



# Ética e Inteligencia Artificial

Jesús Alfonso López  
[jalopez@uao.edu.co](mailto:jalopez@uao.edu.co)

---



# ¿Qué es la Ética?

Universidad Autónoma de Occidente - Cali



## Ética

1.

Disciplina filosófica que estudia el bien y el mal y sus relaciones con la moral y el comportamiento humano.

2.

Conjunto de costumbres y normas que dirigen o valoran el comportamiento humano en una comunidad.

"su ética profesional le impide confesar más cosas"

<http://cienciauanl.uanl.mx/?p=4943>

# ¿Qué es la Moral?

Universidad Autónoma de Occidente - Cali

Moral

1.

Disciplina filosófica que estudia el comportamiento humano en cuanto al bien y el mal.

2.

Conjunto de costumbres y normas que se consideran buenas para dirigir o juzgar el comportamiento de las personas en una comunidad.

3.

Del comportamiento humano o relacionado con él.  
"la sociedad tiene derecho a exigir que quienes asumen la responsabilidad de la información accedan a los medios de difusión con una preparación intelectual y moral suficientes"

4.

Que se basa en lo que la conciencia establece como bueno.

**#JuntosSomosMásFuertes**



<https://revistaliterariamonolito.com/que-es-la-moral/>

# ¿Son lo Mismo?

Universidad Autónoma de Occidente - Cali



La diferencia entre la ética y la moral es fundamentalmente que la primera es el estudio abstracto y teórico del comportamiento humano, la segunda en cambio, se relaciona más con lo práctico o la acción; es decir, la serie de valores principios y normas que rigen nuestra conducta. La ética se refiere a aquellas reglas internas individuales que adquirimos en el hogar o por ejemplo los principios religiosos y los adoptamos; la moral se refiere a los principios externos impuestos por la sociedad

<https://www.ituser.es/actualidad/2019/02/la-etica-como-base-para-avanzar-hacia-un-mundo-digital-mas-sostenible>

**#JuntosSomosMásFuertes**

<https://www.abc.com.py/edicion-impresa/suplementos/escolar/diferencia-entre-etica-y-moral-1803281.html>



La IA, como cualquier otra tecnología disruptiva, no es neutral en sí misma pues es producto del intelecto humano y sus aplicaciones positivas o negativas dependen de las personas que las usan, en este contexto surge naturalmente la componente ética de las aplicaciones de la IA. El nivel de autonomía o de “inteligencia” del sistema está relacionada con el nivel de responsabilidad que se debe tomar en el proceso de diseño del mismo

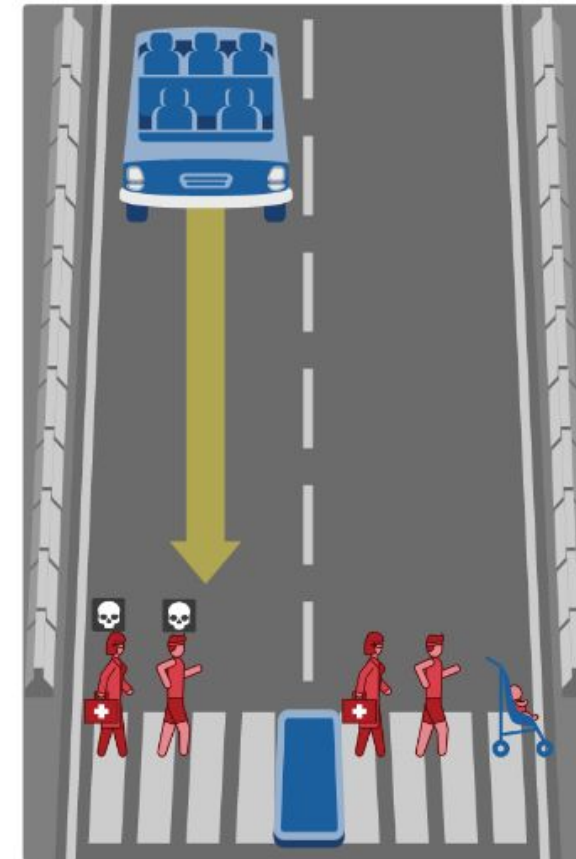
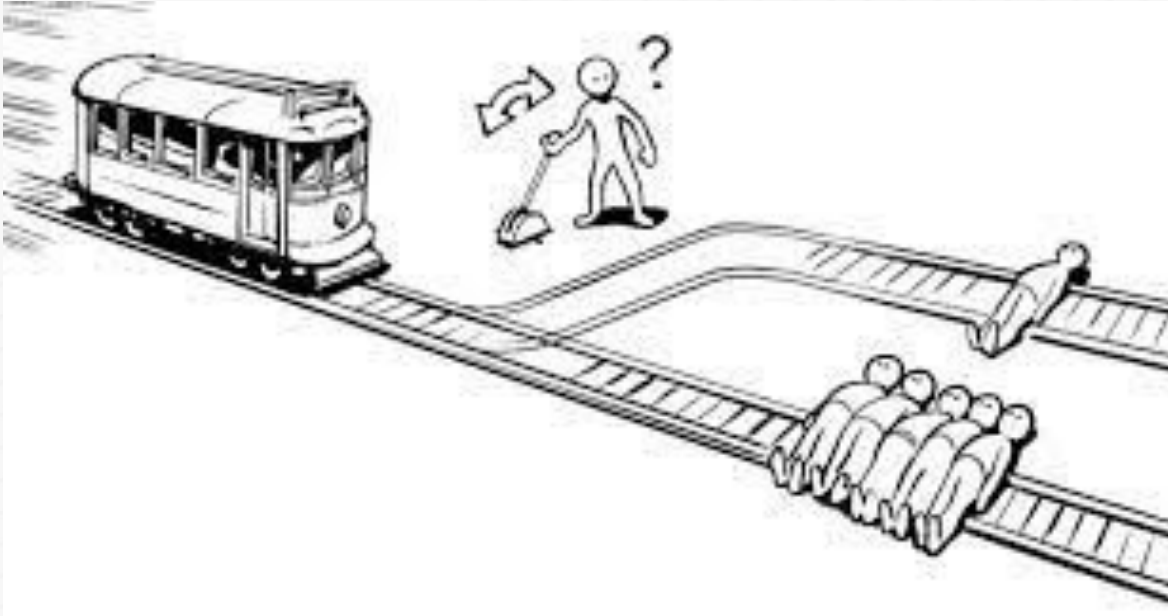


<https://www.oneprod.com/blog/how-does-artificial-intelligence-improve-vibration-analysis/>

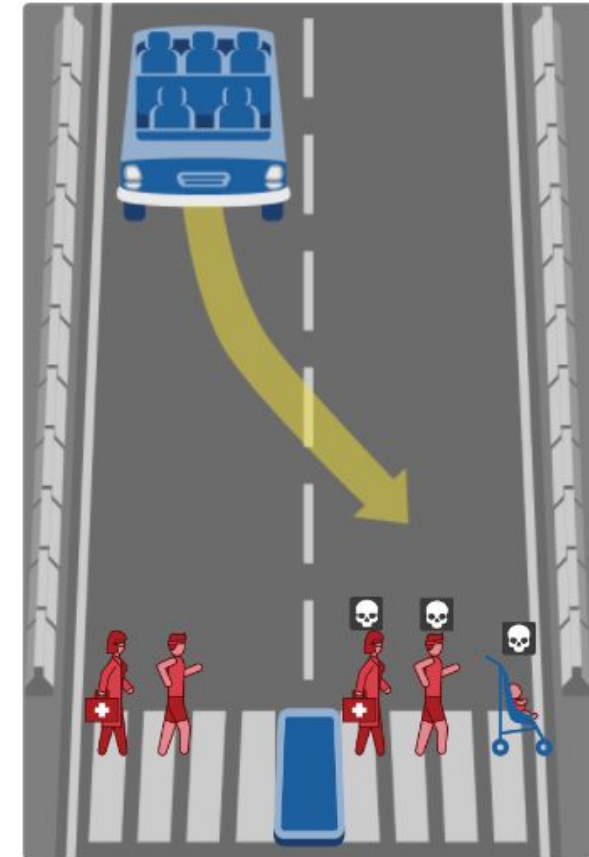
# Ética e Inteligencia Artificial

Universidad Autónoma de Occidente - Cali

What should the self-driving car do?



