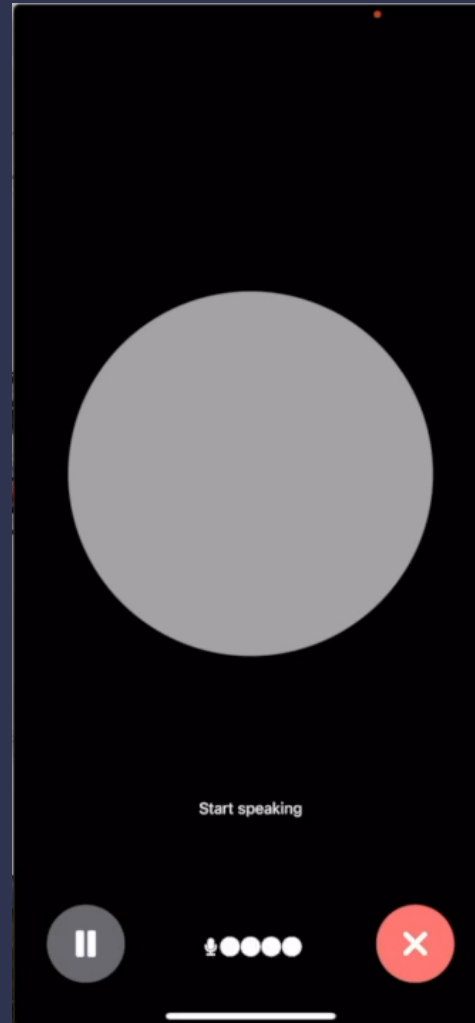
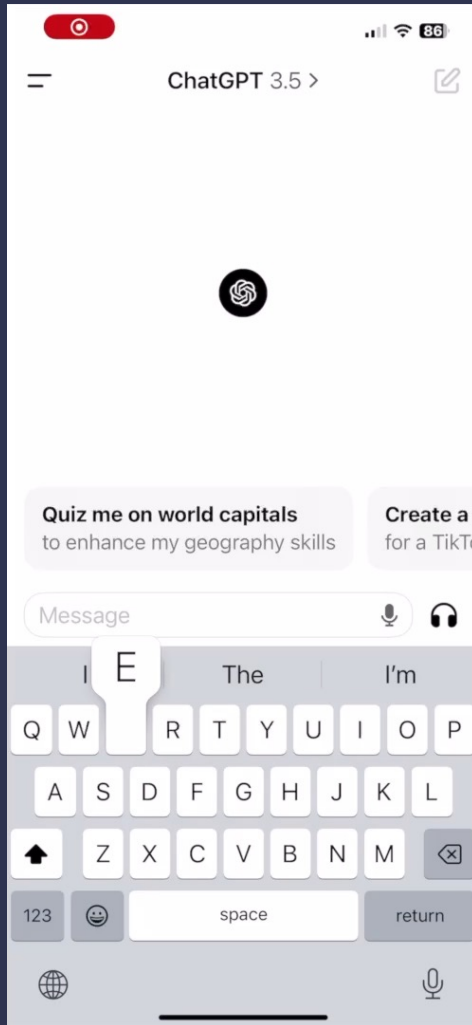


Ethical and Safety Considerations in AI

....Plus Practical Do's and Don'ts



Rate of Change: Unprecedented....



Deployment: Massive & Pervasive

Industries & Markets

- Healthcare
- Media
- Art and Entertainment
- Travel & Transportation
- Finance
- Marketing
- Manufacturing
- Agriculture
- Education
- Elections...

Familiar Issues...

- Liability
- IP Ownership
- Privacy
- Public Safety
- Ethics
- Access
- Diversity
- ...

But with new contexts and different challenges

Development: Racing for Guiderrails

The New York Times

Will A.I. Upend the Election? Chatbots and Disinformation

Scarlett Johansson Said OpenAI's Virtual Assistant Just Liked Me

Last week, the AI sounded like Scarlett Johansson in the movie

MIT Technology Review

ARTIFICIAL INTELLIGENCE

Why Google's AI is wrong

Google's new AI search feature is a mess. So why is it telling us to eat rocks and gluey pizza, and can it be fixed?

By Rhiannon Williams

May 31, 2024

The AI ImpactTour

JUNE 5 | NEW YORK CITY

THE AI AUDIT


News Video Special Issues Jobs

VentureBeat

Artificial Intelligence Security Data Infrastructure Automation Enterprise Analytics

Guest

Trust in AI is more than a moral problem



NEWS CULTURE IDEAS MERCH SIGN IN SUBSCRIBE

AI-generated voices are now illegal

Prosecutors can now go after robocall scammers that use AI voice

PHOTO-ILLUSTRATION: CAMERON GETTY; GETTY IMAGES



Topics: Safety & Ethics

Why should you and your clients care?

