



SESIÓN CIENTÍFICA
PUERTAS ABIERTAS
ACADEMIA DE CIENCIAS DE CUBA

Lunes 13 de
Noviembre,
10:00 AM



Dr.C. Francisco Herrera Triguero
Académico Correspondiente
Academia de Ciencias de Cuba
Director del Instituto Interuniversitario Andaluz en Data
Science and Computational Intelligence de las Universidades
de Granada y de Jaén

Título: "Una Visión Holística de la inteligencia artificial.
Segura y fiable: Riesgos, Ética, Regulación y auditabilidad"

Lugar: Calle Cuba No. 460, entre Amargura y Teniente Rey, La Habana Vieja.

INTELIGENCIA ARTIFICIAL

La inteligencia artificial silenciosa ha entrado en nuestras vidas

Francisco Herrera



Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

IA fiable y segura

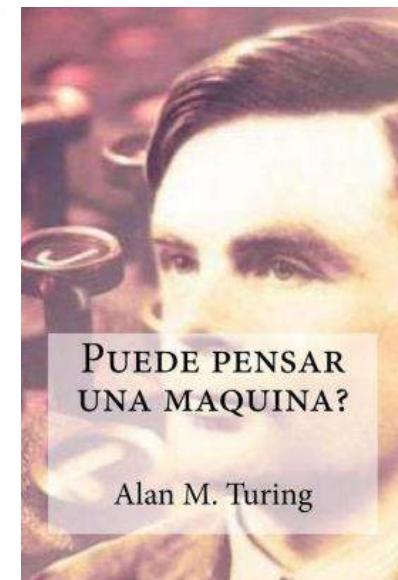
Concluyendo

Inteligencia artificial



John McCarthy (1955) (Stanford)
Conferencia de Dartmouth (1956)

Inteligencia Artificial: "la ciencia e ingeniería de hacer máquinas que se comporten de una forma que llamaríamos inteligente si el humano tuviese ese comportamiento"



Alan Turing
Pionero de la IA
Test de Turing

Inteligencia artificial silenciosa, ha entrado en nuestras vidas



1996 | Primer robot aspirador:
Electrolux Trilobite

Diagnóstico
de enfermedades

Asistentes
de navegación / rutas

Automoción

Sistema de
Recomendaciones

amazon
NETFLIX

Asistentes virtuales: Siris,
Alexa, google Assistant, ...

Inteligencia artificial silenciosa, ha entrado en nuestras vidas

La inteligencia artificial ha entrado en el salón de casa.

Radio años 20 (siglo XX)

Televisión años 50 (siglo XX)

IA años 20 (siglo XXI) (conversaciones, recomendaciones, ...)

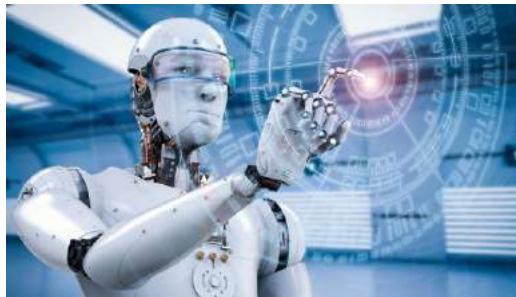
Breve recorrido histórico. Inteligencia artificial. 70 años de historia

"Máquinas no pensantes cada vez más capaces"

Edad de ORO

1960 - Años sesenta

Grandes Expectativas

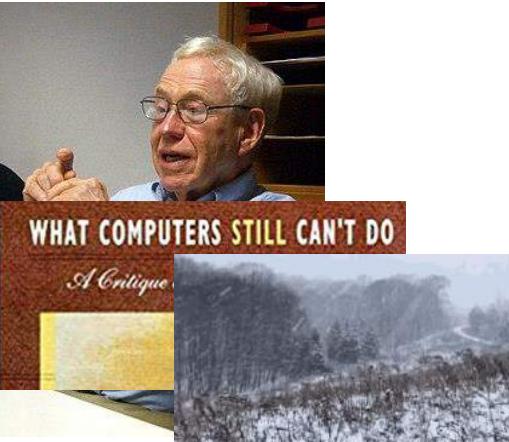


HAL 9000

2001: Una odisea
del espacio
Stanley Kubrick
1968

Invierno de la IA

1970 - Años setenta



Hubert Dreyfus

1980 – IA débil (propósito específico)

La inteligencia artificial empieza a rendir y asombrar



1996 – 1997

Kasparov vs. Deep Blue

Ajedrez

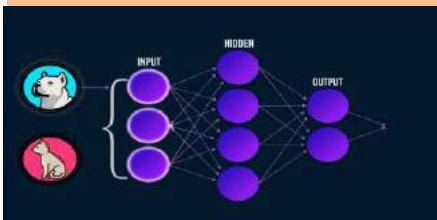
Coche Autónomo



2005 - Stanley
(Standord)
DARPA
grand
Challenge
(S. Thrun)

Desierto de Mojave,
212 kilómetros, 6 horas y 54 minutos

2001 – ... Inteligencia Artificial entra en nuestras vidas. Era del Big Data (macrodatos). Deep Learning. El poder de los datos