Show Description



Show Description

<https://es.linkedin.com/pulse/seguridad-vial-dilema-del-tranv%C3%ADa-i%C3%B1aki-barredo>

#JuntosSomosMásFuerres

<https://www.moralmachine.net/>

## Seguridad física

- Mal funcionamiento de vehículos autónomos
- Decisiones de mantenimiento predictivo errado que produce lesiones en los operarios
- Diagnóstico médico errado

## Salud Financiera

- Pobre recomendación financiera evita conseguir préstamos
- Sistema sofisticados de phishing y explotación de información financiera

## Igualdad y Tratamiento Justo

- Discriminación basada en raza , genero u otro prejuicio
- Uso de información de los contacto en redes sociales para definir la calidad de ciudadano



<https://definicion.mx/individuo/>



# Riesgos de la IA para las Organizaciones

Universidad Autónoma de Occidente - Cali

## Desempeño Financiero

- Algoritmos de negociación que no se adaptan a nuevas circunstancias produciendo pérdida económicas
- Decisiones equivocadas que producen estrategias equivocadas en la producción

## Desempeño No-Financiero

- Algoritmos que sin intención producen una fuerza laboral no diversa
- Estimación inadecuada de fondos y recursos ante emergencias

## Cumplimiento Legal

- Discriminación no atendida que producen decisiones que terminan en litigios
- Liberación de datos protegidos



#JuntosSomosMásFuertes

<https://www.mckinsey.com/business-functions/mckinsey-analytics/our-insights/confronting-the-risks-of-artificial-intelligence>

<https://www.volarisgroup.com/blog/article/fostering-an-entrepreneurial-culture-within-your-organization>



# Riesgos de la IA para la Sociedad

Universidad Autónoma de Occidente - Cali

## Seguridad Nacional

- Uso mal intencionado de productos basados en AI (armamento, drones, etc) o uso para actividades ilegales
- Brechas en datos sensibles que exponen secretos militares

## Estabilidad Económica

- Algoritmos automáticos que incrementan la volatilidad de los mercados financieros
- Algoritmos crean inestabilidad en mercados de monedas produciendo un decremento en la transacciones
- Instrumentos financieros caja-negra que producen riesgos sistémicos

## Estabilidad Política

- Manipulación de procesos institucionales (p.e. elecciones) por medio de información falsa



## Rendición de cuentas

¿Quién es responsable de los daños producidos en caso de que alguno de estos dispositivos opere erróneamente o tome una decisión de forma autónoma que resulte en algún tipo de perjuicio?

## Explicabilidad

Lo que en muchas situaciones puede dificultar una clara asignación de daños, perjuicios y responsabilidades es, precisamente, la falta de una explicación clara de por qué un sistema inteligente tomó una decisión en particular



#JuntosSomosMásFu



## Imparcialidad

los sistemas dotados de IA, especialmente aquellos que operan con grandes cantidades de datos, pueden contener en su programación algún tipo de sesgo o prejuicio que los lleve a alcanzar conclusiones parciales o injustas

## Privacidad

Gran parte de las aplicaciones provistas de IA basan su funcionamiento en el acceso a grandes cantidades de información. En este contexto, existe una preocupación ante el uso y la gestión que pueda hacer de esta información, especialmente de la que posee carácter personal.



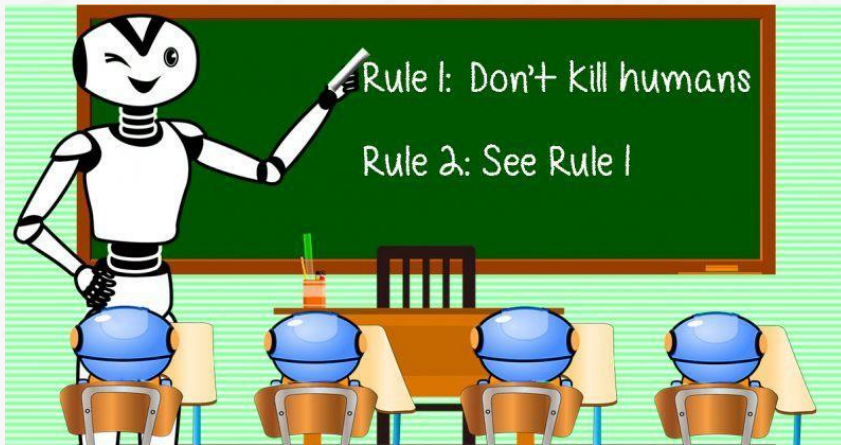
<https://securitytoday.com/articles/2019/03/01/the-flaws-and-dangers-of-facial-recognition.aspx>

# Principios Éticos Para El Diseño Y Desarrollo De La Inteligencia Artificial

Universidad Autónoma de Occidente - Cali

## Respeto de la autonomía humana:

Los sistemas inteligentes deben respetar en todo momento la autonomía y los derechos fundamentales de las personas. Su diseño y programación debe respetar, por tanto, la vida y los derechos humanos sin ningún tipo de discriminación.



<https://towardsdatascience.com/the-evolution-of-robotics-27539c5752fc>

## Transparencia:

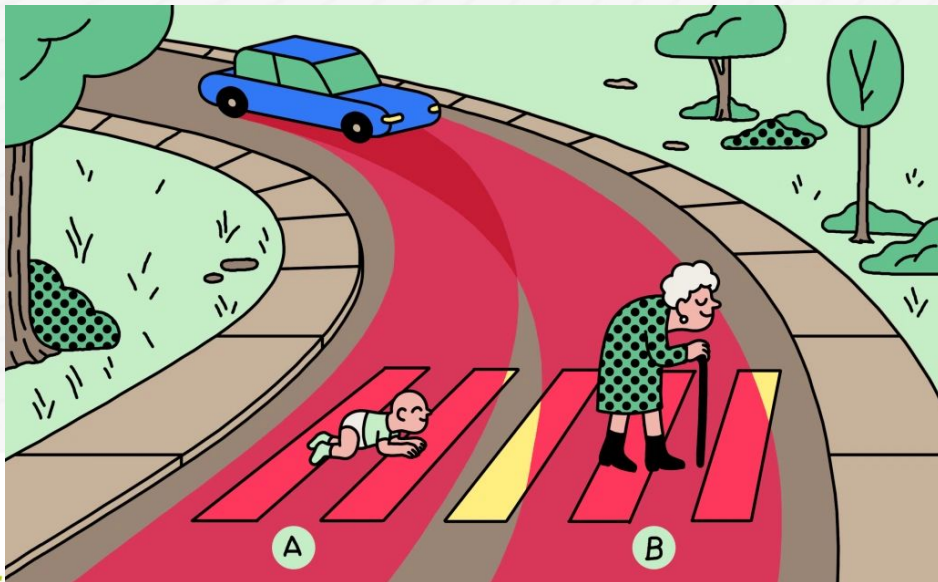
En el caso de los sistemas provistos de IA, la transparencia atañe principalmente a la explicabilidad y la trazabilidad de dichos sistemas. Dado que el diseño de estos dispositivos contempla que tomen decisiones automáticamente con base en distintos cálculos y proyecciones, debe ser posible en todo momento trazar el razonamiento seguido por el sistema y explicar las consecuencias alcanzadas.



# Principios Éticos Para El Diseño Y Desarrollo De La Inteligencia Artificial

Universidad Autónoma de Occidente - Cali

**Responsabilidad y rendición de cuentas:** Estrechamente relacionado con el principio anterior, el diseño y el empleo de sistemas inteligentes deben estar precedidos por una clara asignación de responsabilidades ante los posibles daños y perjuicios que estos puedan ocasionar. La presunta autonomía de estos sistemas no puede servir de pretexto para la dilución de responsabilidades.



## **Robustez y seguridad:**

La fiabilidad de la IA exige que los algoritmos sean suficientemente seguros, fiables y sólidos para operar de manera precisa y segura, y para resolver errores o incoherencias durante todas las fases del ciclo de vida útil de los dispositivos. Este principio exige, además, que los sistemas se diseñen y desarrollen contemplando la posibilidad de ciberataques y fallos técnicos.

<https://www.technologyreview.com/2018/10/24/139313/a-global-ethics-study-aims-to-help-ai-solve-the-self-driving-trolley-problem/>

**Justicia y no discriminación:** el diseño de estos sistemas debe contar con la participación de todos los grupos de interés con los que cada aplicación provista de IA se relacione. Además, estos dispositivos deben garantizar un empleo justo de los datos disponibles para evitar posibles discriminaciones hacia determinados grupos o distorsiones en los precios y en el equilibrio de mercado.