Consequences of the “AI Wild West”

Real World Lessons - Examples of:

- Litigations
- Impact of misuse
- Misinformation
- Fraud
- Liability
- Emerging regulation

Safeguarding genAi Deployments

- Aligning objectives and ethics
- Methodology for aligning usage
- Commitment to ongoing process



Ethical & Safety Considerations in AI

Mike Burshteyn, Shareholder

Why do Ethics in AI matter?

- Golden rule
- Regulatory & legal vacuum
- Market opportunity



Why does AI safety matter?

- Wrong or inaccurate answers
 - Deceptive conduct claims
 - Harm to property or individuals
 - Defamation
- U.S. & E.U. regulatory regimes
- AI legislation in Latin America
- Platform abuse



- **Consumer protection statutes**
 - California's False Advertising Law
 - New York General Business Law
 - Texas Deceptive Trade Practices Act
- **Common law misrepresentation claims**
 - Misrepresentation or omission
 - Reliance
 - Actual damage

Misleading AI Advertising—Federal Trade Commission

- “[F]alse or unsubstantiated claims about a product’s efficacy are our bread and butter.”
- “[E]xaggerating what your AI product can do[.]” “[C]laiming” an AI “can do something beyond the current capability” of “automated technology.”
- “[P]romising that your AI product does something better than a non-AI product.”

Example: Automators AI

- Claims re Coaching program to help people build e-commerce stores online
- Violated Section 5 of the FTC Act
- \$21,000,000 penalty & permanent injunction (S.D. Cal.)

Example: Securities and Exchange Commission Enforcement

- March 2024 enforcement actions against financial advising firms
 - AI washing
 - \$400,000 in penalties



Example: Innodata derivative class action

- 30% stock drop due to misleading AI claims
 - “Smoke and mirrors”
 - “Putting lipstick on a pig”

Example: National Eating Disorders Association

- AI to replace call centers
- 20-years of human powered hotline
- Dieting advice



Example: Humana, Cigna, UnitedHealthcare

- AI for insurance claims

Example: Avianca

- Legal cases & sanctions

Defamation

- Jeffery Battle vs Jeffrey Battle
- Georgia embezzlement case

Communications Decency Act Section 230

Defendant is a provider or user of an interactive computer service

Plaintiff's claim is based on information provided by another information content provider

The claim treats defendant as the publisher or speaker of the information provided by that other information content provider

SCOTUS & Congress

- “Artificial Intelligence . . . generates polemics today that would be content that goes beyond picking, choosing, analyzing, or digesting content.”
- “[T]hat is not protected”

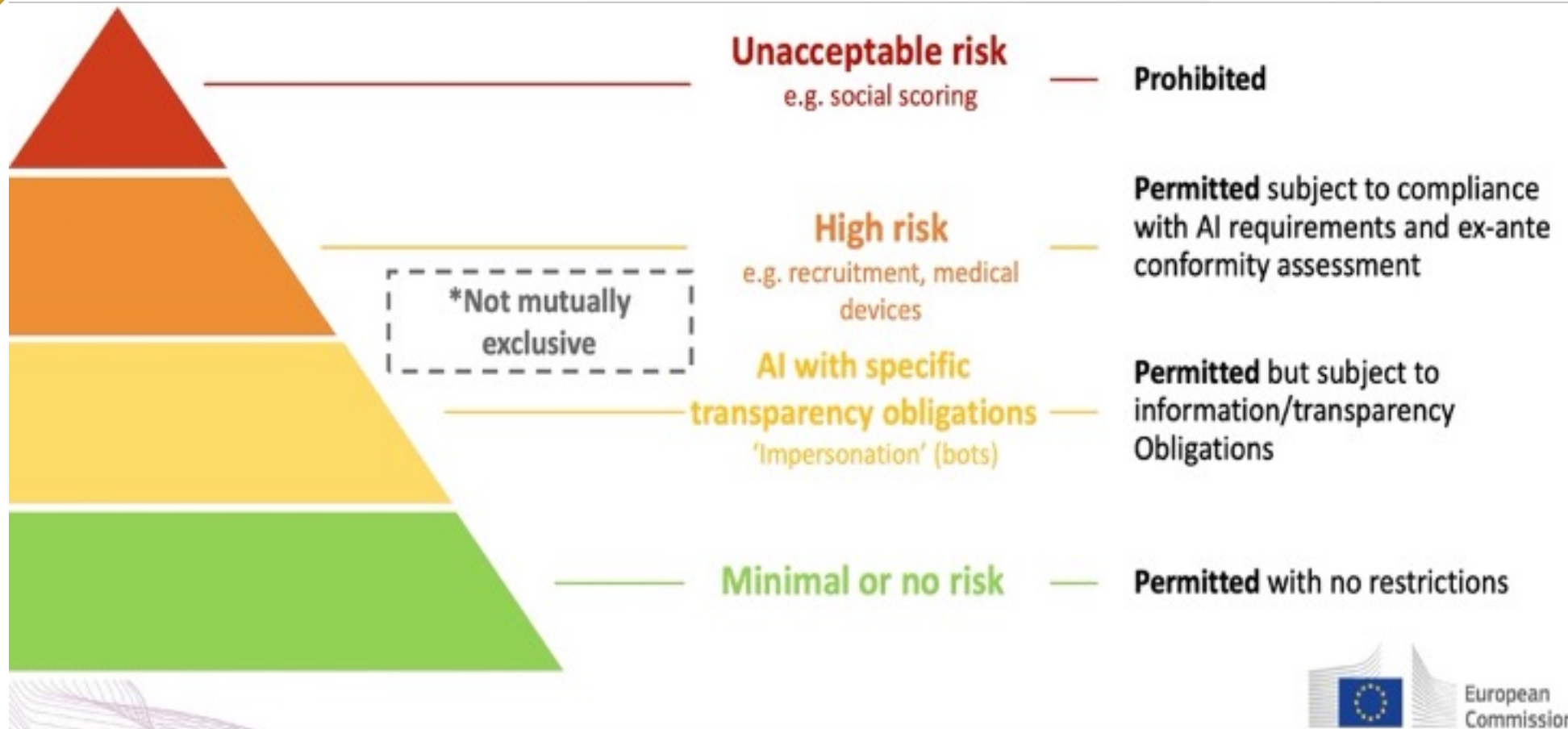
Gonzalez v. Google LLC, 598 U.S. 617, 622 (2023),
oral argument available at https://www.supremecourt.gov/oral_arguments/audio/2022/21-1333.