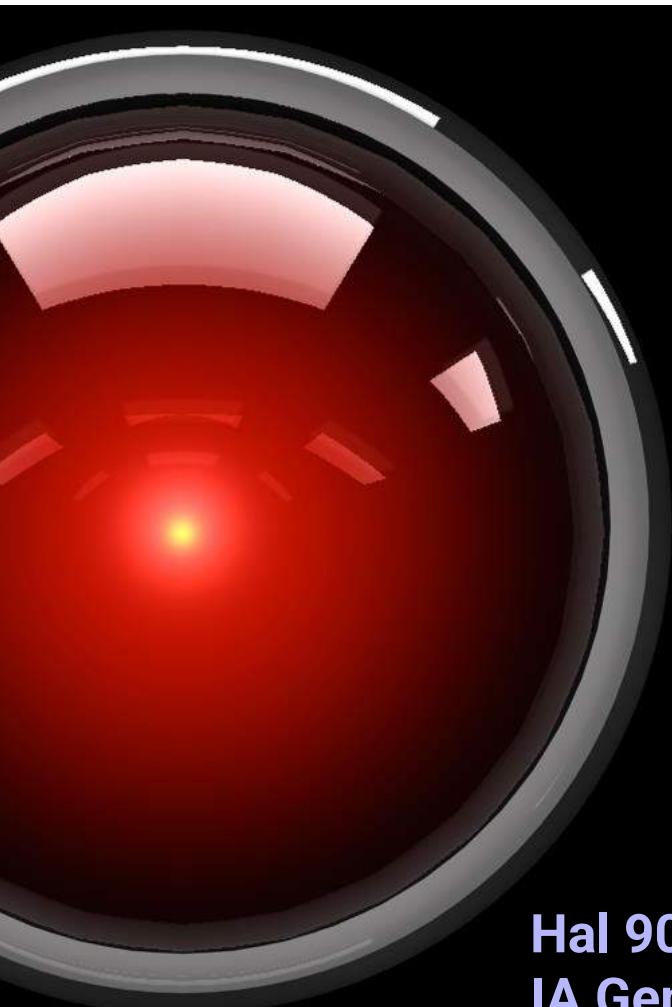


2022 - Era de la IA Generativa



Inteligencia artificial. Invierno de la IA

Expectativas no cumplidas en los 60, “IA General o fuerte”, grandes retos, asistentes conversacionales y sistemas inteligentes al nivel de los humanos, ej. Hal 9000



**Hal 9000: ejemplo de
IA General o IA Fuerte**

El invierno de la IA

A collage of images related to the "winter of AI". At the top left is a portrait of Marvin Minsky. To his right is the title card for the movie "2001: A Space Odyssey" with the tagline "a space odyssey". Below the title card is a red abstract graphic with diagonal stripes. The overall theme is the disappointment of AI expectations in the 1960s.

Hal 9000

2001: Una odisea del espacio (Stanley Kubrick, 1968)

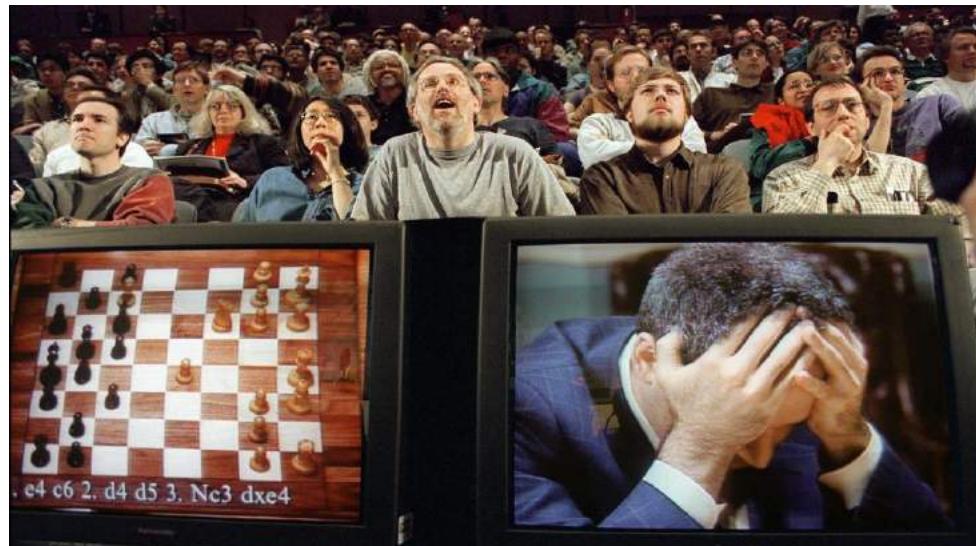
Asesor: Marvin Minsky

Inteligencia artificial. Impulso (IA de propósito específico, IA débil)
Años 80 y 90, impulso en base a abordar problemas complejos concretos,
IA débil. Gran hito ajedrez 96-97 (DeepBlue, ej. de “IA débil”, solo juega al ajedrez)

El día que los microchips vencieron a las neuronas:

10/2/1996

Deep Blue (una máquina que analizaba 100 millones de jugadas por segundo) derrotó al entonces campeón mundial en la primera partida del match; el humano se repuso y ganó el duelo. Cómo recuerda aquel encuentro el ajedrecista nacido en Bakú



■ Ajedrez

La máquina despiadada

'Deep Blue' vence a Kasparov en diecinueve movimientos

EFE • NUEVA YORK

El ordenador 'Deep Blue', que ayer conducía las piezas blancas, derrotó y humilló al campeón del mundo de la Asociación Profesional de Ajedrez (PCA), el ruso Gari Kasparov, en 19 movimientos, en la sexta y última partida del encuentro que han disputado en Nueva York.

El estupor fue general. Público, aficionados, maestros, y seguidores de la partida a través de Internet en todo el mundo, no podían dar crédito a lo que estaba pasando en la sexta partida del encuentro entre la máquina y el hombre. La partida comenzó con la variante Nimzowitch de la antigua defensa Caro Kan, desarrollada por los mues-

tos Caro y Kan a finales del siglo pasado y que es una de las preferidas de Anatoli Karpov y otros jugadores de élite.

Sorpresa

La sorpresa llegó en la octava jugada con un sacrificio de la máquina que, aunque es teórico, no se suele practicar en el ajedrez de alto

nivel por los riesgos que conlleva de partida abierta y de difícil cálculo. Pero la segunda sorpresa, aún mayor, llegó en el movimiento 17 cuando Kasparov se dejó la dama a cambio de una torre y un alfil pero lo peor es que no había contrajuego debido a la mala situación del rey de Kasparov en el centro del tablero. La jugada 19 de 'Deep Blue' -c4- dejó al campeón del mundo sin esperanza alguna, ni siquiera de aspirar a las tablas. Una derrota sorprendentemente el mejor del mundo. ¿Error humano, cansancio, 'stress', demasiada confianza en sí mismo?



Kasparov muestra su desesperación frente al ordenador.