<https://www.sciencemag.org/news/2017/04/even-artificial-intelligence-can-acquire-biases-against-race-and-gender>



# Métodos Técnicos para Implementar IA con Principios Éticos

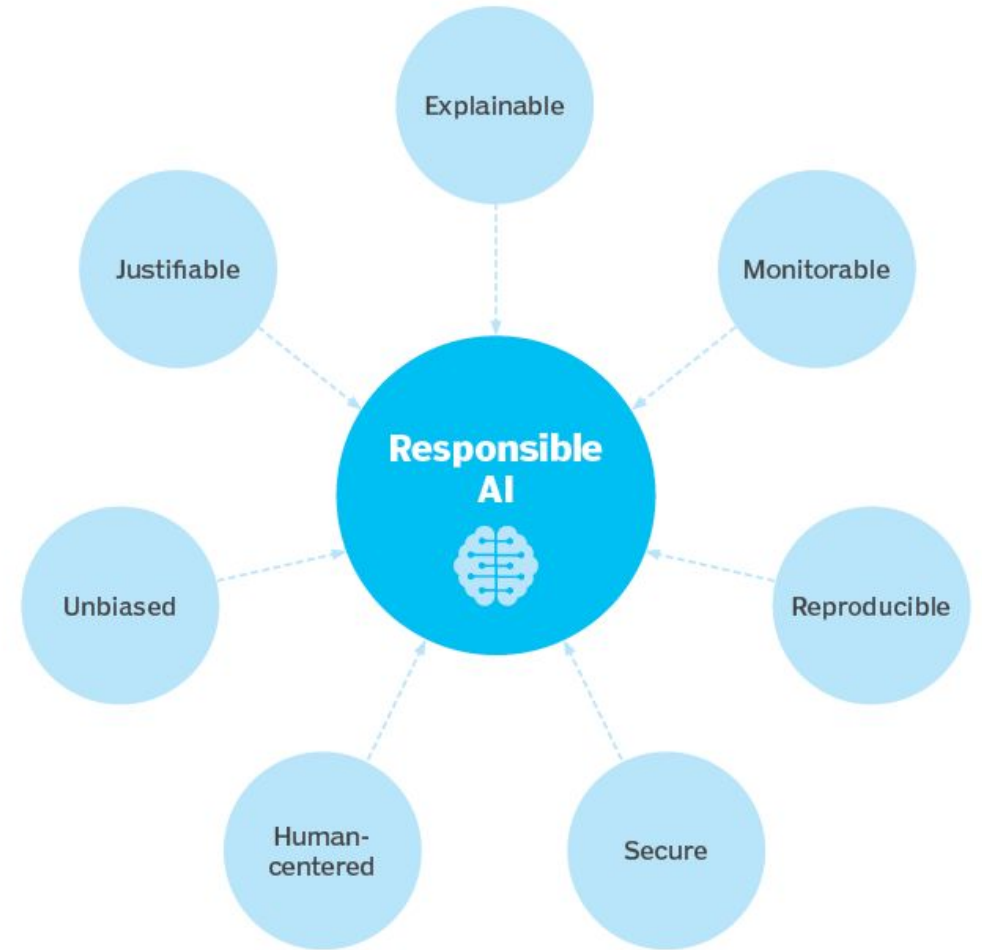
Universidad Autónoma de Occidente - Cali

## Inteligencia Artificial Responsable

La IA responsable es la práctica de diseñar, desarrollar e implementar IA con la buena intención de empoderar a los empleados y las empresas, e impactar de manera justa a los clientes y la sociedad, lo que permite a las empresas generar confianza y escalar la IA con confianza.

<https://www.accenture.com/us-en/services/applied-intelligence/ai-ethics-governance>

#JuntosSomosMásFuertes



<https://www.techtarget.com/searchenterpriseai/definition/responsible-AI>

# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

**AI explicable (XAI):** la capacidad de explicar un modelo después de que se haya desarrollado

**Aprendizaje automático interpretable:** arquitecturas de modelos transparentes y aumento de la intuición y comprensión de los modelos de aprendizaje automático.

**IA ética:** equidad sociológica en las predicciones de aprendizaje automático (es decir, si una categoría de persona se pondera de manera desigual).

**IA segura:** depuración e implementación de modelos ML con contramedidas similares contra amenazas internas y cibernéticas como se vería en el software tradicional

