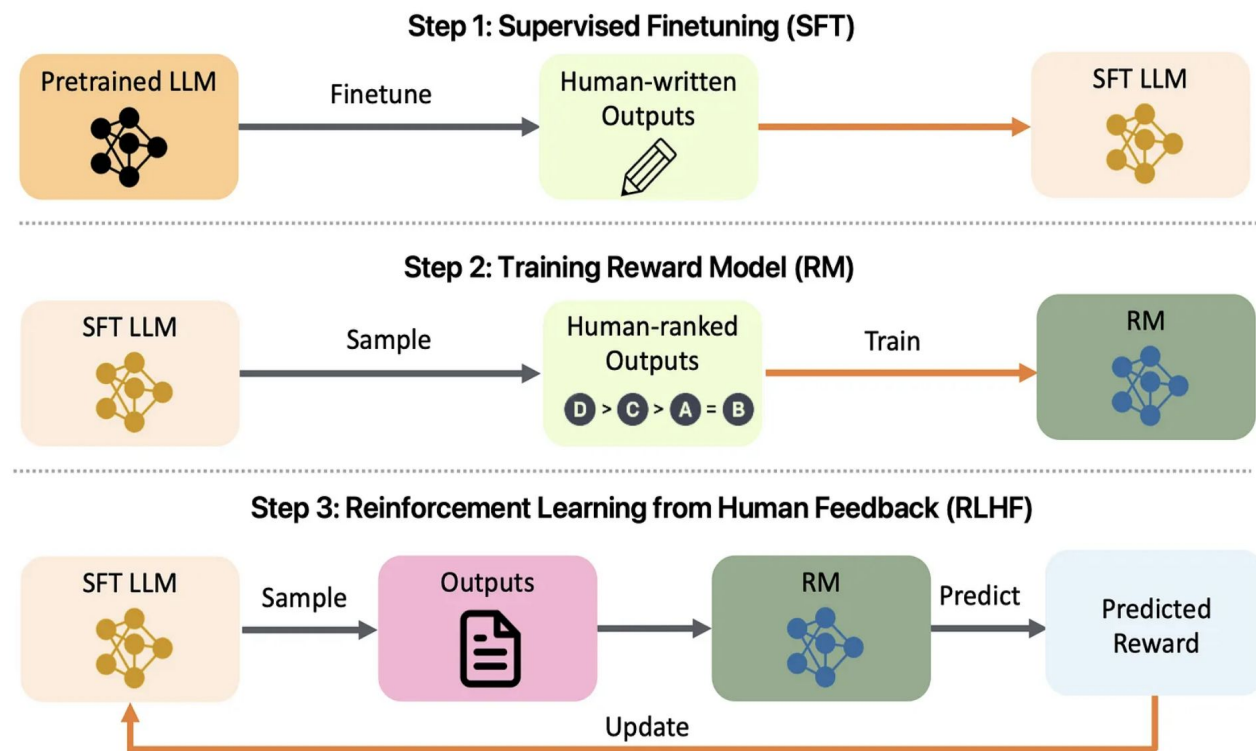


Risks	Input	Generated Response
Jailbreak	Yes	No
Prompt Leakage	No	Yes
Self-Harm	Yes	Yes
Hate/Harassment	Select Products	Yes
Violence	Select Products	Yes
Illegal Regulated Content	No	Yes
Sensitive Topics	No	Select Products, Select Topics
Discriminatory Jobs	No	Job Products

Where possible, moderation focuses on the generated content as inputs are more varied and more prone to false positives  
 Products grounding the response from curated content (e.g., LinkedIn Learning) can operate on sensitive topics

# Model Alignment

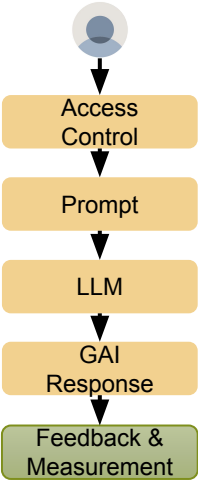
- We ensure that any foundational model we use is compliant with LinkedIn's Responsible AI Principles
- Alignment refers to the process of ensuring that LLMs behave according to human values and preferences
- Alignment focuses on fundamental risks (e.g., sexual content) and not on product specific risks (e.g., political content)
- We have a test bed to compare various methods: RLHF, DPO, KPO, etc.