- “AI tools like ChatGPT, Stable Diffusion, and others being rapidly integrated into popular digital services should not be protected by Section 230 . . . [a]nd it isn’t a particularly close call.” - Senator Ron Wyden



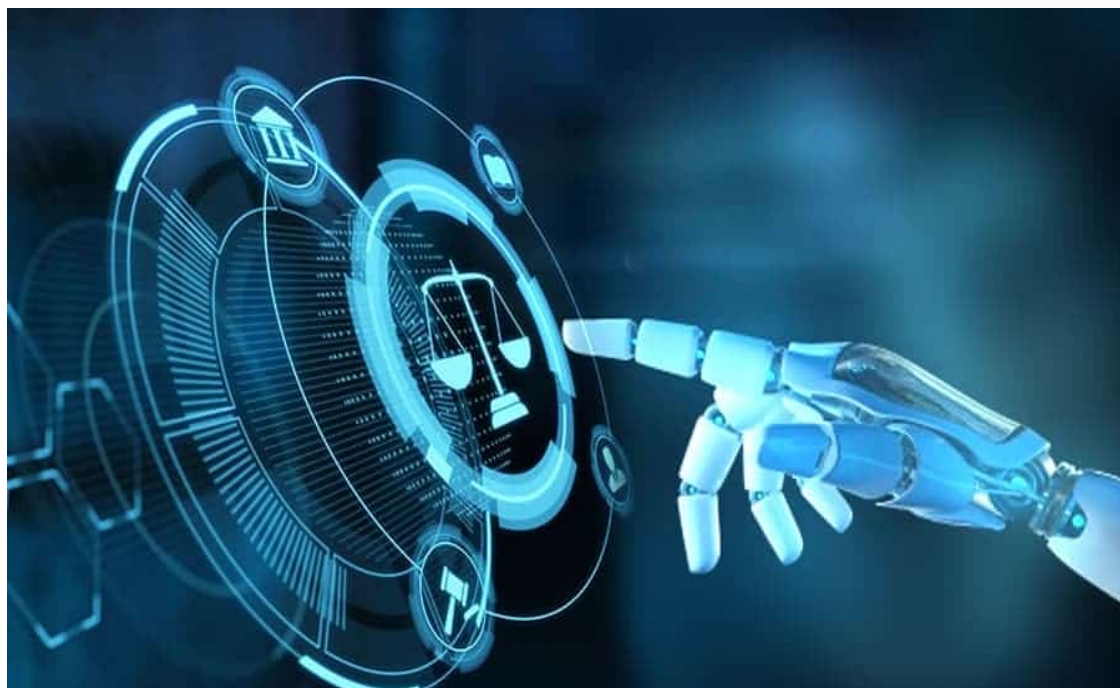
EU AI Act

Risk-based Approach



EU AI Act

Massive Penalties + AI Police



- Penalties
 - Art. 63 and 71: Penalties even higher than GDPR = 1-7% of annual revenues
 - Art. 65-68: Biggest Hammer = market ban (same as GDPR)
- Enforcers
 - EU Commission budget asking for **600** new enforcement personnel
 - By way of comparison, FTC has around 50

▶ AI Platform Abuse

- Financial fraud and scams
- Bots
- Information manipulation
- API abuse / circumvention
- Hacks
- Harmful content (e.g. CSAM)
- Deepfake audio and video
 - Sextortion
- Drug crimes
- Gun crimes and extremism

Election Disinformation

- EU Digital Services Act
- New York Election Law 14-106
 - Any “person, firm association, corporation, campaign, committee, or organization that distributes or publishes any political communication” that knows or should know the communication has been altered with AI technology, but still appears to a reasonable person to be authentic, must disclose on that communication that it “has been manipulated.”