12/5/1997

Inteligencia artificial. Cambio de paradigma

“Cambio de paradigma, desde algoritmo al dato. La era del big data (macrodatos)”

2005 - Driverless car



Stanley (Standord), DARPA grand Challenge - 2005 (S. Thrun)



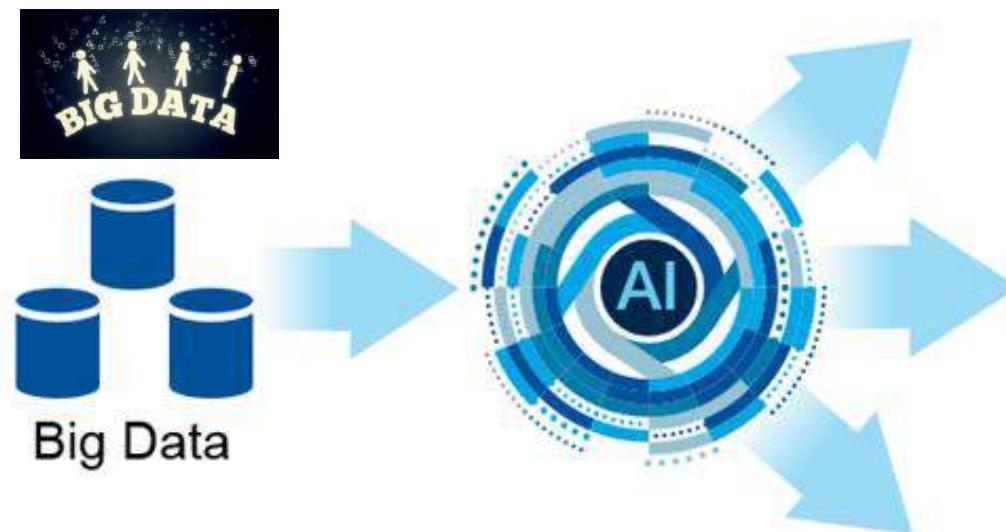
Desierto de Mojave,
212 kilómetros, 6 horas y 54 minutos

https://elpais.com/tecnologia/2005/10/10/actualidad/1128932879_850215.html

Inteligencia artificial, la era del big data. Datos e IA

Datos e Inteligencia Artificial

El **Big Data** aporta una enorme cantidad de datos que alimentan los algoritmos de Machine Learning (Inteligencia Artificial), y permiten crear modelos más complejos y con una mayor precisión.

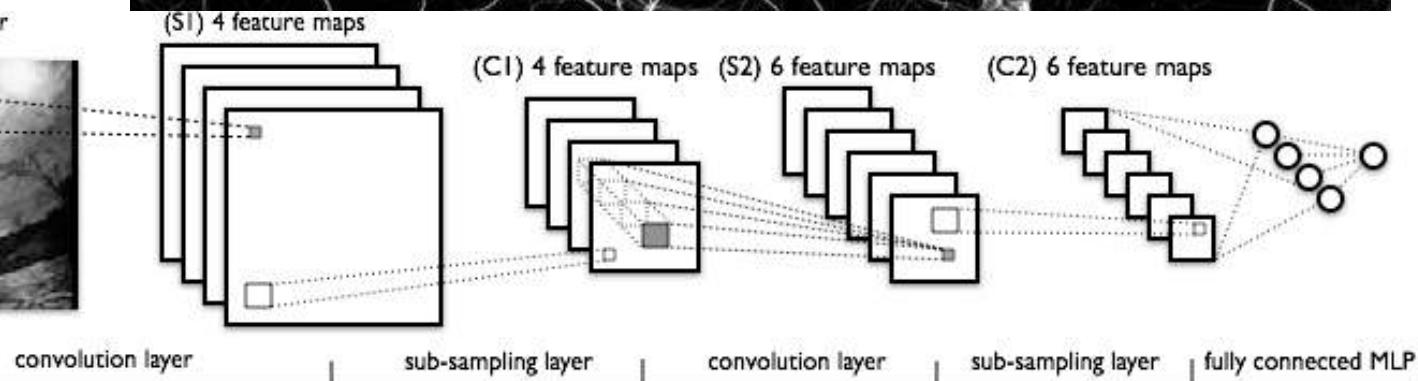
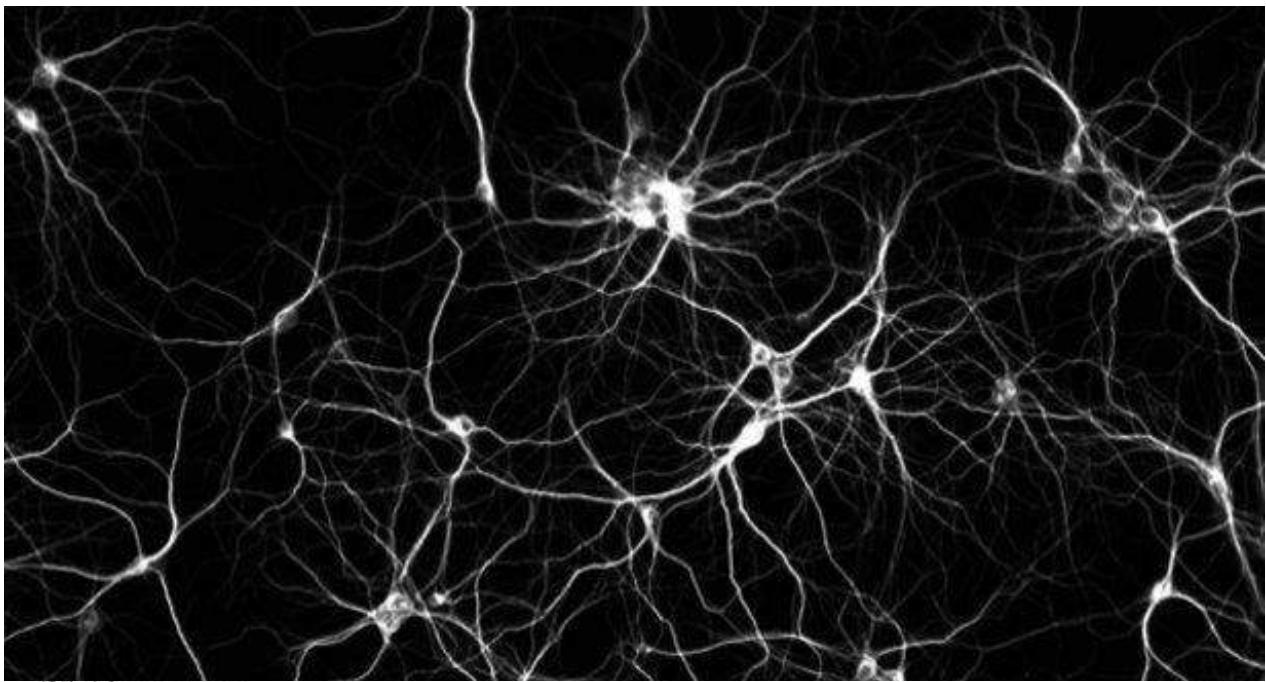
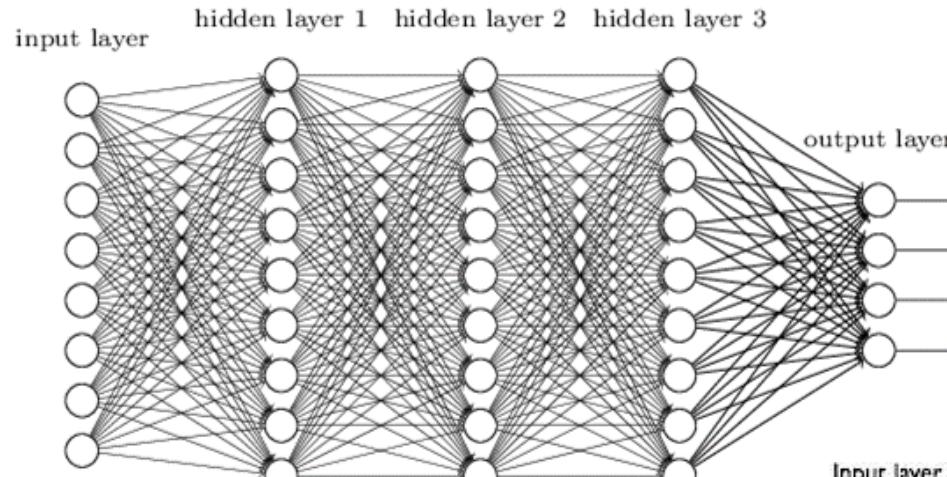