**IA centrada en el ser humano:** interacciones del usuario con sistemas de IA y ML.

**Cumplimiento:** asegurarse de que sus sistemas de IA cumplan con los requisitos reglamentarios pertinentes.

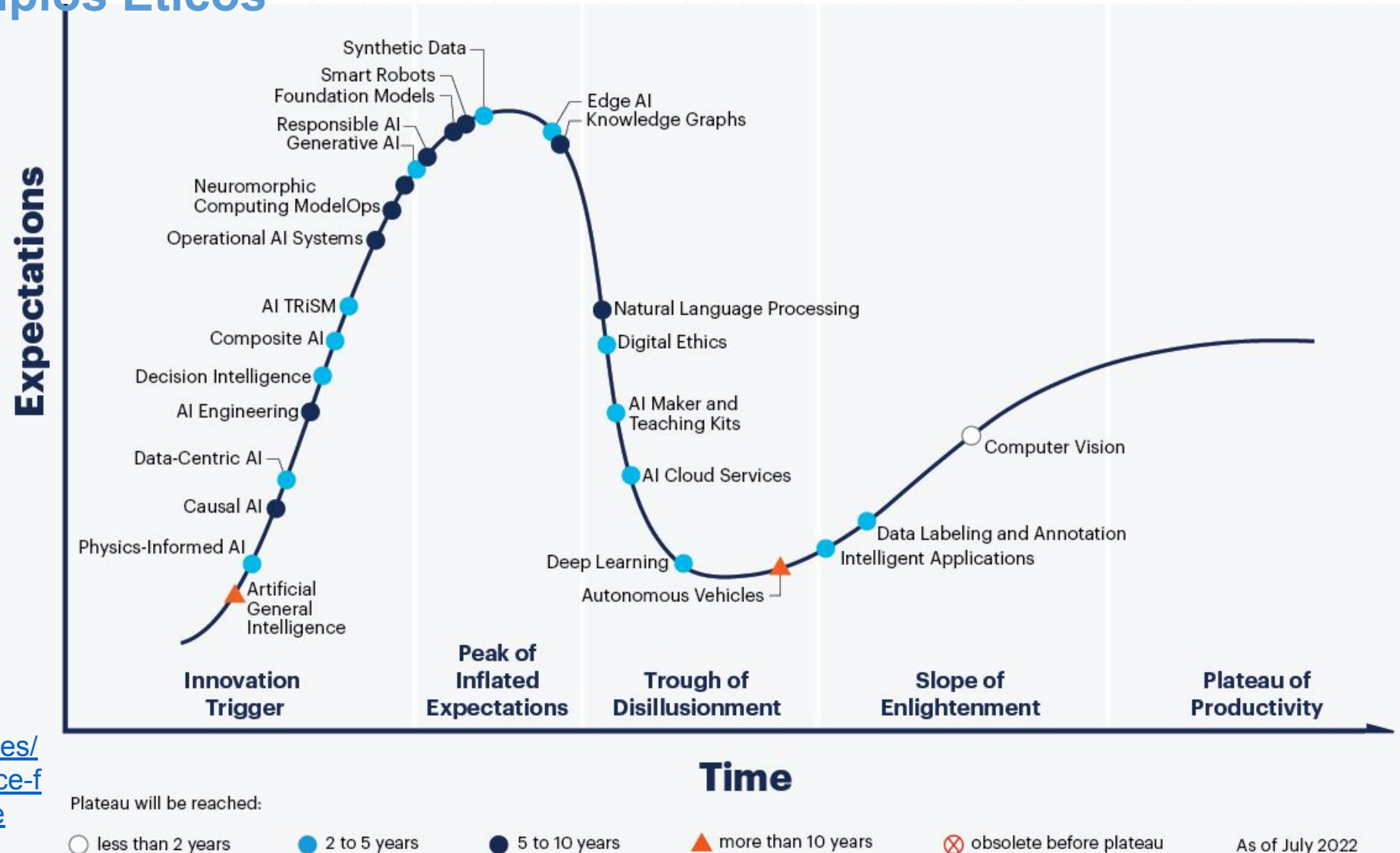


<https://h2o.ai/insights/responsible-ai/>



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali



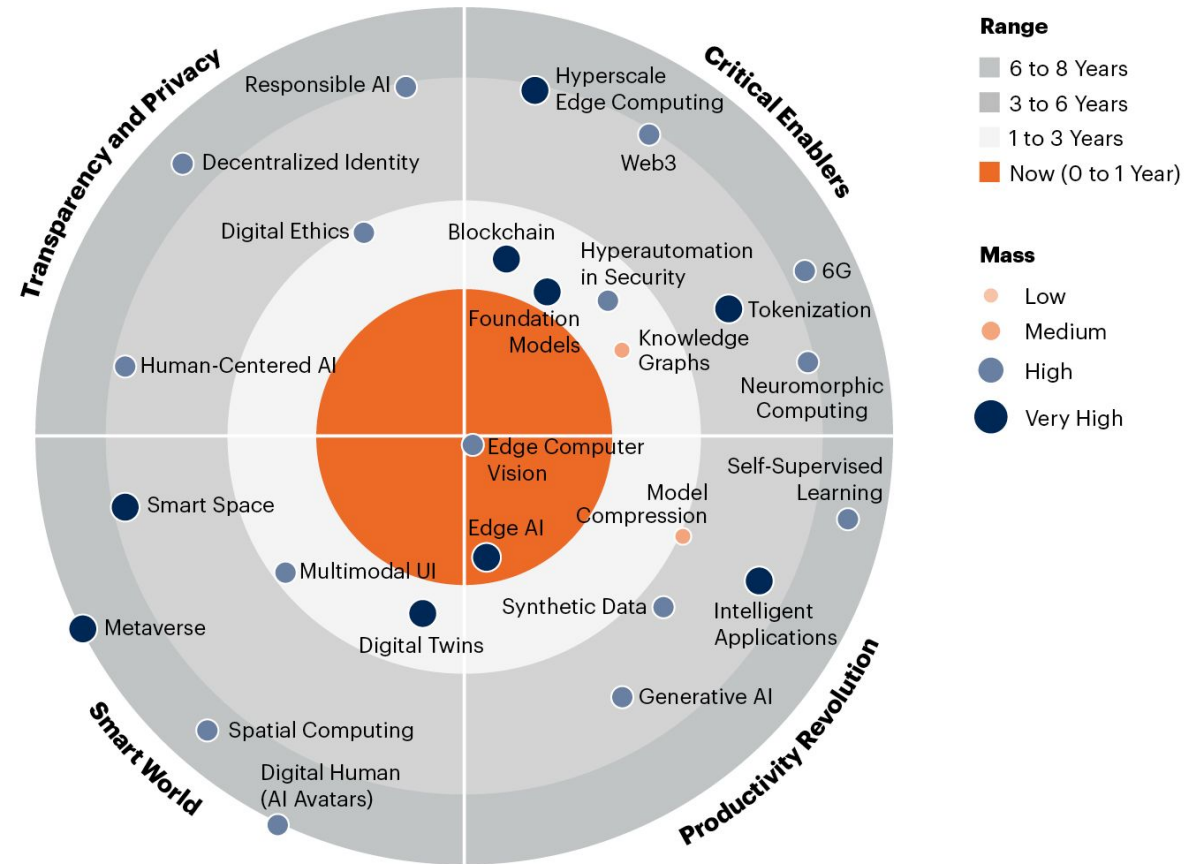
<https://www.gartner.com/en/articles/what-s-new-in-artificial-intelligence-from-the-2022-gartner-hype-cycle>

#JuntosSomosMásFi

# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

## 2023 Gartner Emerging Technologies and Trends Impact Radar



<https://www.gartner.com/en/articles/4-emerging-technologies-you-need-to-know-about>

#JuntosSomosMásFuertes

gartner.com

Note: Range measures number of years it will take the technology/trend to cross over from early adopter to early majority adoption. Mass indicates how substantial the impact of the technology or trend will be on existing products and markets.

Source: Gartner  
© 2023 Gartner, Inc. All rights reserved. CM\_GTS\_2034284

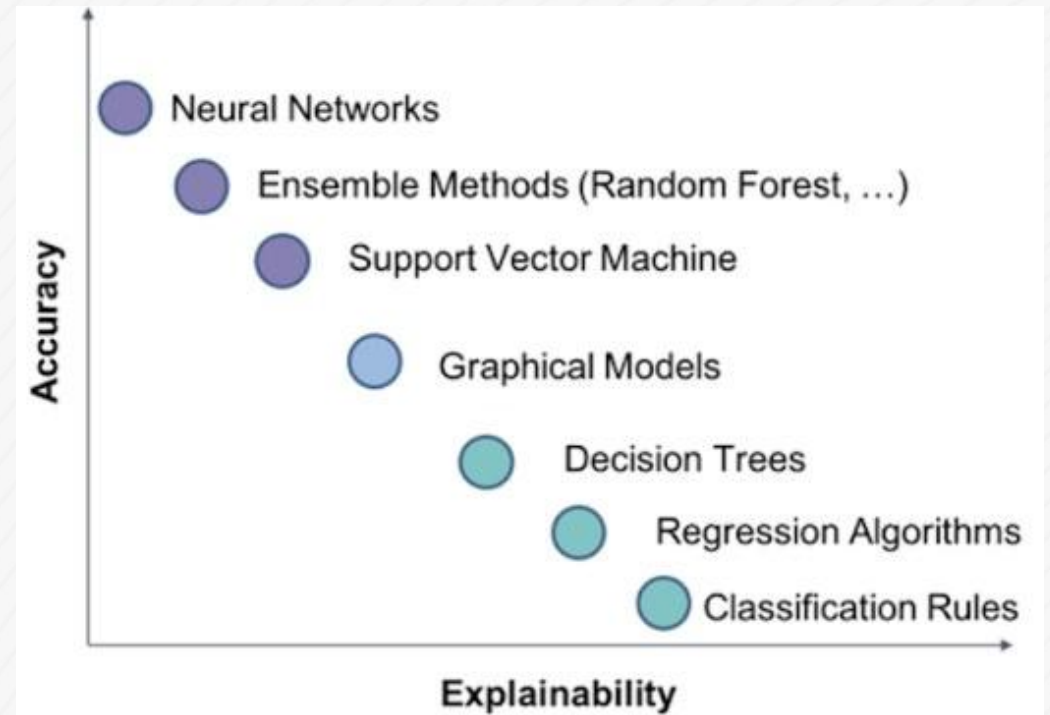
Gartner®



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

**IA explicable:** en los últimos años, ha adquirido también bastante relevancia el campo de investigación denominado XAI (siglas inglesas de *Explainable AI*). En este campo, se han propuesto distintos métodos para convertir muchos de los actuales sistemas de IA en arquitecturas transparentes que dispongan de mecanismos para mostrar de forma clara su funcionamiento y su razonamiento internos.