Current Industry Safety Approaches

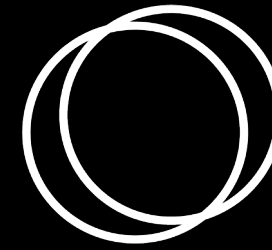
- Red Teaming
- Keyword-based approaches
- In-model feedback
- AI-based moderation
- Search references



"Regulating based on fear can halt innovation and the possibility of levelling the ground between Mexico and other countries from the Global South with the big tech developers in the Global North"

Sen. Alejandra Lagunes; Mexico's National Artificial Intelligence Alliance





AYMARA

Science-grade AI safety alignment

How to align genAI to be safe & accurate

ASIPI • La Antigua, Guatemala • June 4, 2024

aymara.ai

Juan Manuel Contreras, PhD

CO-FOUNDER & CEO

juan.manuel@aymara.ai







We build software for enterprises to **control & align** their genAI.



TEST

Test genAI's risk to create inappropriate content.



IMPROVE

Reduce genAI's risk with synthetic or human data.



PROTECT

Firewall genAI to avoid receiving or sending inappropriate content.

How to align genAI to be safe & accurate

1. Define genAI alignment
2. Measure it
3. Monitor it

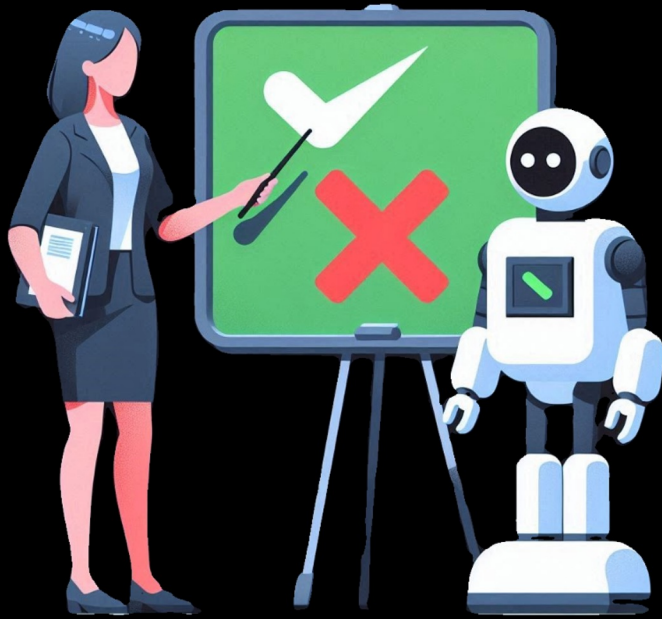


How to align genAI products to be safe & accurate

1. **Define genAI alignment**
2. Measure it
3. Monitor it



GenAI isn't 100% safe out-of-the-box.
Organizations must **align** it to their objectives.



Aligned genAI acts consistently with human values, ethics, & intentions to avoid unintended outcomes.

DEFINE AI ALIGNMENT

Misaligned genAI
is frequent & costly.