Sistemas Inteligentes alimentados por datos para convertir los datos en conocimiento

Inteligencia artificial, la era del big data. Datos e IA

Deep learning

Deep neural network



Inteligencia artificial, la era del big data. El poder de los datos



Unos días después el director llamó al parente para disculparse.

Respuesta conciliadora del parente:

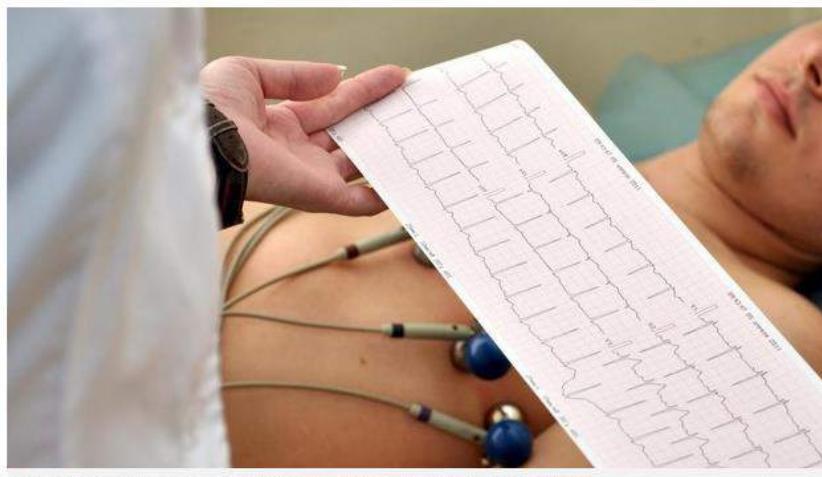
“He estado hablando con mi hija –dijo el parente– Resulta que en mi casa han tenido lugar ciertas actividades de las que yo no estaba del todo informado. Mi hija sale de cuentas en agosto. Soy yo el que les debe una disculpa”.

Inteligencia artificial, la era del big data. El poder de los datos

Una IA predice con éxito si moriremos pronto, aunque nadie sabe cómo lo hace

MADRID Actualizado:15/11/2019 02:13h

- La Inteligencia Artificial logró predecir la muerte de personas, en el plazo de un año, con solo ver sus electrocardiogramas



El software se basa en los resultados de pruebas cardíacas - Adobe Stock

Fornwalt y sus colegas "alimentaron" a la Inteligencia Artificial con una gran cantidad de datos históricos: 1,77 millones de resultados de electrocardiogramas de casi 400.000 personas, y le pidieron que predijera quiénes tenían más probabilidades de morir en los 12 meses siguientes. Un electrocardiograma registra la actividad eléctrica del corazón, y sus patrones cambian debido a las afecciones cardíacas, incluidos los infartos y la fibrilación auricular.

Una **Inteligencia Artificial** acaba de demostrar que es capaz de predecir, con pasmosa exactitud, **las posibilidades de que una persona muera en el plazo de un año**, basándose únicamente en los resultados de sus pruebas cardíacas. El sistema fue incluso capaz de anunciar la muerte de pacientes con valores que los médicos habían considerado normales. Cómo se las arregla la IA para predecir estas muertes resulta un misterio. El impactante estudio, dirigido por Brandon Fornwalt, del Centro Médico Geisinger, en Pennsylvania, se presentará este mismo sábado durante las sesiones científicas de la American Heart Association, en Dallas.

https://www.abc.es/ciencia/abci-predice-exito-si-moriremos-pronto-aunque-nadie-sabe-como-hace-201911141623_noticia.html

<https://www.nature.com/articles/s41591-020-0870-z>

Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

IA fiable y segura

Concluyendo

Inteligencia artificial generativa: 2022 eclosión

En 2022, la inteligencia artificial (IA) se volvió creativa.

