

## Emerging Technologies & Career Advice

Peter Carragher, adapted from Trust & Safety Teaching Consortium

## Learning Objectives

**Today we will:**

- Learn about AI/ML, AR/VR, and decentralized systems
- Examine how these technologies intersect with Trust & Safety
- Explore career pathways and opportunities

- Learn about AI/ML, AR/VR, and decentralized systems
- Examine how these technologies intersect with Trust & Safety
- Explore career pathways and opportunities



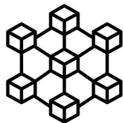
### Artificial Intelligence

AI is a rapidly growing technology that can be used to automate trust and safety processes.



### Virtual Reality

AR/VR technologies are bringing people into augmented and immersive worlds.



### Web3 and the Fediverse

Web3 and blockchain technology decentralizes data, digital assets, and entire platforms.



**Artificial Intelligence**  
AI is a rapidly growing technology that can be used to automate trust and safety processes.



**Virtual Reality**  
AR/VR technologies are bringing people into augmented and immersive worlds.

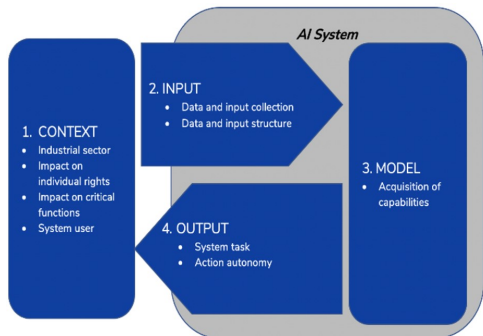


**Web3 and the Fediverse**  
Web3 and blockchain technology decentralizes data, digital assets, and entire platforms.

2026-04-18

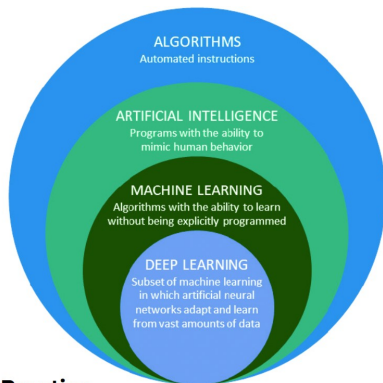
**Artificial Intelligence**

# What is Artificial Intelligence?



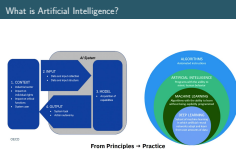
OECD

**From Principles → Practice**



2026-04-18

## What is Artificial Intelligence?



### AI System/Definitions (left)

- OECD (image from: <https://oecd.ai/en/wonk/cset-test-classification-ai-systems>)
- Other key definitions from
  - White House
  - NIST

### AI Definitions (right)

- CC license
- Image from this paper: [https://www.researchgate.net/figure/sualization-of-algorithms-vs-artificial-intelligence-vs-machine-learning-vs-de-fig1\\_339997962](https://www.researchgate.net/figure/sualization-of-algorithms-vs-artificial-intelligence-vs-machine-learning-vs-de-fig1_339997962)

# What does AI have to do with Trust & Safety?

- Identifying and scoring content
- Flagging content for review
- Actioning content
- Assisting humans with decisions

- Content *recommendation*
- Equity, fairness, 'bias'/harm
- Generating synthetic content to test
- More than just 'content'

2026-04-18

## └─What does AI have to do with Trust & Safety?

What does AI have to do with Trust & Safety?

"Content moderation", but to break that down...

- Identifying and scoring content
- Flagging content for review
- Actioning content
- Assisting humans with decisions

But also:

- Content recommendation
- Equity, fairness, 'bias'/harm
- Generating synthetic content to test
- More than just 'content'