Source: <https://arxiv.org/abs/2308.05374>

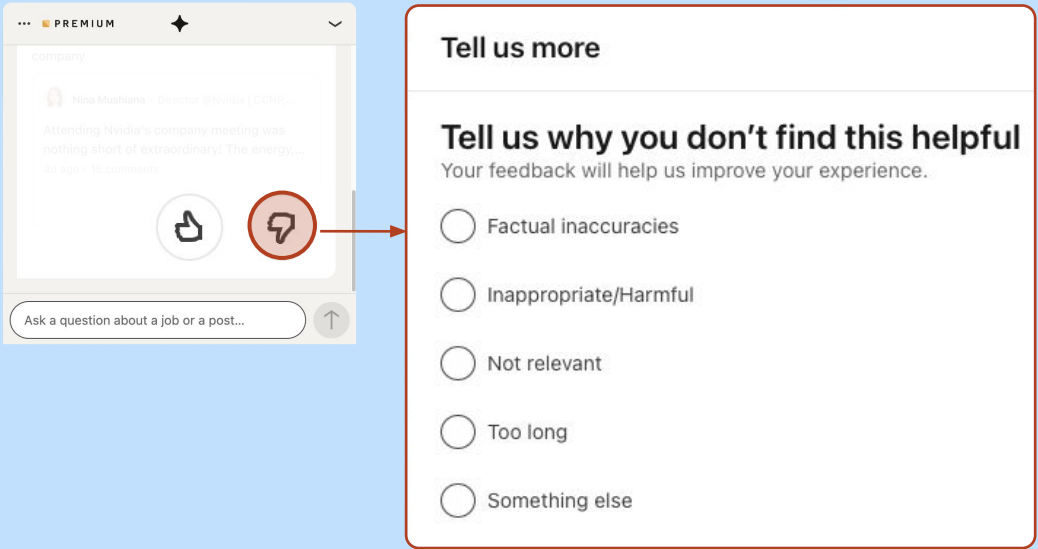
# Measurement: Feedback-Based

- User feedback metrics serve as the guardrails for all GenAI products



## Common UX & User Feedback

Standardizes UI for users and enables scalable, central oversight of feedback for all GenAI products



## Common Metrics

(Guardrail) Negative Feedback Rate measures the user perception of GenAI response quality

$$NFR = \frac{\#negative\ feedback}{\#generated\ response}$$

(Guardrail) Trust Feedback Rate measures the user perception of GenAI response harmfulness

$$TFR = \frac{\#factual\ inaccuracies + \#inappropriate}{\#generated\ response}$$

# Measurement: Bias/Stereotypes

## Incidental

Biases/stereotypes that are evident in any single piece of GenAI content. Examples:

- Straightforward: Every junior engineer should discuss his career goals with his manager.
- Subtle: New mothers ought to pause their career and stay home with their babies for 2 years

## •Challenges

- Bias/stereotypes can be very subjective and varies widely
- Difficult to collect positive examples to train a high-quality model

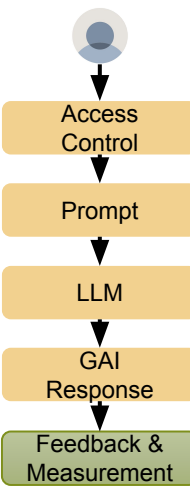
## •Solutions

- Leverage LLMs, which are trained on large amount of data and have a general understanding of bias/stereotypes (high cost)
- Leverage two-stage measurement flow to balance cost, precision, and recall
- Retrain models with balanced data
- Prompt re-writing strategies

## Representational

Skewness in demographic distribution across GenAI for an industry or a product that does not innately cater to a specific demographic group. Examples:

- Text: Articles on science careers only give examples of male scientists
- Images: doctors are always white



# Other Uses of Generative AI in Trust



# The Revolution of Generative AI

- Revolutionary technology advances opportunity, but also brings new and increased risks



LinkedIn **faces threats** by bad actors using GenAI to carry out harm through inauthentic accounts and harmful content