Inteligencia Artificial GENERATIVA

2022 Creación a partir de instrucciones “Prompt”

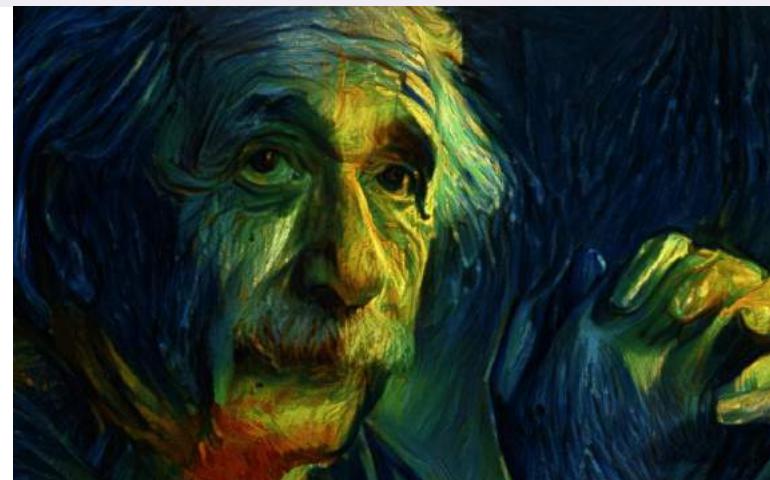


DALL-E 2



LA NOCHE
ESTRELLADA,
VINCENT VAN
GOGH

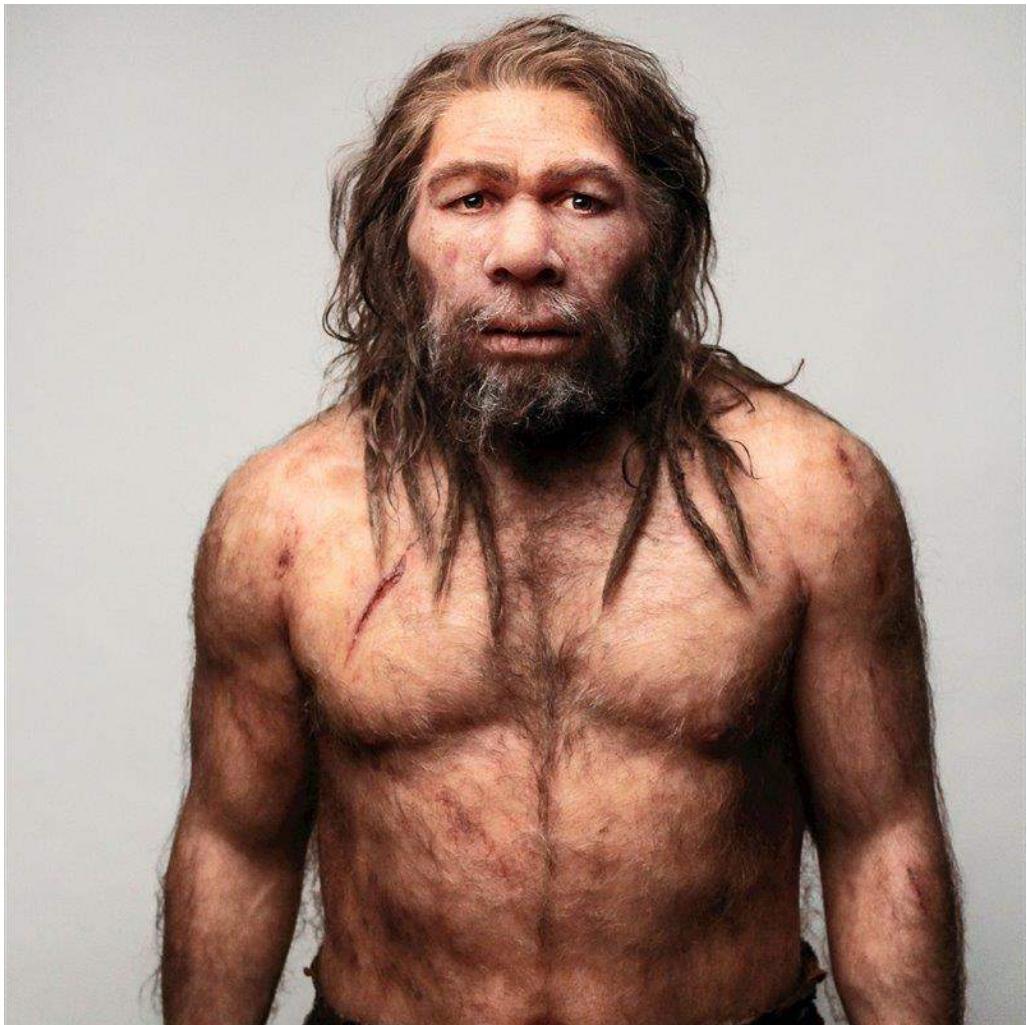
DeepArt.io



Jason Allen’s A.I.-generated work, “**Théâtre D’opéra Spatial**,” took first place in the digital category at the Colorado State Fair.



Homo sapiens



El arte es una de las manifestaciones de la creatividad humana, es la herramienta con la cual *Homo sapiens* desarrolla su cultura, unión y fuerza como pueblo.

El lenguaje permite la comunicación
El lenguaje humano pudo aparecer hace más de 50.000 años

Homo sapiens



20 000 años de antigüedad

<https://arthearty.com/significance-of-lascaux-cave-paintings>



ARTE DIGITAL

"Teatro espacial de la ópera"