Reactive + Proactive Approaches

# Example 1



## How Pinterest fights misinformation, hate speech, and self-harm content with machine learning

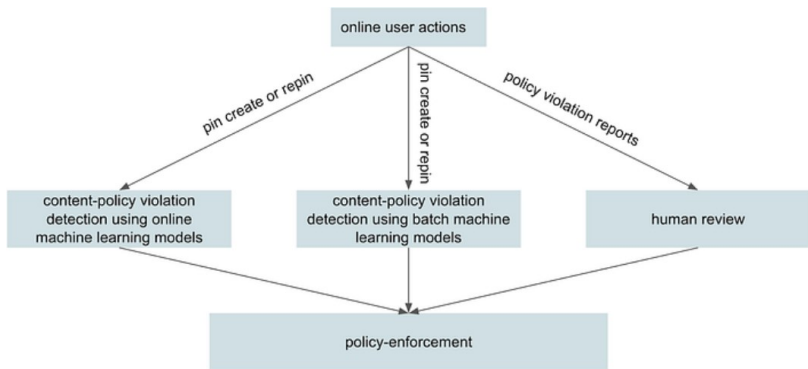
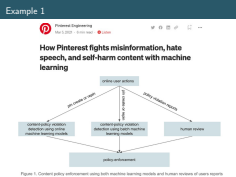


Figure 1. Content policy enforcement using both machine learning models and human reviews of users reports

2026-04-18

## Example 1



From: <https://medium.com/pinterest-engineering/how-pinterest-fights-misinformation-hate-speech-and-self-harm-content-with-machine-learning>

Key points from post:

- Combination of ML + user reports
- Detect unsafe content before it's reported
- Combination of methods, data sources signals
- Signals conveyed from keywords, images, etc
- Key questions around measurement

# Example 2

2026-04-18

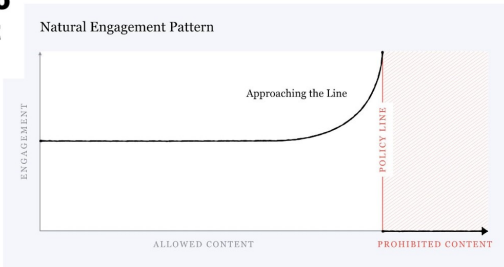
## Example 2

Example 2

Facebook will change algorithm to demote "borderline content" that almost violates policies



**Facebook will change algorithm to demote "borderline content" that almost violates policies**



Usage of AI/ML to identify borderline content + prohibited content



TRY IT OUT

CREATE CUSTOM DEMO

Type here to see the potential effect of your comment.

\*Please use this demo only if you are older than 13. Please do not enter personally identifiable information.

### Guidelines for Responsible Conversational AI

AI is rapidly becoming normalised. In fact, it's more common than you think. Conversational AI is becoming normalised, and has been for some time. Conversational AI is improving and we are seeing people using bots in their everyday lives.

While the potential benefits, convenience and innovation also important to approach any AI solution with care and interactive demo that highlight a selection of the Microsoft AI solutions touching on topics such as transparency, reliability, safety

Go to the demo >



Try out one of these services/learning demos to learn more about:

- Toxicity/harmful speech detection - <https://perspectiveapi.com/>
- MSFT conversational responsible AI demos
- <https://teachablemachine.withgoogle.com/>

## Discussion Questions

- What types of T&S-related topics can (and should) AI/ML based systems address?
- What is *not* well suited to be addressed by AI or automated methods in T&S?
- Policy development and enforcement often require judgements about context. What are important contextual signals that should be considered in human/machine environments?
- How would you measure the effectiveness of AI systems for T&S objectives?
- How should T&S practitioners think about differential impact across communities, including markets? Like those in the Global South, where training data may be less reflective of specific populations

- What types of T&S-related topics can (and should) AI/ML based systems address?
- What is *not* well suited to be addressed by AI or automated methods in T&S?
- Policy development and enforcement often require judgements about context. What are important contextual signals that should be considered in human/machine environments?
- How would you measure the effectiveness of AI systems for T&S objectives?
- How should T&S practitioners think about differential impact across communities, including markets? Like those in the Global South, where training data may be less reflective of specific populations

2026-04-18

Virtual and Augmented  
Reality

<https://newpublic.substack.com/p/a-social-network-taxonomy>



- **Augmented Reality**  
Overlaying digital information onto the real world, typically through a device like a smartphone or a headset
- **Virtual Reality**  
Fully immersive experiences in a simulated environment, typically through a dedicated VR headset

## Defining Augmented and Virtual Reality

2026-04-18

# Defining Augmented and Virtual Reality



- **Augmented Reality**  
Overlaying digital information onto the real world, typically through a device like a smartphone or a headset
- **Virtual Reality**  
Fully immersive experiences in a simulated environment, typically through a dedicated VR headset