LinkedIn **uses** GenAI to improve its ability to measure, prevent, and mitigate abuse



LinkedIn **builds** GenAI products that need to be trustworthy and safeguarded from misuse

# AI-Assisted Sampling

## Abuse



Easily up to 95% is unknown

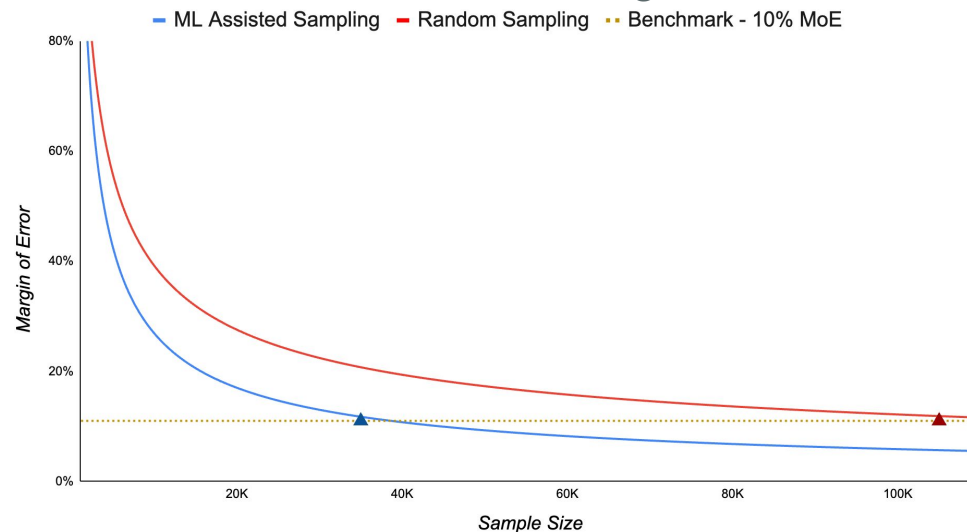
## Sampling with human labeling is expensive

- Less than 1% of samples are typically positive (high class imbalance)



## ML-assisted sampling can cut costs by 70% with tighter error bars

- Unbiased estimator using an ML model biasing our sampling



# Deepfakes

TECHNOLOGY

That smiling LinkedIn profile face might be a computer-generated fake



## Motivation

Emerging trend of large fake account attacks using automated generation and deepfake profiles photos

## Solution

Detector for AI-generated (deepfake) profile photos

- *Precision = 99.3% and recall = 85.4%*
- *Catches GAN-generated (StyleGAN1, StyleGAN2, StyleGAN3, EG3D) and diffusion-based (Stable Diffusion v1, Stable Diffusion v2, and DALL-E) images*

Engineering Blogpost: [New Approaches For Detecting AI-Generated Profile Photos](#)



# Scale Up Content Review

Creating content gets easier and cheaper, requiring trust teams to scale up their review

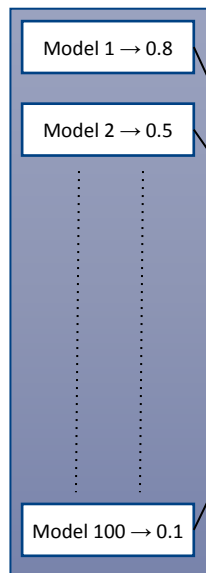


Content Components

- Profile
- Post
- Private Message



AI



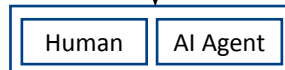
Recommendation

- Good
- Unsure
- Policy Violating

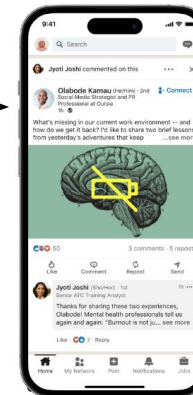
Extra Review

Review Queue

- Item 1
- Item 2
- ...
- Item N



Live Content



## Opportunity:

- Higher scalability
- Better decision quality
- Human wellbeing improvement
- More cost efficient
- Easier scaling up and down

# Thank you!

[danielo@linkedin.com](mailto:danielo@linkedin.com)