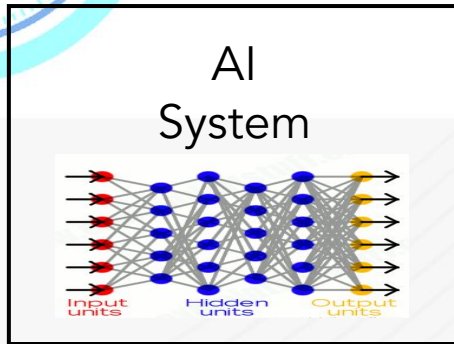


<https://www.kdnuggets.com/2019/01/explainable-ai.html>

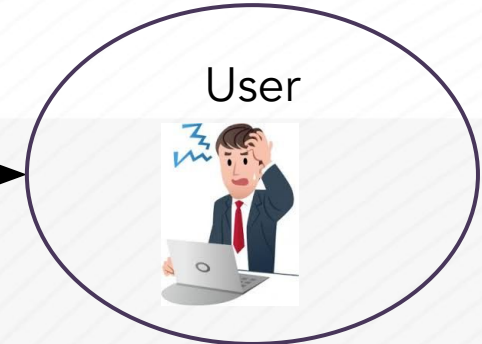
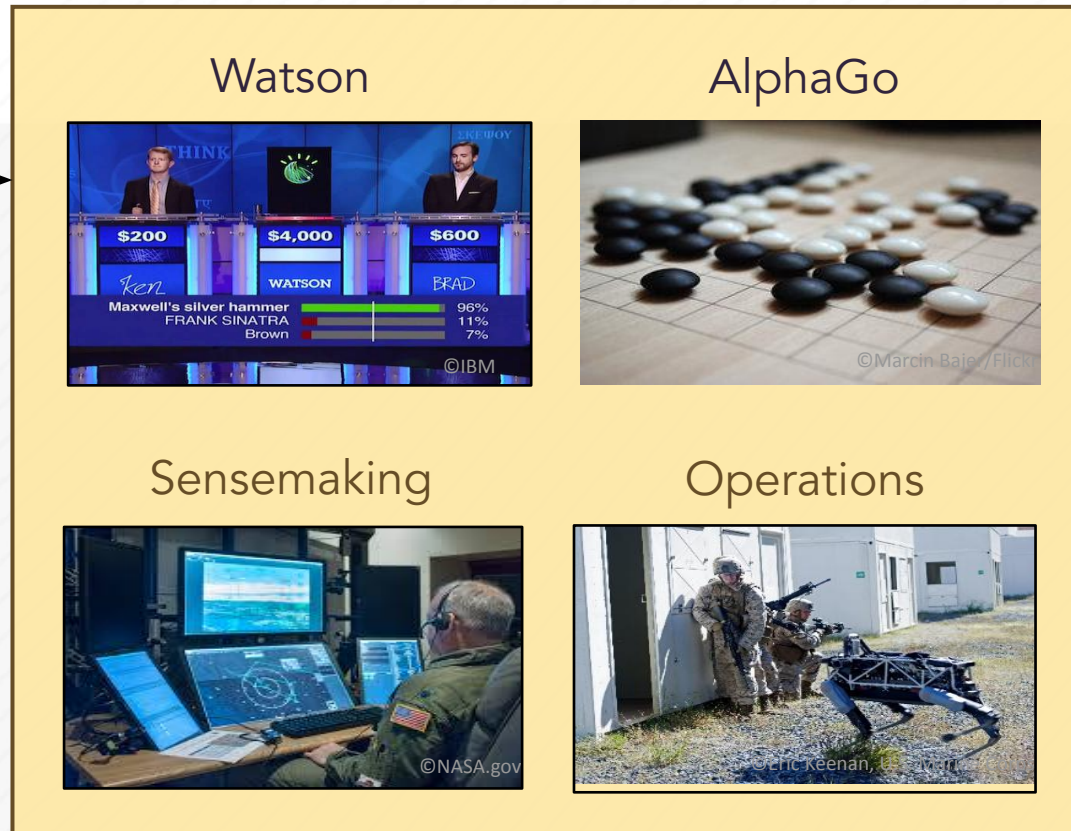
#JuntosSomosMásFuertes

# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali



- We are entering a new age of AI applications
- Machine learning is the core technology
- Machine learning models are opaque, non-intuitive, and difficult for people to understand



- Why did you do that?
- Why not something else?
- When do you succeed?
- When do you fail?
- When can I trust you?
- How do I correct an error?

- The current generation of AI systems offer tremendous benefits, but their effectiveness will be limited by the machine's inability to explain its decisions and actions to users.
- Explainable AI will be essential if users are to understand, appropriately trust, and effectively manage this incoming generation of artificially intelligent partners.

#JuntosSomosMásFuertes

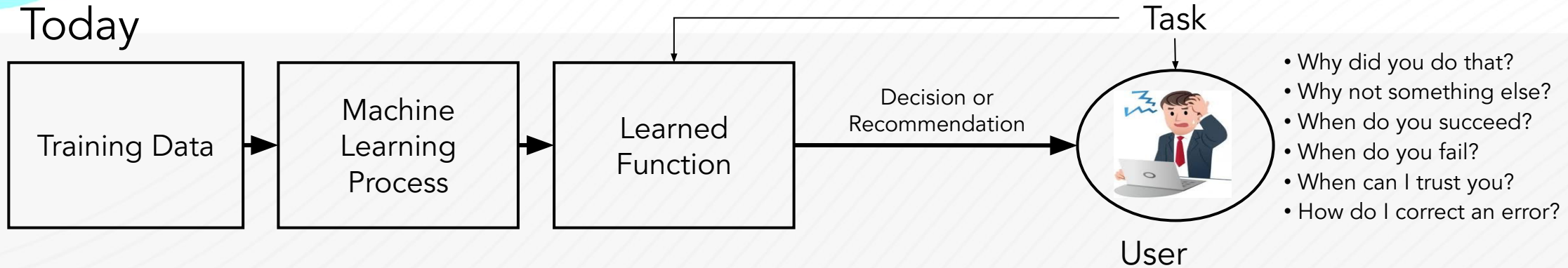
Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)



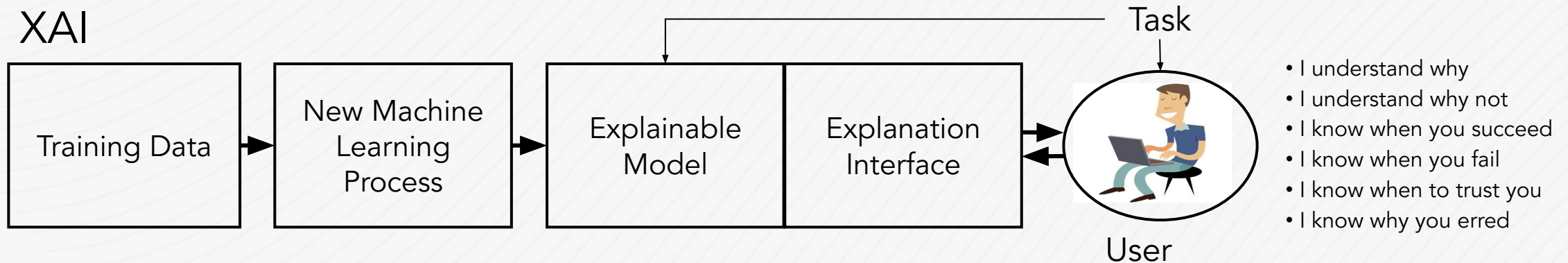
# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

## Today



## XAI



[https://www.darpa.mil/attachments/XAIIndustryDay\\_Final.pptx](https://www.darpa.mil/attachments/XAIIndustryDay_Final.pptx)

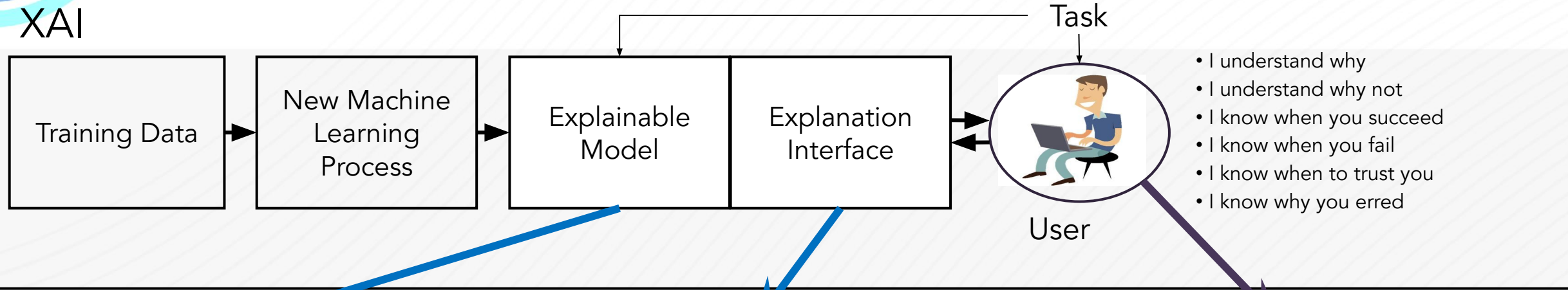
#JuntosSomosMásFuertes

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)

# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

XAI



## Explainable Models

- develop a range of new or modified machine learning techniques to produce more explainable models

## Explanation Interface

- integrate state-of-the-art HCI with new principles, strategies, and techniques to generate effective explanations

## Psychology of Explanation

- summarize, extend, and apply current psychological theories of explanation to develop a computational theory

[https://www.darpa.mil/attachments/XAIIndustryDay\\_Final.pptx](https://www.darpa.mil/attachments/XAIIndustryDay_Final.pptx)