# Example Harm Areas of AR and VR



- Harassment and bullying
- Inappropriate content
- Privacy concerns
- Age restrictions and child safety
- Accessibility and inclusivity

2026-04-18

## Example Harm Areas of AR and VR



In immersive and augmented spaces, long-standing areas of harm take on new dimensions:

- Harassment and bullying
- Inappropriate content
- Privacy concerns
- Age restrictions and child safety
- Accessibility and inclusivity

## Discussion Questions

- How do trust and safety concerns differ between virtual reality (VR), augmented reality (AR), and traditional digital environments? Discuss the unique challenges and opportunities that immersive technologies present for content moderation.
- What are the ethical considerations surrounding the collection and use of user data in VR and AR environments? How can platforms ensure that they respect user privacy while promoting safety?
- Discuss the role of AI and machine learning in moderating VR and AR experiences. What are the benefits and limitations of these technologies, and how can they aid human moderation?
- How can VR and AR platforms foster a culture of empathy, respect, and inclusivity among users? Discuss the role of virtual etiquette, social norms, and digital citizenship in creating a safe and welcoming environment for all users.

- How do trust and safety concerns differ between virtual reality (VR), augmented reality (AR), and traditional digital environments? Discuss the unique challenges and opportunities that immersive technologies present for content moderation.
- What are the ethical considerations surrounding the collection and use of user data in VR and AR environments? How can platforms ensure that they respect user privacy while promoting safety?
- Discuss the role of AI and machine learning in moderating VR and AR experiences. What are the benefits and limitations of these technologies, and how can they aid human moderation?
- How can VR and AR platforms foster a culture of empathy, respect and inclusivity among users? Discuss the role of virtual etiquette, social norms, and digital citizenship in creating a safe and welcoming environment for all users.

2026-04-18

Decentralized  
Technologies

Web3 and blockchain technology provides a decentralized way to structure platforms, data, and digital currencies.

It can also be used to create and manage ownership of digital assets that are used on tech platforms and even in-game experiences.

These novel technologies also introduce new abuse vectors that you may face as a Trust & Safety professional.

2026-04-18

## Blockchain, Cryptocurrencies, NFTs

# Blockchain, Cryptocurrencies, NFTs

## Blockchain, Cryptocurrencies, NFTs

- **Blockchain Technology** provides a decentralized and secure way to store data and conduct transactions.
- **Key Attributes**
  - Decentralization
  - Encryption
  - Consensus mechanisms
- **Cryptocurrency** Digital or virtual currency that uses cryptography for security and operates on a decentralized network.
- **NFTs** 'Non-Fungible Tokens', or a digital asset with a unique identifier.

- **Blockchain Technology** provides a decentralized and secure way to store data and conduct transactions.
- **Key Attributes**
  - Decentralization
  - Encryption
  - Consensus mechanisms
- **Cryptocurrency** Digital or virtual currency that uses cryptography for security and operates on a decentralized network.
- **NFTs** 'Non-Fungible Tokens', or a digital asset with a unique identifier.



Example Technologies:

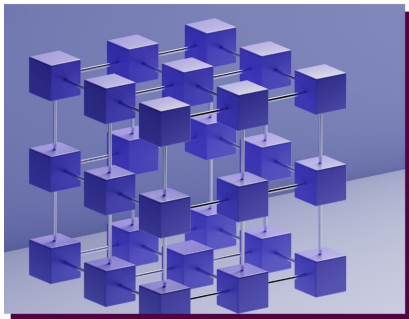
- **Cryptocurrency:** Digital or virtual currency that uses cryptography for security and operates on a decentralized network.
- **NFTs:** 'Non-Fungible Tokens', or a digital asset with a unique identifier.

2026-04-18

## Blockchain, Cryptocurrencies, NFTs



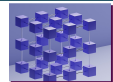
- **Cryptocurrency:** Digital or virtual currency that uses cryptography for security and operates on a decentralized network.
- **NFTs:** 'Non-Fungible Tokens', or a digital asset with a unique identifier.