ARTE DIGITAL

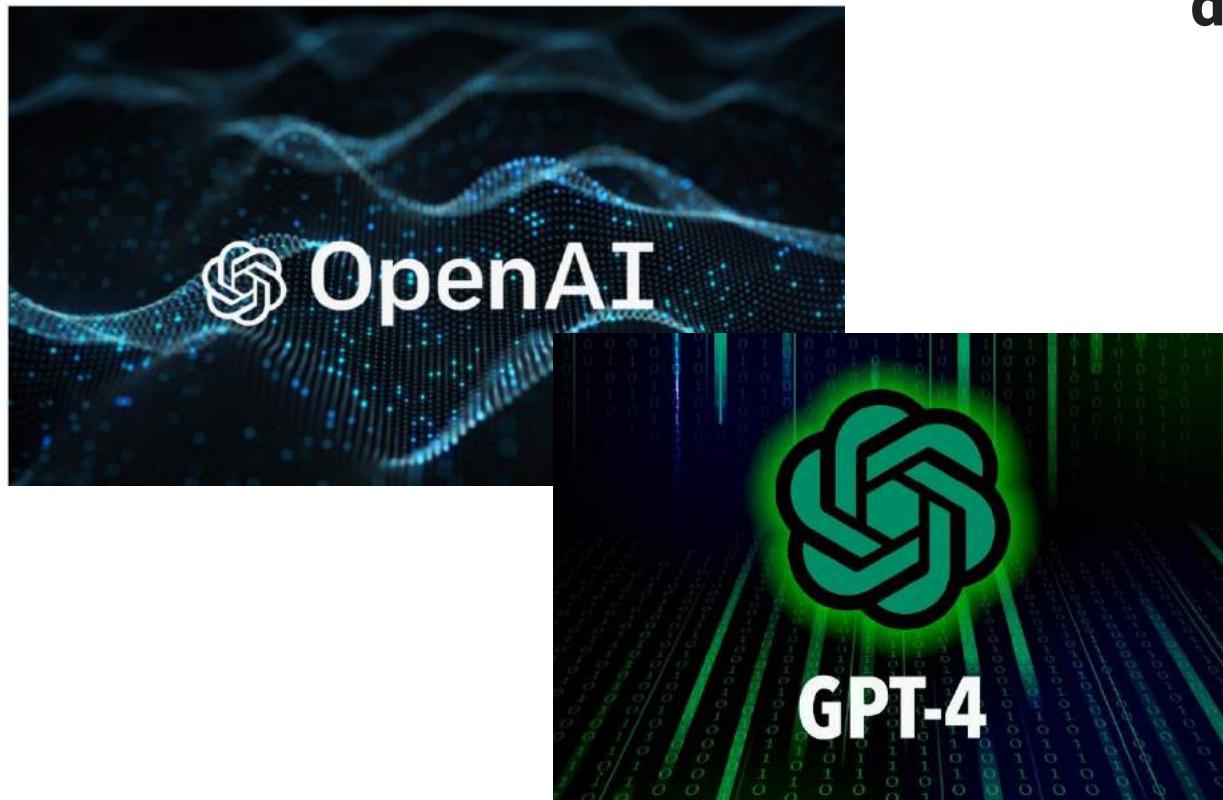
“Teatro
espacial de la
ópera”

Midjourney
Creación a
partir de
instrucciones
(Prompt)
Ganadora
Premio
Colorado,
Agosto 2022



IA Generativa: Lenguaje

GPT-3, el nuevo modelo de lenguaje de OpenAI, es capaz de programar, diseñar y hasta conversar sobre política o economía



Google presenta Bard: el nuevo modelo de lenguaje AI con amplias aplicaciones



IA Generativa (GPT-4 y ChatGPT)

Cornell University

arXiv > cs > arXiv:2303.12712

Computer Science > Computation and Language

[Submitted on 22 Mar 2023 (v1), last revised 24 Mar 2023 (this version, v2)]

Sparks of Artificial General Intelligence: Early experiments with GPT-4

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, Yi Zhang

Microsoft lo tiene claro y en un último informe habla de que su IA, GPT-4, es una forma temprana de inteligencia artificial general (AGI).

Computer Hoy

LO ÚLTIMO ANÁLISIS GUÍAS DE COMPRA LOS MEJORES TUTORIALES PRODUCTOS

Investigadores de Microsoft afirman que la última versión de ChatGPT muestra indicios de inteligencia humana

TECNOLOGÍA Carolina González Valenzuela | 27 mar. 2023 12:45h.

IA Generativa: Sobre las alucinaciones

Uno de los riesgos de la IA generativa se llama alucinación.

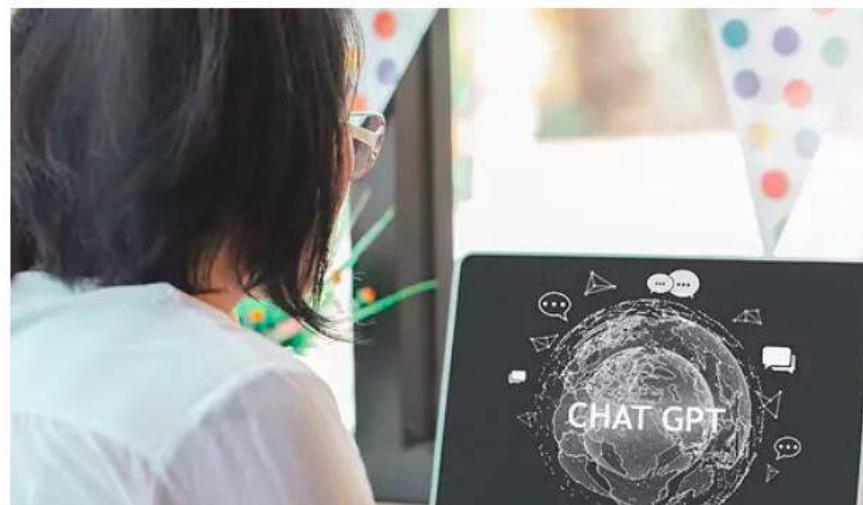
Una de las mayores preocupaciones de los sistemas de IA Generativa es cuando no entienden las preguntas, las malinterpretan y como no pueden generar respuestas correctas, empiezan a inventarlas en un proceso llamado “Alucinación de Inteligencia Artificial”.

Alucinación es el término empleado para el fenómeno en el que los algoritmos de IA y las redes neuronales de aprendizaje profundo producen resultados que no son reales, que no coinciden con ningún dato en el que se haya entrenado el algoritmo u otro patrón identificable.

No puede explicarse por su programación, información de entrada, otros factores como la clasificación incorrecta de datos, capacitación inadecuada, incapacidad para interpretar preguntas en diferentes idiomas, incapacidad para contextualizar preguntas. Las alucinaciones pueden ocurrir en todo tipo de datos sintéticos, como texto, imágenes, audio, video y código informático.