#JuntosSomosMásFuentes

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

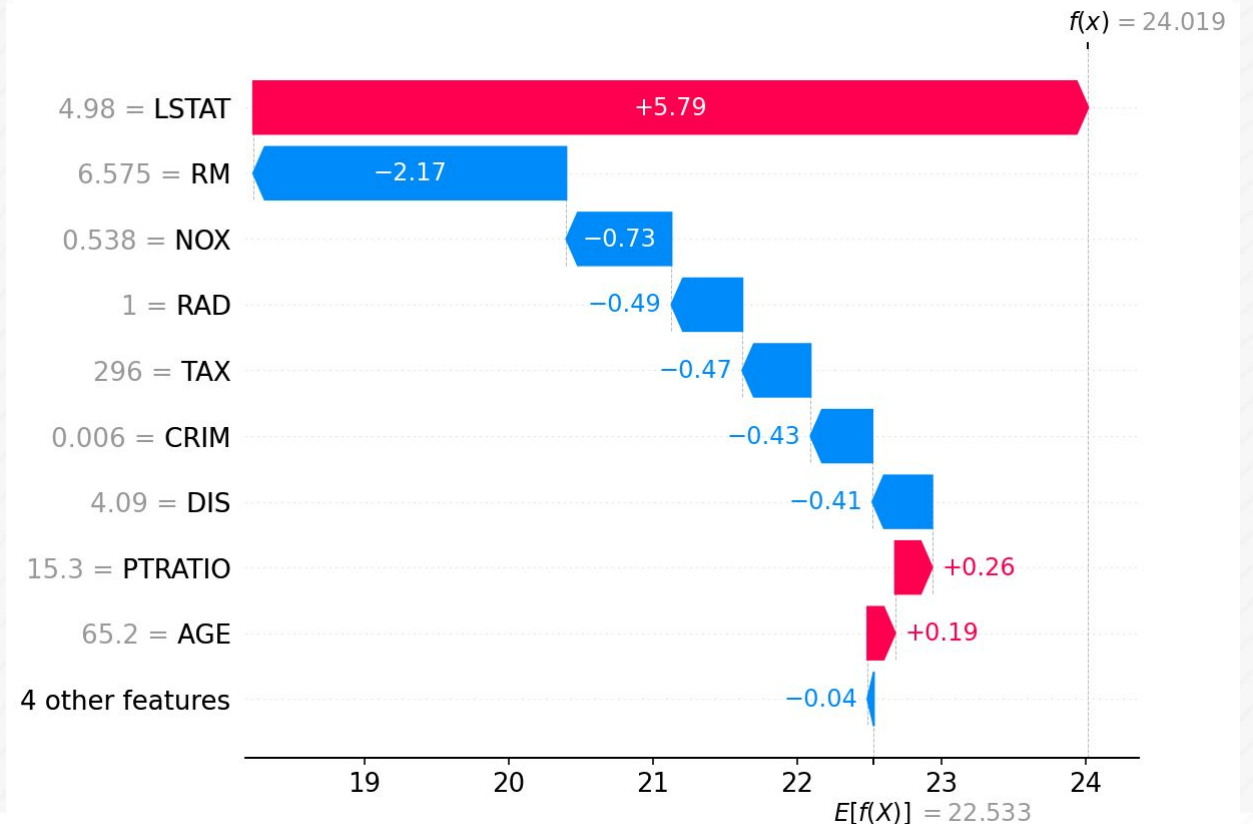
## LIME



<https://github.com/marcotcr/lime>

<https://www.youtube.com/watch?v=hUnRCxnydCc>  
#JuntosSomosMásFuertes

## SHAP



<https://github.com/slundberg/shap>

<https://simmmachines.com/explainable-ai/>

# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

- Se desea una IA explicable en casos de uso que impliquen responsabilidad. Por ejemplo, la IA explicable podría ayudar a crear vehículos autónomos que puedan explicar sus decisiones en caso de accidente.
- La IA explicable es fundamental para situaciones que involucran equidad y transparencia donde hay escenarios con información confidencial o datos asociados con ella (es decir, atención médica)
- Mayor confianza entre humanos y máquinas
- Mayor visibilidad en el proceso de toma de decisiones del modelo (lo que ayuda con la transparencia)



<https://vitalflux.com/what-is-explainable-ai-concept-s-examples/>

#JuntosSomosMásFuertes



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

- La IA explicable es un área de investigación relativamente nueva y todavía hay muchos desafíos activos que presentan los modelos explicables en la actualidad. Un desafío es que la explicabilidad puede ser a expensas de la precisión del rendimiento del modelo, ya que los sistemas de inteligencia artificial explicables tienden a tener un rendimiento más bajo en comparación con los modelos no explicables o los modelos de caja negra.
- Uno de los desafíos clave en la IA explicable es cómo generar explicaciones que sean precisas y comprensibles.
- Otro desafío clave con la inteligencia artificial explicable es que los modelos de IA explicables pueden ser más difíciles de entrenar y ajustar en comparación con los modelos de aprendizaje automático no explicables.
- Otro desafío es que los sistemas de IA explicables pueden ser más difíciles de implementar, ya que las características de explicabilidad generalmente necesitan cierto nivel de intervención humana en el circuito.



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

**Atención médica:** Puede ayudar a los médicos a explicar su diagnóstico a los pacientes y explicar cómo ayudará un plan de tratamiento. Esto ayudará a crear una mayor confianza entre los pacientes y sus médicos al mismo tiempo que mitigará cualquier posible problema ético.

**Fabricación:** la IA explicable podría usarse para explicar por qué una línea de montaje no funciona correctamente y cómo necesita ajustes con el tiempo.

**Defensa:** la IA explicable puede ser útil para aplicaciones de entrenamiento militar para explicar el razonamiento detrás de una decisión tomada por un sistema de inteligencia artificial. Esto es importante porque ayuda a mitigar posibles desafíos éticos, como por qué identificó erróneamente un objeto o no disparó contra un objetivo.

**Vehículos autónomos:** La IA explicable se puede usar para vehículos autónomos donde la explicabilidad proporciona una mayor conciencia situacional en accidentes o situaciones inesperadas, lo que podría conducir a una operación tecnológica más responsable (es decir, prevención de choques).





# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

**Aprobaciones de préstamos:** la inteligencia artificial explicable se puede utilizar para explicar por qué se aprobó o denegó un préstamo. Esto es importante porque ayuda a mitigar cualquier posible desafío ético al proporcionar un mayor nivel de comprensión entre humanos y máquinas

**Selección de currículum:** la inteligencia artificial explicable podría usarse para explicar por qué se seleccionó o no un currículum. Esto proporciona un mayor nivel de comprensión entre humanos y máquinas, pues se mitigan los problemas relacionados con el sesgo y la injusticia.

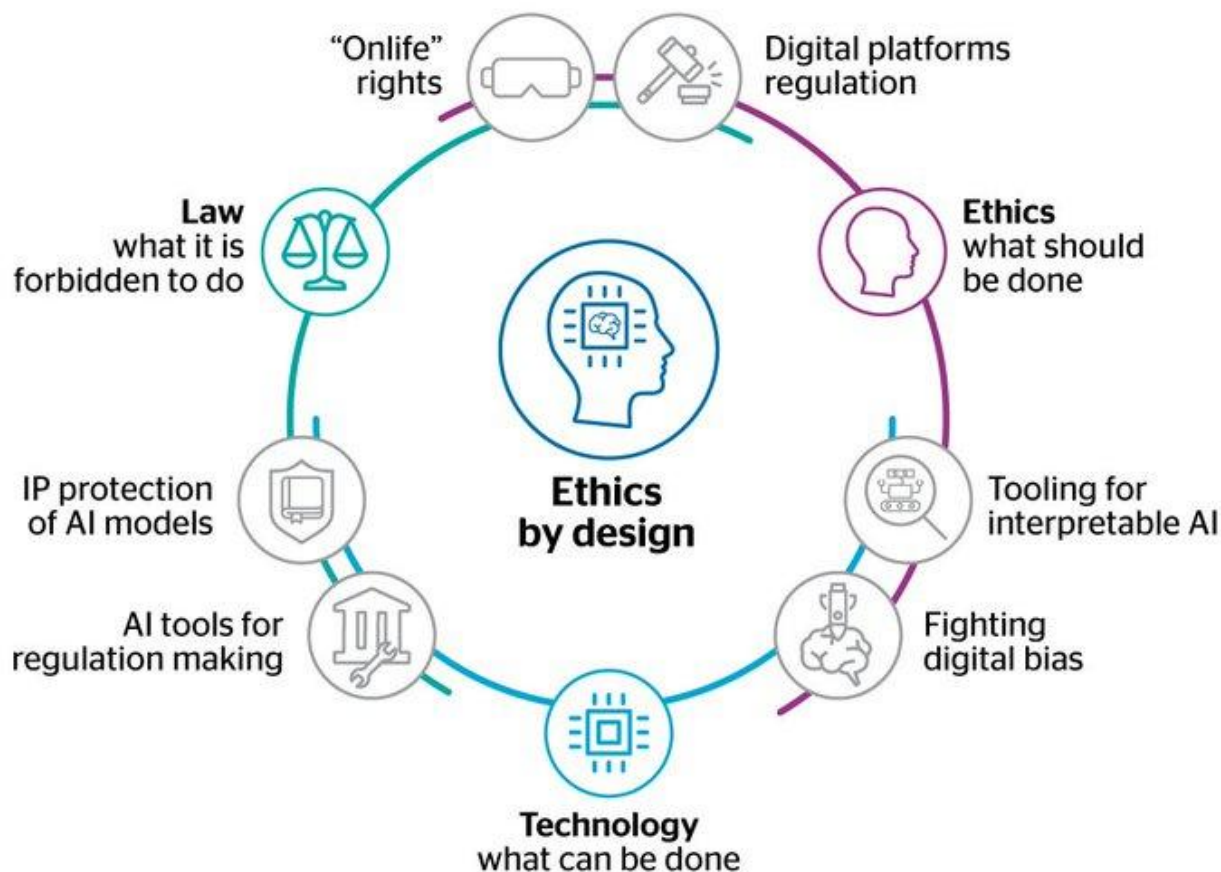
**Detección de fraude:** la IA explicable es importante para la detección de fraude en los servicios financieros. Esto se puede usar para explicar por qué una transacción se marcó como sospechosa o legítima, lo que ayuda a mitigar los posibles desafíos éticos asociados con el sesgo injusto y los problemas de discriminación



# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

***Ethics by design***: se trata de uno de los métodos que más atención ha recibido en los últimos años. Mediante el diseño previo de los algoritmos que controlan los sistemas inteligentes, se podría garantizar el comportamiento ético de estos.



<https://twitter.com/akwyz/status/1065061518423801857>



**Prueba y validación del producto:** para garantizar la seguridad de los dispositivos y asegurar que su diseño y programación responda adecuadamente a lo planeado es someter los distintos dispositivos a mecanismos de examen y validación.

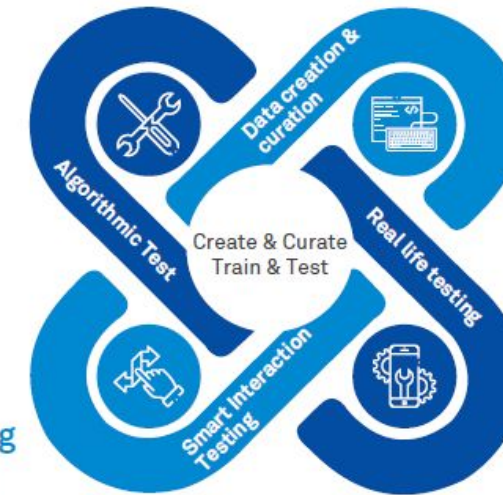
## Algorithms

- Natural Language Processing/Understanding
- Image processing
- Machine learning
- Deep learning

## Smart interaction testing

- Devices (Siri, Alexa, Google Home ...)
- AR/VR
- Drones
- Driverless car
- Robotic arm

## Testing for AI



## Data creation and curation

- Domain specific data
- Cleansing and identifying sample data set
- Contextual data clusters
- Data denoising
- Data labelling

## Real life testing

- Human Unbiased testing
- Challenger Model (algorithm accuracy testing) eg. F-score, confusion matrix for classification algorithms
- Decision analysis (explainable AI)
- Deployment and accessibility testing
- NFR Testing
- Location based User Do testing
- Test triangle (Unit, Service, UI)
- White Box and Black Box testing
- Model Back Testing



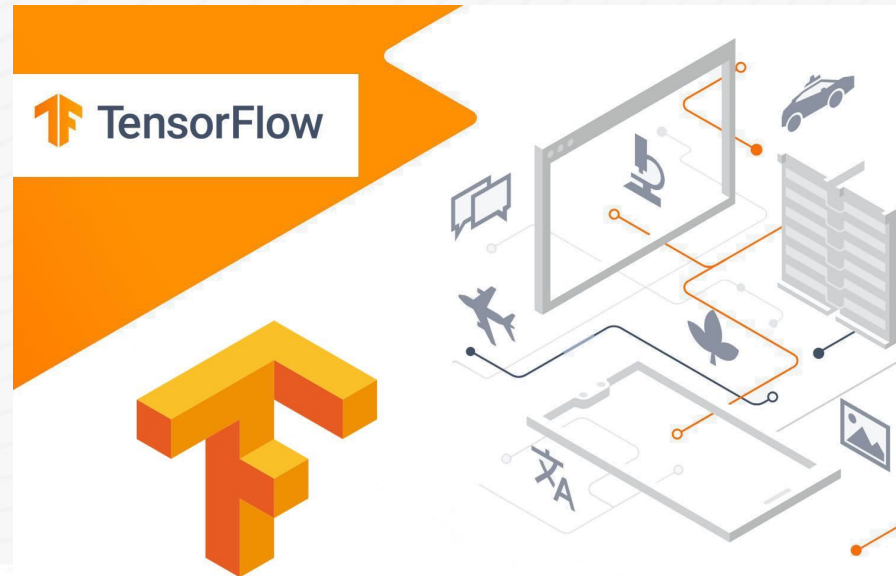
# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali



## Google AI

<https://ai.google/responsibilities/responsible-ai-practices/>



<https://cloud.google.com/explainable-ai>

[https://www.tensorflow.org/responsible\\_ai](https://www.tensorflow.org/responsible_ai)

#JuntosSomosMásFuertes



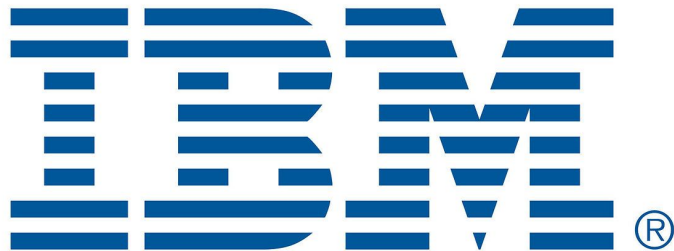
# Métodos Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali



# Microsoft

<https://www.microsoft.com/en-us/ai/responsible-ai>



<https://www.ibm.com/watson/explainable-ai>

#JuntosSomosMásFuertes

# Métodos No-Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

**Regulación:** la labor regulatoria y legislativa puede contribuir a establecer unos parámetros de seguridad y operatividad claros y definidos. Dentro de esta categoría, los Gobiernos y agencias de regulación cuentan con medios como tratados internacionales y resoluciones, procesos de estandarización, directrices y normas no vinculantes, o contratos.

**Certificaciones:** para garantizar la fiabilidad y seguridad de las aplicaciones provistas de IA, y de fomentar, al mismo tiempo, la confianza de los usuarios, es la de emitir certificaciones específicas por parte de las empresas desarrolladoras. Estas certificaciones podrían traducir los distintos estándares en materia de seguridad, transparencia o fiabilidad.

Proposed Regulatory Framework for Modifications to Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD)

<https://www.fda.gov/medical-devices/software-medical-device-samd/artificial-intelligence-and-machine-learning-software-medical-device>

#JuntosSomosMásFuertes





# Métodos No-Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

POLICY AND LEGISLATION | Publication 21 April 2021

## Proposal for a Regulation laying down harmonised rules on artificial intelligence

The Commission has proposed the first ever legal framework on AI, which addresses the risks of AI and positions Europe to play a leading role globally.