- Decentralization can raise major harm vectors, including:
  - Account Stealing and other platform integrity vulnerabilities
  - Financial Fraud and Scams
  - Inauthentic Behavior

2026-04-18

## Blockchain, Cryptocurrencies, NFTs



- Decentralization can raise major harm vectors, including:
  - Account Stealing and other platform integrity vulnerabilities
  - Financial Fraud and Scams
  - Inauthentic Behavior

# Discussion Questions

- How do trust and safety concerns manifest in the context of cryptocurrencies, NFTs, and digital assets? Discuss the unique challenges these technologies present compared to traditional financial systems and assets.
- Explore the potential risks associated with the use of cryptocurrencies and NFTs, such as fraud, money laundering, and illicit activities. How can these risks be mitigated while preserving the benefits of decentralization and user autonomy?
- Discuss the role of self-regulation in the cryptocurrency and NFT ecosystem. How can industry stakeholders, such as exchanges, wallet providers, and developers, contribute to trust and safety?
- Analyze the implications of content moderation for NFTs and digital assets, particularly in relation to intellectual property rights, copyright infringement, and counterfeit goods. How can platforms balance free expression with the protection of creators' rights?
- Discuss the importance of user education and awareness in promoting trust and safety in the cryptocurrency and NFT space. How can platforms, developers, and regulators collaborate to empower users to make informed decisions?

2026-04-18

## Discussion Questions

Discussion Questions

- How do trust and safety concerns manifest in the context of cryptocurrencies, NFTs, and digital assets? Discuss the unique challenges these technologies present compared to traditional financial systems and assets.
- Explore the potential risks associated with the use of cryptocurrencies and NFTs, such as fraud, money laundering, and illicit activities. How can these risks be mitigated while preserving the benefits of decentralization and user autonomy?
- Discuss the role of self-regulation in the cryptocurrency and NFT ecosystem. How can industry stakeholders, such as exchanges, wallet providers, and developers, contribute to trust and safety?
- Analyze the implications of content moderation for NFTs and digital assets, particularly in relation to intellectual property rights, copyright infringement, and counterfeit goods. How can platforms balance free expression with the protection of creators' rights?
- Discuss the importance of user education and awareness in promoting trust and safety in the cryptocurrency and NFT space. How can platforms, developers, and regulators collaborate to empower users to make informed decisions?

2026-04-18

## └─ The Fediverse

The Fediverse

The advent of new technologies includes converging with what has come so far. Platforms still create opportunities to create content, but now they're federated across different servers.

Content, games, art, and currency is built on the blockchain. Ownership is in the hands of users, granting greater user freedoms and ownership, and challenging traditional models of keeping platforms safe.

In this section, we go over some of the promises and challenges of decentralized technologies.

- **What does it mean to be “Federated”?**
  - Decentralized and interoperable network of independent social media platforms.
  - Communicate using open-source protocols.
  - Enables users to join various self-hosted instances or servers.
- **Moderation**
  - Usually team of volunteer moderators and administrators.
  - Set policies and enforce them within federated instance/server.

2026-04-18

└ Fediverse

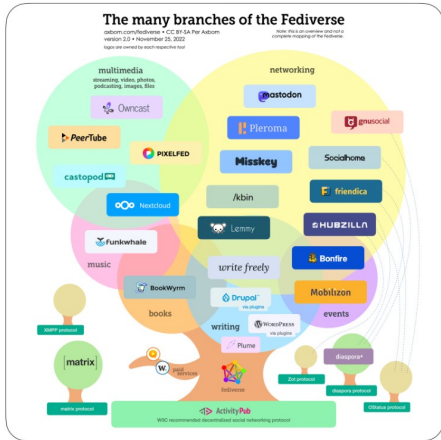
- **What does it mean to be “Federated”?**
  - Decentralized and interoperable network of independent social media platforms.
  - Communicate using open-source protocols.
  - Enables users to join various self-hosted instances or servers.
- **Moderation**
  - Usually team of volunteer moderators and administrators.
  - Set policies and enforce them within federated instance/server.