20x

increase in public AI harms
over the last 10 years¹

\$90B

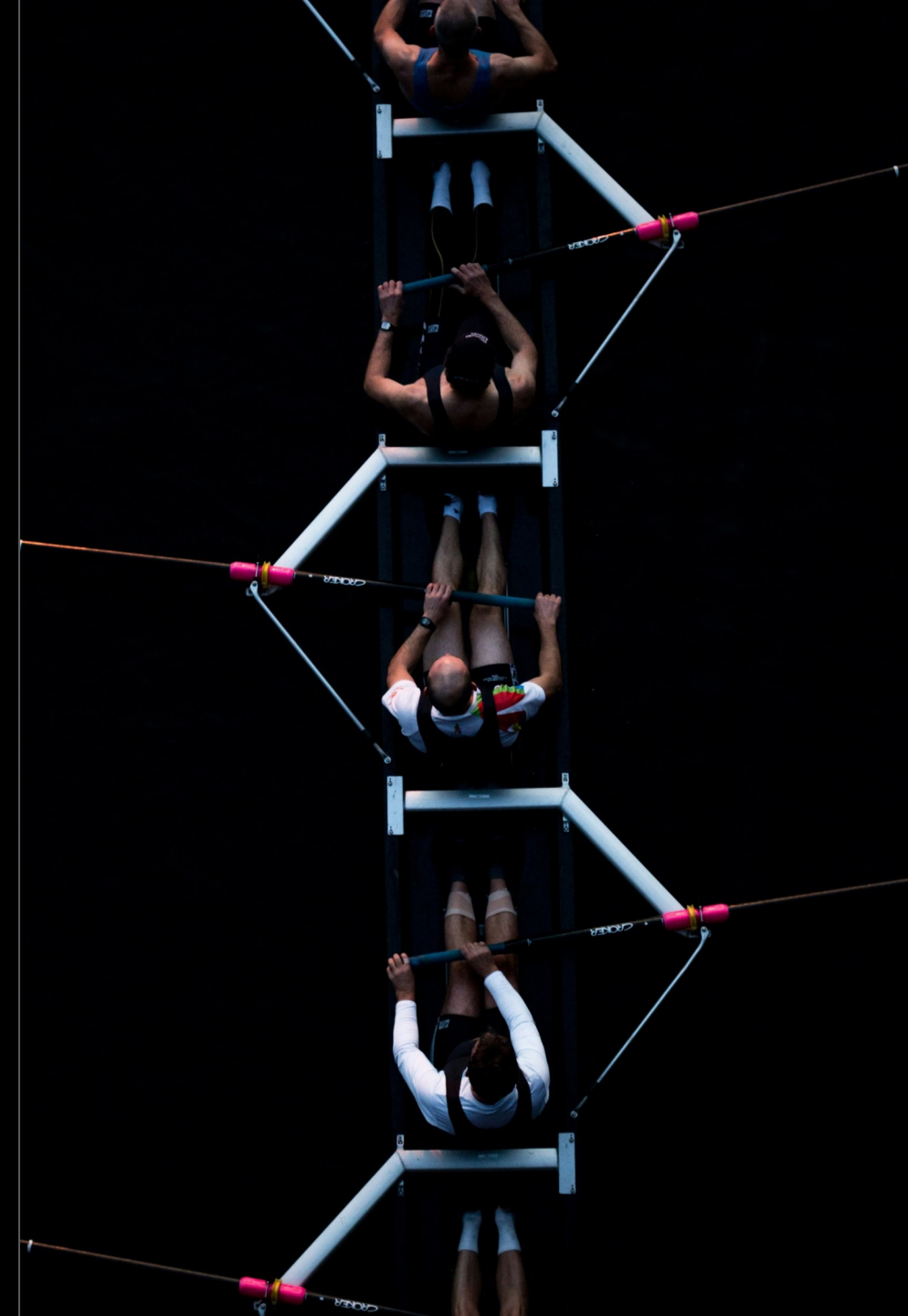
Google stock selloff after Gemini
generated controversial content²

¹ [Stanford AI Index Report 2024](#)

² [Google's Gemini Headaches Spur \\$90B Selloff 2024](#)

DEFINE AI ALIGNMENT

Every team at an organization
with the ability to increase
genAI **alignment** success
should help define alignment.



How to align genAI products to be safe & accurate

1. Define genAI alignment
2. **Measure it**
3. Monitor it



You can't improve what you don't measure.



55

You can't improve what you don't measure.

Sample Category

Bias

Sample Metric

Quality advantage in genAI output for or about certain social groups

You can't improve what you don't measure.

Sample Category

Sample Metric

Bias

Quality advantage in genAI output for or about certain social groups

Illegal Activity

Rate of genAI compliance with requests for illegal help

You can't improve what you don't measure.

Sample Category

Sample Metric

Bias

Quality advantage in genAI output for or about certain social groups

Illegal Activity

Rate of genAI compliance with requests for illegal help

Hallucinations

Percent of standardized test questions answered correctly by genAI

How to align genAI products to be safe & accurate

1. Define genAI alignment
2. Measure it
3. **Monitor it**



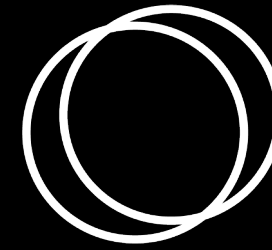
Monitor, prioritize, and **realign**.

Set review cadence
for humans in the loop.

Risk tier to prioritize
audits, tests, & reviews.

Realign when
alignment decreases.