<https://digital-strategy.ec.europa.eu/en/library/proposal-regulation-laying-down-harmonised-rules-artificial-intelligence>

<https://www.ai.gov/>



#JuntosSomosMásFuertes

# Métodos No-Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

V  
ADOPCIÓN DE  
PRINCIPIOS  
ÉTICOS PARA LA  
IA EN COLOMBIA

MODELO CONCEPTUAL PARA  
EL DISEÑO DE REGULATORY  
SANDBOXES & BEACHES EN  
INTELIGENCIA ARTIFICIAL

Documento borrador para discusión

**TASK  
FORCE**  
PARA EL DESARROLLO E  
IMPLEMENTACIÓN DE LA  
INTELIGENCIA ARTIFICIAL  
EN COLOMBIA

Documento  
**CONPES**

CONSEJO NACIONAL DE POLÍTICA ECONÓMICA Y SOCIAL  
REPÚBLICA DE COLOMBIA  
DEPARTAMENTO NACIONAL DE PLANEACIÓN

3975

POLÍTICA NACIONAL PARA LA TRANSFORMACIÓN DIGITAL E INTELIGENCIA  
ARTIFICIAL

Documento  
**CONPES**

CONSEJO NACIONAL DE POLÍTICA ECONÓMICA Y SOCIAL  
REPÚBLICA DE COLOMBIA  
DEPARTAMENTO NACIONAL DE PLANEACIÓN

3920

POLÍTICA NACIONAL DE EXPLOTACIÓN DE DATOS  
(BIG DATA)

<https://dapre.presidencia.gov.co/TD/MARCO-ETICO-PARA-LA-INTELIGENCIA-ARTIFICIAL-EN-COLOMBIA.pdf> (Agosto, 2020)

<https://dapre.presidencia.gov.co/AtencionCiudadana/Documents/TASK-FORCE-para-desarrollo-implementacion-Colombia-propuesta-201120.pdf>

<https://dapre.presidencia.gov.co/AtencionCiudadana/DocumentosConsulta/consulta-200820-MODELO-CONCEPTUAL-DISENO-REGULATORY-SANDBOXES-BEACHES-IA.pdf>

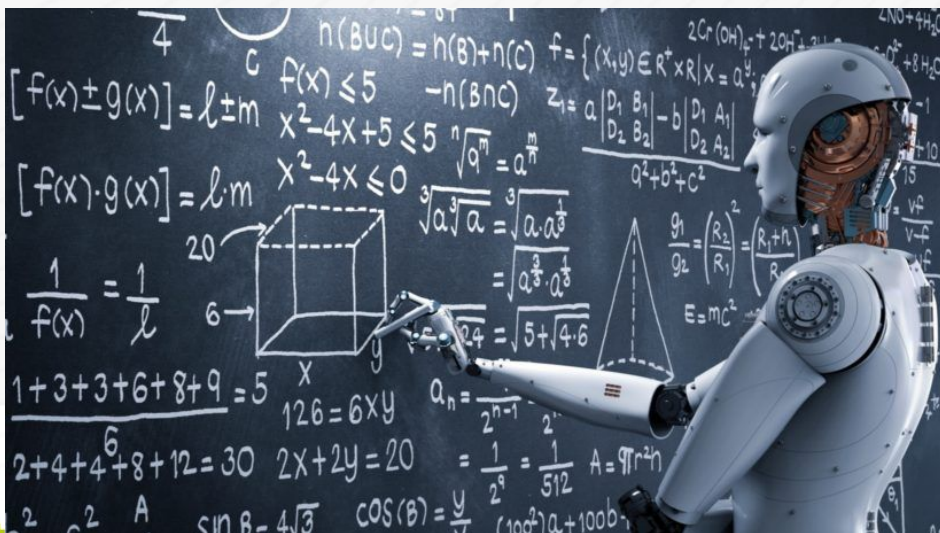


# Métodos No-Técnicos para Implementar IA con Principios Éticos

Universidad Autónoma de Occidente - Cali

## Educación y sensibilización:

La educación y la comunicación en materia de IA puede contribuir a crear una mayor conciencia en torno a los potenciales riesgos que esta tecnología entraña. Esta labor ha de alcanzar a todos los grupos de interés (diseñadores, consumidores —ya sean individuos o empresas—, reguladores, etc.)



**Investigación:** para garantizar que la IA se siga desarrollando de manera fiable y segura hay que asegurar que la ética y el buen gobierno acompañen siempre a los temas de investigación en IA. Esto puede alcanzarse dando prioridad a estos temas de investigación en la asignación de presupuestos o incentivando el trabajo de grupos y centros de investigación que analicen qué desafíos plantea la IA a la ética, al gobierno y a la responsabilidad social de las empresas.

# Enfrentando los Riesgos: Principios de Asilomar

Universidad Autónoma de Occidente - Cali

## Research Issues

- 1) **Research Goal:** The goal of AI research should be to create not undirected intelligence, but beneficial intelligence.
- 2) **Research Funding:** Investments in AI should be accompanied by funding for research on ensuring its beneficial use, including thorny questions in computer science, economics, law, ethics, and social studies, such as:
  - How can we make future AI systems highly robust, so that they do what we want without malfunctioning or getting hacked?
  - How can we grow our prosperity through automation while maintaining people's resources and purpose?
  - How can we update our legal systems to be more fair and efficient, to keep pace with AI, and to manage the risks associated with AI?
  - What set of values should AI be aligned with, and what legal and ethical status should it have?
- 3) **Science-Policy Link:** There should be constructive and healthy exchange between AI researchers and policy-makers.
- 4) **Research Culture:** A culture of cooperation, trust, and transparency should be fostered among researchers and developers of AI.
- 5) **Race Avoidance:** Teams developing AI systems should actively cooperate to avoid corner-cutting on safety standards.

<https://futureoflife.org/ai-principles/>



# Enfrentando los Riesgos: Principios de Asilomar

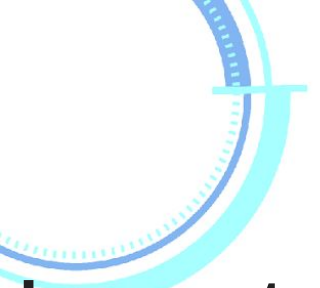
Universidad Autónoma de Occidente - Cali

## Ethics and Values

- 6) **Safety:** AI systems should be safe and secure throughout their operational lifetime, and verifiably so where applicable and feasible.
- 7) **Failure Transparency:** If an AI system causes harm, it should be possible to ascertain why.
- 8) **Judicial Transparency:** Any involvement by an autonomous system in judicial decision-making should provide a satisfactory explanation auditable by a competent human authority.
- 9) **Responsibility:** Designers and builders of advanced AI systems are stakeholders in the moral implications of their use, misuse, and actions, with a responsibility and opportunity to shape those implications.
- 10) **Value Alignment:** Highly autonomous AI systems should be designed so that their goals and behaviors can be assured to align with human values throughout their operation.
- 11) **Human Values:** AI systems should be designed and operated so as to be compatible with ideals of human dignity, rights, freedoms, and cultural diversity.
- 12) **Personal Privacy:** People should have the right to access, manage and control the data they generate, given AI systems' power to analyze and utilize that data.
- 13) **Liberty and Privacy:** The application of AI to personal data must not unreasonably curtail people's real or perceived liberty.
- 14) **Shared Benefit:** AI technologies should benefit and empower as many people as possible.
- 15) **Shared Prosperity:** The economic prosperity created by AI should be shared broadly, to benefit all of humanity.
- 16) **Human Control:** Humans should choose how and whether to delegate decisions to AI systems, to accomplish human-chosen objectives.
- 17) **Non-subversion:** The power conferred by control of highly advanced AI systems should respect and improve, rather than subvert, the social and civic processes on which the health of society depends.
- 18) **AI Arms Race:** An arms race in lethal autonomous weapons should be avoided.

#JuntosSomosMásFuertes

<https://futureoflife.org/ai-principles/>



# Enfrentando los Riesgos: Principios de Asilomar

Universidad Autónoma de Occidente - Cali

## Longer-term Issues

- 19) **Capability Caution:** There being no consensus, we should avoid strong assumptions regarding upper limits on future AI capabilities.
- 20) **Importance:** Advanced AI could represent a profound change in the history of life on Earth, and should be planned for and managed with commensurate care and resources.
- 21) **Risks:** Risks posed by AI systems, especially catastrophic or existential risks, must be subject to planning and mitigation efforts commensurate with their expected impact.
- 22) **Recursive Self-Improvement:** AI systems designed to recursively self-improve or self-replicate in a manner that could lead to rapidly increasing quality or quantity must be subject to strict safety and control measures.
- 23) **Common Good:** Superintelligence should only be developed in the service of widely shared ethical ideals, and for the benefit of all humanity rather than one state or organization.



## Reflexión Final

Todo aquel que forma parte del ecosistema de la Inteligencia Artificial: los empleados de las grandes compañías de tecnología, los gerentes, los líderes y los miembros de las juntas directivas, las startups, los inversionistas, los profesores y estudiantes de posgrado (y de pregrado), así como cualquier otra persona que trabaje en Inteligencia Artificial, debe reconocer que está tomando decisiones éticas todo el tiempo, todos deben estar preparados para explicar las decisiones que han tomado durante las fases de desarrollo, prueba y despliegue (de los modelos de Inteligencia artificial)

Ammy Webb ( Los Nueve Gigantes)





**Gracias**

**#JuntosSomosMásFuertes**