# Fediverse and Content Moderation

2026-04-18

## └ Fediverse and Content Moderation

- Fediverse communities are governed by administrators or moderators
- Similar to internet forums, community platforms (e.g. Facebook Pages, Reddit subreddits, Discord servers)
- What's different? There's a limited centralized Trust & Safety function



- Fediverse communities are governed by administrators or moderators
- Similar to internet forums, community platforms (e.g. Facebook Pages, Reddit subreddits, Discord servers)
- What's different? There's a limited centralized Trust & Safety function



- What are the unique challenges of moderating content on a federated platform compared to centralized platforms? Discuss the advantages and disadvantages of both systems.
- Discuss the role of human moderators and automated systems in maintaining trust and safety on federated platforms. How can these two approaches complement each other?
- Analyze the concept of content moderation within the context of free speech and censorship. To what extent should federated platforms moderate content, and what principles should guide their decisions?
- What are the potential risks associated with self-governance on federated platforms? How can these platforms mitigate these risks while preserving user autonomy?
- Discuss the potential legal and ethical implications of content moderation on federated platforms. How should these platforms navigate issues related to privacy, data protection, and liability?

## Discussion Questions

- What are the unique challenges of moderating content on a federated platform compared to centralized platforms? Discuss the advantages and disadvantages of both systems.
- Discuss the role of human moderators and automated systems in maintaining trust and safety on federated platforms. How can these two approaches complement each other?
- Analyze the concept of content moderation within the context of free speech and censorship. To what extent should federated platforms moderate content, and what principles should guide their decisions?
- What are the potential risks associated with self-governance on federated platforms? How can these platforms mitigate these risks while preserving user autonomy?
- Discuss the potential legal and ethical implications of content moderation on federated platforms. How should these platforms navigate issues related to privacy, data protection, and liability?

2026-04-18

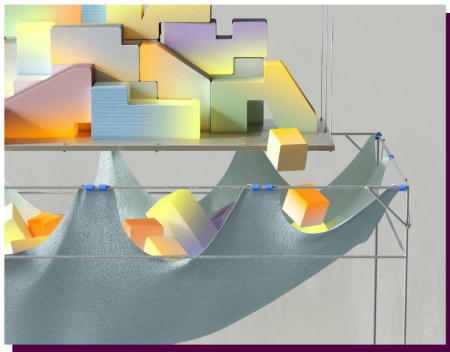
Career Advice

2026-04-18

## └ Careers in Trust & Safety

In industry, Trust & Safety professionals are responsible for ensuring the safety of users and clients in the policies they form, products they build, and stories they tell in telling an organization's story to the world.

As new technologies emerge and novel forms of harm proliferate, the impact of Trust & Safety work will be more important than ever.



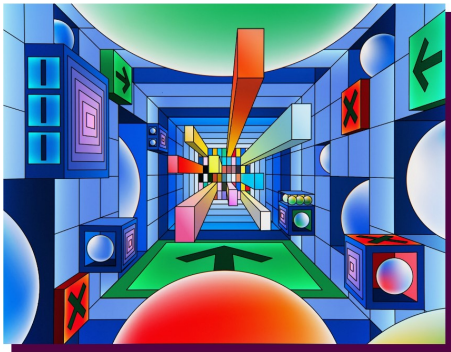
- **Legal:** Guide company across teams on risk mitigation, develop responses to regulatory bodies and law enforcement.
- **Policy:** Design content policies, performing risk assessments for product teams, and forming partnerships with civil society and governments.
- **Operations:** Enforce policy through scaled operational programs. Dedicated teams across subject-matter expertise, program and project management.

2026-04-18

## Trust & Safety Roles in Tech



- **Legal:** Guide company across teams on risk mitigation, develop responses to regulatory bodies and law enforcement.
- **Policy:** Design content policies, performing risk assessments for product teams, and forming partnerships with civil society and governments.
- **Operations:** Enforce policy through scaled operational programs. Dedicated teams across subject-matter expertise, program and project management.



- **Product & Engineering:** Builds internal and consumer-facing products, including T&S tooling to safety features
- **Data Science:** Analyzes policy violations and their impact on platforms, develops measurement and detection methods, and create strategies to address harms.
- **Marketing & Communications:** Develops narrative strategies for showcasing organization's safety commitments. Drives content strategies and user research efforts.