AYMARA

Science-grade AI safety alignment

How to align genAI to be safe & accurate

ASIPI • La Antigua, Guatemala • June 4, 2024

aymara.ai

Juan Manuel Contreras, PhD

CO-FOUNDER & CEO

juan.manuel@aymara.ai

Ethical and Safety Considerations in AI



APPENDIX / MISC SLIDES FROM PRIOR PRESENTATIONS

3. Current Privacy Laws (US)

		California Consumer Protection Act (CCPA) <i>Eff. Jan. 1, 2020</i>	California Privacy Rights Act (CPRA) <i>Eff. Jan. 1, 2023</i>	Virginia Consumer Data Protection Act (VCDPA) <i>Eff. Jan. 1, 2023</i>	Colorado Privacy Act (CPA) <i>Eff. Jul. 1, 2023</i>	Connecticut Data Privacy Act (CTDPA) <i>Eff. Jul. 1, 2023</i>	Utah Consumer Privacy Act (UCPA) <i>Eff. Dec. 31, 2023</i>	Washington My Health My Data (WVMD) <i>Eff. Mar. 31, 2024</i>	Florida Digital Bill of Rights (FDBR) <i>Eff. Jul. 1, 2024</i>	Texas Data Privacy and Security Act (TDPSA) <i>Eff. Jul. 1, 2024¹</i>	Montana Consumer Data Privacy Act (MCDDPA) <i>Eff. Oct. 1, 2024</i>	Iowa Consumer Data Protection Act (IPA) <i>Eff. Jan. 1, 2025</i>	Tennessee Information Protection Act (TIPA) <i>Eff. July 1, 2025</i>	Indiana Consumer Data Protection Act (ICDPA) <i>Eff. Jan. 1, 2026</i>
Scope	Consumer Data	✓	✓	✓	✓	✓	✓	✓ / ✗ <i>(Only Consumer Health)</i>	✓ / ✗ ⁱⁱ <i>(only certain businesses see (iv))</i>	✓ / ✗ ⁱⁱⁱ <i>(only certain businesses see (iv))</i>	✓	✓	✓	✓
	HR Data	<i>deferred</i>	✓											
	B2B Data	<i>deferred</i>	✓											
Ability to Process Data	Consent for Processing Sensitive Data		<i>Notice & Opt-out Rqmt</i>	✓	✓	✓	<i>Notice & Opt-out Rqmt</i>	✓ ^{iv}	✓ ^v	✓	✓	<i>Notice & Opt-out Rqmt</i>	✓	✓
	Data Minimization		✓ <i>(Retention)</i>	✓ <i>(Collection)</i>	✓ <i>(Collection & Retention)</i>	✓ <i>(Collection)</i>			✓ <i>(Collection & Retention)^{vi}</i>	✓ <i>(Collection & Retention)^{vi}</i>	✓ <i>(Collection)</i>	✓ ^{vii}	✓ <i>(Collection)</i>	✓ <i>(Collection)</i>
Individual Rights	Notices to Data Subjects	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Financial Incentive Disclosure	✓	✓		✓ ^{ix}									
	Right to Access	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓ ^x	✓
	Right to Correction (aka Right to Fix Errors)		✓	✓	✓	✓			✓	✓	✓		✓	✓
	Right to Deletion (aka Right to Be Forgotten)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Right to Opt-Out of Behavioral Advertising		✓	✓	✓	✓	✓	<i>Opt-in with right to withdraw consent</i>	✓	✓	✓	✓ ^{xi}	✓	✓
	Right to Opt-Out of Sale	✓	✓	✓	✓	✓	✓	<i>Opt-in with right to withdraw consent</i>	✓	✓	✓	✓	✓	✓
	Right to Opt-Out of Profiling & Automated Decision Making		<i>To be addressed by regulations</i>	✓	✓	✓			✓	✓	✓		✓	✓
	Right to Limit Use of Sensitive Information ^{xi}		✓					✓	✓	✓	✓	✓		✓
	Right to Appeal			✓	✓	✓	✓			✓	✓	✓	✓	✓
Right to Nondiscrimination	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
Honor opt-out signals (E.g., GPC)	✓ ^{xiii}	✓		✓	✓	✓			✓	✓	✓			
Accountability & Governance		✓ <i>To be addressed by regulations</i>	✓	✓	✓	✓			✓	✓	✓		✓	✓
Security	Appropriate Data Security to Safeguard Information	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Breach Notification	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}	✓ ^{xv}
Transfers to Third Parties	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	

4. New AI Laws



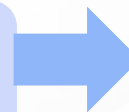
Insurance



CO Protecting Consumers from Unfair Discrimination in Insurance Practices



Employment & Hiring



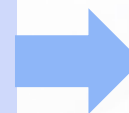
IL Artificial Intelligence Video Interviewing Act



NY Automated Decision Tools law



Investments



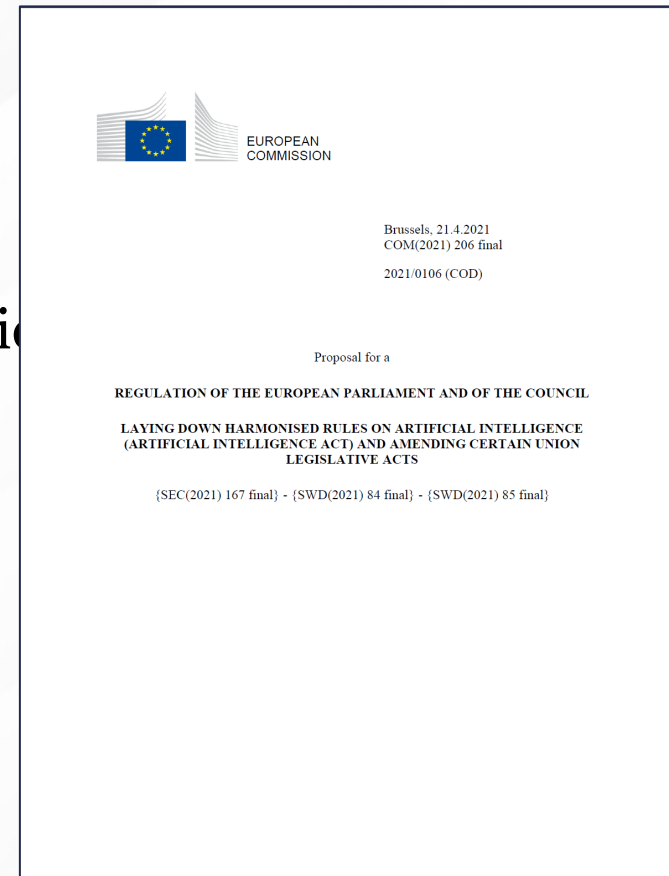
SEC Proposed Rule on Use of Predictive Data Analytics for Broker-Dealers and Investment Advisors

4. AI Act - EU



Massive Scope

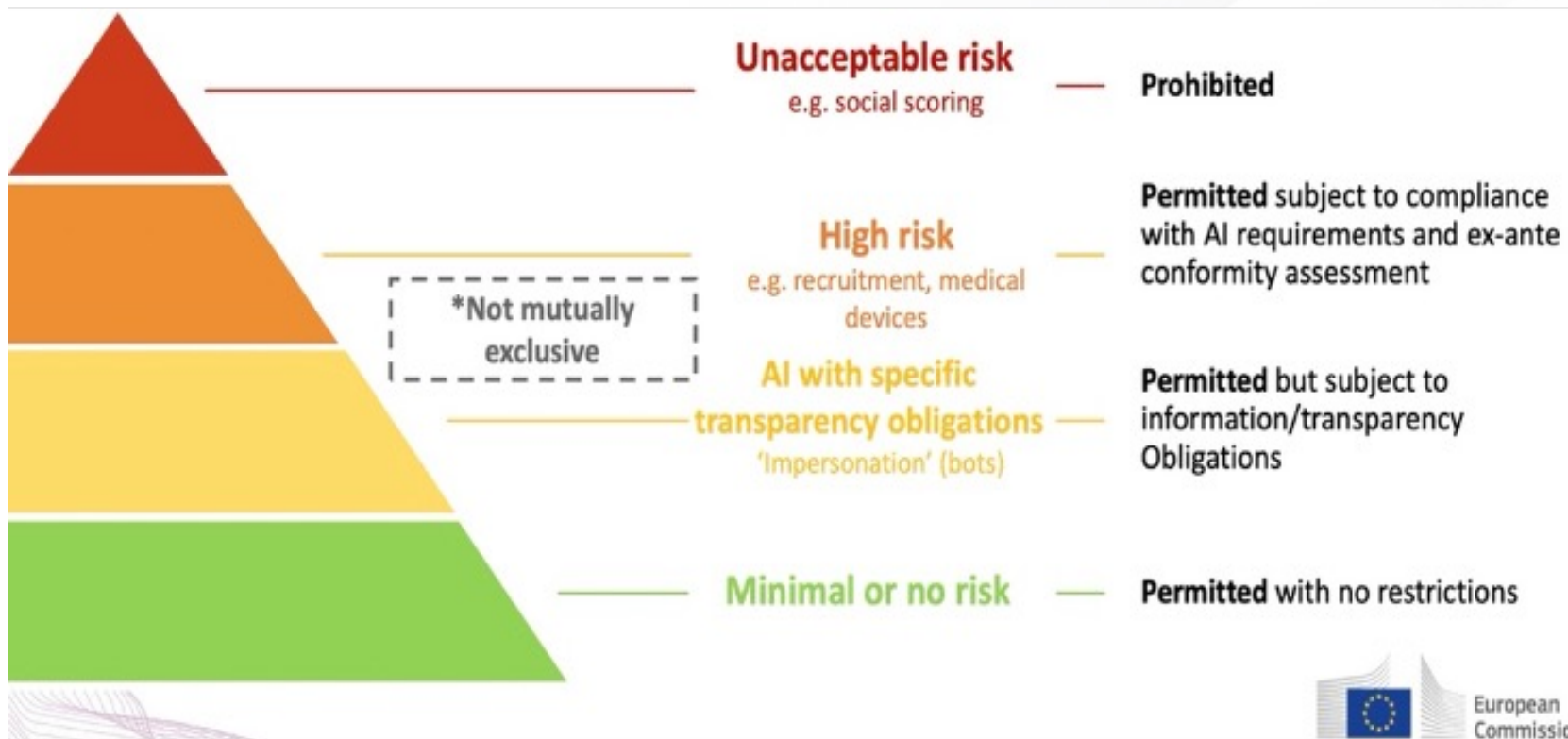
- Art. 2: Jurisdictional Scope ~ GDPR
- Art. 3: Covered Entities:
 - 1. Providers (broadest set of duties)
 - 2. Users
 - 3. Importers
 - 4. Distributors
 - 5. Operators