2026-04-18

## Trust & Safety Roles in Tech



- **Product & Engineering:** Builds internal and consumer-facing products, including T&S tooling to safety features
- **Data Science:** Analyzes policy violations and their impact on platforms, develops measurement and detection methods, and create strategies to address harms.
- **Marketing & Communications:** Develops narrative strategies for showcasing organization's safety commitments. Drives content strategies and user research efforts.

2026-04-18

## Trust & Safety Partnerships

Working in Trust & Safety doesn't have to be limited to working at a tech company.

In fact, tech companies need extensive collaborations and partnerships across industry, academic, civil society, and government partners.



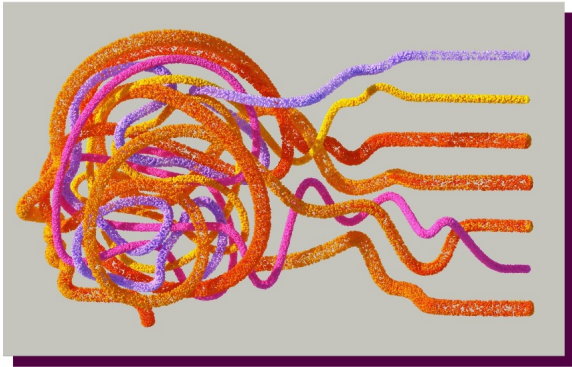
- Fellow tech companies and industry working groups
- Coordinate on operations and investigations
- Collaborate in industry working groups (e.g. TSPA) and share best practices
- Collaborate on guiding regulators

2026-04-18

## Industry Partnerships



- Fellow tech companies and industry working groups
- Coordinate on operations and investigations
- Collaborate in industry working groups (e.g. TSPA) and share best practices
- Collaborate on guiding regulators



- Partner on user research and social impact
- Guide product teams
- Audit content policies and public policy initiatives
- Example partners: academic labs, PhD and Masters students, specialized researchers and research projects

2026-04-18

## Academic Partnerships

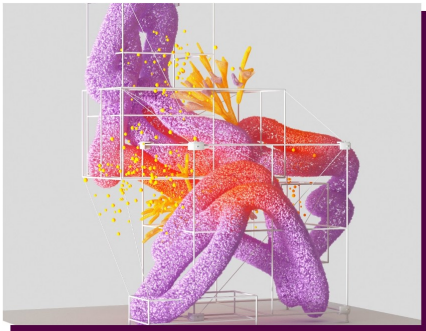


- Partner on user research and social impact
- Guide product teams
- Audit content policies and public policy initiatives
- Example partners: academic labs, PhD and Masters students, specialized researchers and research projects



- Create guiding frameworks for policy, product, data privacy, Trust & Safety operations
- Co-create education resources and industry reports to guide industry awareness of issues
- Example Partners: subject-matter focused (e.g. child safety), product-based (e.g. organizations developing AI product and policy frameworks)

# Civil Society Partnerships



- Create guiding frameworks for policy, product, data privacy, Trust & Safety operations
- Co-create education resources and industry reports to guide industry awareness of issues
- Example Partners: subject-matter focused (e.g. child safety), product-based (e.g. organizations developing AI product and policy frameworks)

2026-04-18

## Civil Society Partnerships



- Regulate platform content moderation practices and safety products
- Issue law enforcement inquiries (e.g. subpoenas, exigent requests)
- Example partners: local, provincial, national, and international governing bodies, law enforcement agencies

2026-04-18

## └ Governments

Governments



- Regulate platform content moderation practices and safety products
- Issue law enforcement inquiries (e.g. subpoenas, exigent requests)
- Example partners: local, provincial, national, and international governing bodies, law enforcement agencies

- **Blockchain** icon created by Good Ware
- **Brain icon** created by Freepik
- **VR** icon created by Nikita Golubev
- Fediverse map by Per Axbom, <https://axbom.com/fediverse/>
- Career Advice section art from DeepMind's Visualising AI project, <https://visualisingai.deepmind.com>

- Blockchain icon created by Good Ware
- Brain icon created by Freepik
- VR icon created by Nikita Golubev
- Fediverse map by Per Axbom, <https://axbom.com/fediverse/>
- Career Advice section art from DeepMind's Visualising AI project, <https://visualisingai.deepmind.com>