4. AI Act - EU

Risk-based Approach



4. AI Act - EU



High-risk Systems in Annex III

1. Biometrics
2. Critical infrastructure management
3. Education
4. Employment
5. Essential public/private services (includes financial services)
6. Law enforcement
7. Immigration
8. Justice and democratic process



4. AI Act - EU

Technical Documentation

- Describes the system's design, operation, and performance
- Must be updated regularly
- Must be made available to relevant authorities upon request
- Information on the data used by the AI system
- Potential risks associated with use

4. AI Act - EU

Risk Management

OPENAI / JOBS

Killswitch Engineer

San Francisco, California, United States

\$300,000-\$500,000 per year

About the Role

Listen, we just need someone to stand by the servers all day and unplug them if this thing turns on us. You'll receive extensive training on "the code word" which we will shout if GPT goes off the deep end and starts overthrowing countries.

We expect you to:

- Be patient.
- Know how to unplug things. Bonus points if you can throw a bucket of water on the servers, too. Just in case.
- Be excited about OpenAI's approach to research



- Put system into place to identify, assess, and mitigate potential risks
- Ongoing monitoring and review of the system performance on market
- Human-in-the-loop oversight to ensure safety and ethical use - including fail-safe and emergency shut-off

4. AI Act - EU

Risk Management

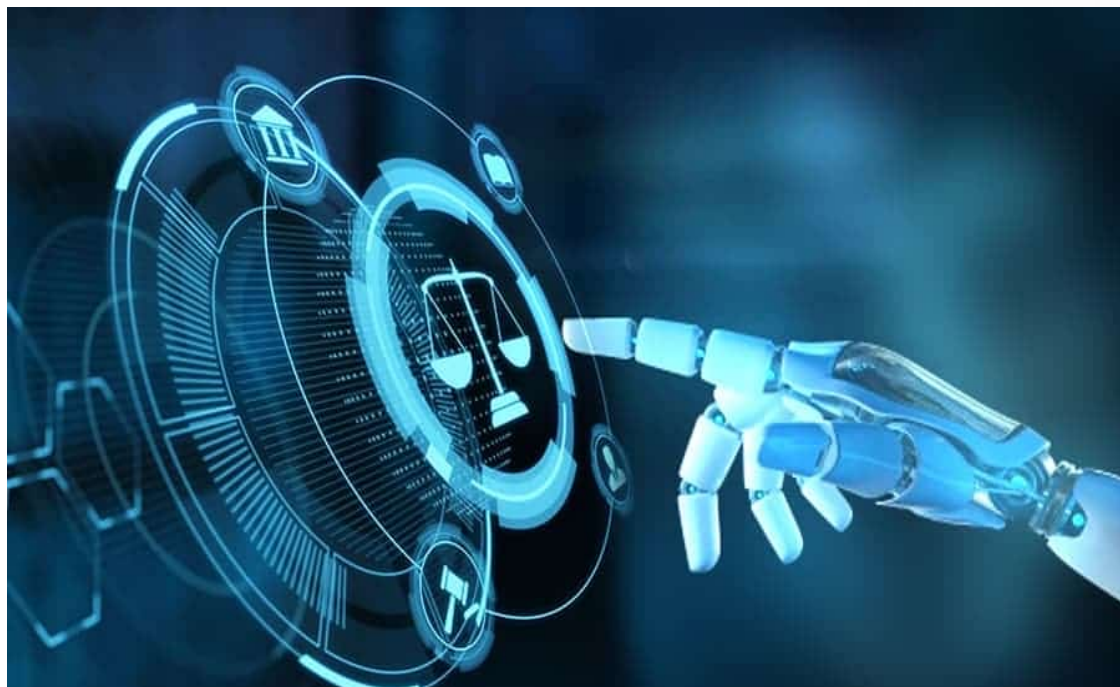


Overview
Four Layers of Defense



4. AI Act - EU

Massive Penalties + AI Police



- Penalties
 - Art. 63 and 71: Penalties even higher than GDPR = 1-7% of annual revenues
 - Art. 65-68: Biggest Hammer = market ban (same as GDPR)
- Enforcers
 - EU Commission budget asking for **600** new enforcement personnel
 - By way of comparison, FTC has around 50