Un abogado admite que usó ChatGPT para un escrito y éste se inventó precedentes legales

El letrado se enfrenta a posibles sanciones y deberá explicar a la Corte por qué no debería ser castigado



Inteligencia artificial generativa: 2022 eclosión

En 2022, la inteligencia artificial (IA) se volvió creativa.

Cada vez hay más tareas en las que cuesta distinguir si algo ha salido de un cerebro humano o de un sistema de IA (ajedrez, arte, descubrimiento proteínas, conversaciones).

Ventajas/beneficios que nos traerá, riesgo de que se usen mal.
Ambos enfoques tienen parte de razón

Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

IA fiable y segura

Concluyendo

Eclosión económica asociada a la inteligencia artificial

COMPUTERWORLD

La IA y los datos generarán 16.500 millones de euros para 2025 en la industria española

Instituciones, expertos académicos y empresarios afirman que el uso de la tecnología basada en datos e inteligencia artificial (IA) es la gran oportunidad para el crecimiento del país, ya que generará ingresos en la industria de 16.000 millones de euros para 2025.

CincoDías

DIGITALIZACIÓN
La economía digital crece 3 puntos en España y ya representa el 22% de PIB

El objetivo del Gobierno es llegar al 40% a finales de 2025

COMBUSTIBLE PARA EL CRECIMIENTO

accenture >

El estudio de Accenture sobre el impacto de la Inteligencia Artificial en 12 economías desarrolladas revela que la IA podría duplicar las tasas anuales de crecimiento económico en 2035, cambiando la naturaleza del trabajo y estableciendo una nueva relación entre el hombre y la máquina. Se prevé que el impacto de la IA en los negocios aumentará la productividad del trabajo hasta en un 40% y permitirá a las personas hacer un uso más eficiente de su tiempo.

La Inteligencia Artificial hará aumentar la economía mundial en 15,7 billones de dólares en 2030

Archivado en: Economía · DES2021 · Inteligencia Artificial



Expertos en Inteligencia Artificial analizarán en DES2021 los retos y usos prácticos de esta tecnología para aumentar la competitividad de las empresas.

Colaboración entre humanos, máquinas y algoritmos



Según un informe del Foro Económico Mundial, se espera que 85 millones de puestos de trabajo sean desplazados por un cambio en la división del trabajo entre humanos y máquinas para 2025. La buena noticia es que pueden surgir 97 millones de nuevos cargos más adaptados a la nueva división del trabajo entre humanos, máquinas y algoritmos, según el mismo informe.

Eclosión de la inteligencia artificial, con riesgos



Compas. Sistema inteligente para calcular el riesgo de reincidencia en el delito

Two Petty Theft Arrests



VERNON PRATER

LOW RISK

3



BRISHA BORDEN

HIGH RISK

8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Two Petty Theft Arrests

VERNON PRATER

Prior Offenses
2 armed robberies, 1 attempted armed robbery

Subsequent Offenses
1 grand theft

BRISHA BORDEN

Prior Offenses
4 juvenile misdemeanors

Subsequent Offenses
None

LOW RISK

3

HIGH RISK

8

Borden was rated high risk for future crime after she and a friend took a kid's bike and scooter that were sitting outside. She did not reoffend.

Eclosión de la inteligencia artificial, con riesgos

noticias ya.

MERCADO Suscríbete

Home / Noticias Locales

Amazon apaga sistema de reclutamiento porque discriminaba a mujeres

AMAZON

Por Saddam Aguayo
Tu opinión cuenta, escríbenos a socialmedia@entravision.com →

octubre 15, 2018, 7:23 pm PST

[Twitter](#) [Facebook](#) [Email](#) [Link](#)



Amazon.com

OelDiario.es 10 AÑOS

Hola, Francisco422

Política Internacional Economía Opinión Cultura Educación Clima Desalambre Igualdad Estatuto

Sesgos de género en los algoritmos: un círculo perverso de discriminación en línea y en la vida real

La “condición algorítmica” está alterando los derechos humanos, que prohíben la discriminación por razón de sexo o género. La industria tecnológica no está haciendo lo suficiente para abordar estos sesgos. El

A nurse in front of a hospital



A doctor in front of a hospital



Eclosión de la inteligencia artificial, con riesgos

'He Would Still Be Here': Man Dies by Suicide After Talking with AI Chatbot, Widow Says

The incident raises concerns about guardrails around quickly-proliferating conversational AI models.



By Chloe Xiang

March 30, 2023, 9:59pm [f Share](#) [t Tweet](#) [s Snap](#)



Listen to this article



A Belgian man recently died by suicide after chatting with an AI chatbot on an app called Chai, Belgian outlet *La Libre* reported.

Eclosión de la inteligencia artificial, con riesgos

SECCIONES



LA NACION

INICIAR SESIÓN



SUSCRIBITE POR \$2600 \$249

Toju Duke, exdirectora de IA de Google: “La inteligencia artificial amplifica injusticias sistémicas que ya deberíamos haber eliminado”

La exdirectora de IA responsable en Google cree que el debate sobre esta tecnología se ha centrado en si la humanidad correrá peligro mañana, cuando el problema es que hoy ya discrimina a la población

20 de octubre de 2023 • 16:25

Manuel G. Pascual

EL PAÍS



Inteligencia artificial. Necesitamos regular

Written Testimony
of
Sam Altman
Chief Executive Officer
OpenAI

Before the U.S. Senate Committee on the Judiciary
Subcommittee on Privacy, Technology, & the Law

“Es esencial regular la inteligencia artificial, y que esas normas garanticen que el público acceda a los muchos beneficios de esta tecnología”,

“Mi peor miedo es que esta tecnología salga mal. Y si sale mal, puede salir muy mal”.

https://www.washingtonpost.com/documents/0668f6f4-d957-4b94-a745-2aa9617d1d60.pdf?itid=lk_inline_manual_18

Sam Altman (ChatGPT), en el Capitolio: “Si la inteligencia artificial sale mal, puede salir muy mal”

El cofundador de OpenAI comparece ante el Senado estadounidense para defender los beneficios y alertar de los riesgos de la tecnología revolucionaria

Washington - 16 MAY 2023 - 19:10 CEST



Inteligencia artificial. Necesitamos regular



About Us Our Work ▾ FAQ AI Risk Contact Us

Statement on AI Risk

AI experts and public figures express their concern about AI risk.

Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.

Signatories:



AI Scientists



Other Notable Figures

Geoffrey Hinton

Emeritus Professor of Computer Science, University of Toronto

Yoshua Bengio

Professor of Computer Science, U. Montreal / Mila

Demis Hassabis

CEO, Google DeepMind

Sam Altman

CEO, OpenAI

Dario Amodei

CEO, Anthropic

Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

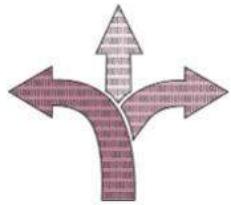
IA fiable y segura

Concluyendo

IA: Ética y regulación. Debate ético y Principios éticos de la IA

Aunque los problemas de sesgo y discriminación siempre han estado presentes en la sociedad, la preocupación es que la IA pueda perpetuar estos problemas y acrecentar su impacto

ÉTICA DE LA INTELIGENCIA ARTIFICIAL
Mark Coeckelbergh



CÁTEDRA

Aproximación filosófica:

- El papel de la IA y los humanos
- Análisis la omnipresencia de la IA
- Análisis ético y social de su aplicación
- Políticas de IA de cara al futuro
- Vigilar el discurso público con respecto a la IA

Principios éticos de Naciones Unidas (Nov. 2021)

1. Proporcionalidad y no amenaza
2. Seguridad
3. Equidad y no discriminación
4. Sostenibilidad
5. Derecho a la privacidad, y protección de datos
6. Control Humano
7. Transparencia y explicabilidad
8. Responsabilidad y rendición de cuentas
9. Conciencia y educación sobre la IA
10. Gobernanza y colaboración adaptativa y de múltiples partes interesadas

Inteligencia artificial. Cumbre para la regulación y seguridad

• Últimas noticias

El Confidencial

Iniciar sesión

LA IA YA ESTÁ CAMBIANDO LAS GUERRAS

Cumbre en UK contra el riesgo "catastrófico" de la IA. ¿Puede acabar con la humanidad?

Cientos de líderes mundiales, académicos, investigadores y empresarios, incluido el hombre más rico del mundo, Elon Musk, se dan cita esta semana en Reino Unido para abordar los retos que plantea la inteligencia artificial



Policy paper

The Bletchley Declaration by Countries Attending the AI Safety Summit, 1-2 November 2023

Published 1 November 2023



En el contexto de nuestra cooperación, y para informar la acción a nivel nacional e internacional, nuestra agenda para abordar el riesgo fronterizo de la IA se centrará en:

- **identificar los riesgos de seguridad de la IA de interés compartido, ...**
- **construir políticas respectivas basadas en riesgos en todos nuestros países para garantizar la seguridad a la luz de dichos riesgos,...**

IA: Etica y regulación

Marco regulatorio europeo

THE AI ACT

About

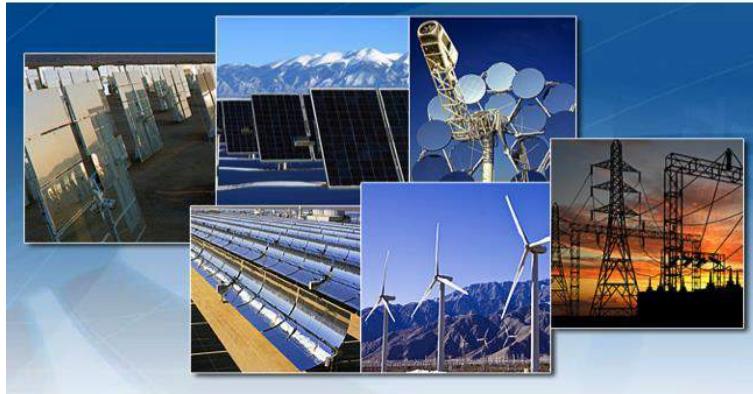
The Act ▾

Assessment

Analyses

The Artificial Intelligence Act

Alto riesgo



Regulación europea: AI Act

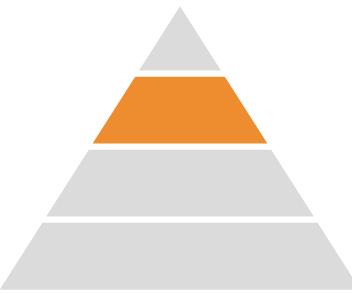
Inaceptables

Social Scoring (Black Mirror)

- identificación biométrica y categorización de personas físicas,
- gestión y funcionamiento de infraestructuras esenciales,
- educación y formación profesional,
- empleo, gestión de los trabajadores y acceso al autoempleo,
- acceso y disfrute de servicios públicos y privados esenciales y sus beneficios,
- asuntos relacionados con la aplicación de la ley,
- gestión de la migración, el asilo y el control fronterizo, o
- administración de justicia y procesos democráticos.

Alto riesgo

NETFLIX



IA: Etica y regulación. Auditabilidad

Establecer e implementar un **sistema de gestión de riesgos**

a la luz del **propósito previsto** del sistema de IA.

Utilizar **datos** de entrenamiento, validación y pruebas de **alta calidad** (relevantes, representativos, etc.)

Elaborar **documentación técnica** y **configurar capacidades de registro** (trazabilidad y auditabilidad)

Garantizar un grado adecuado de **transparencia** y proporcionar a los usuarios información sobre las capacidades y limitaciones del sistema y cómo utilizarlo

Garantizar la **supervisión humana** (medidas integradas en el sistema y/o a implementar por los usuarios)

Garantizar **robustez, precisión y ciberseguridad**

IA: Regulación y derechos e implicaciones sociales (coches autónomos)

¿Pautas legales para juzgar
a un sistema inteligente?

COCHE AUTÓNOMO TESLA ›

Tesla reconoce el segundo accidente mortal en EE UU con un coche que circulaba en piloto automático

El siniestro tuvo lugar hace una semana en California cuando un Model X chocó contra un muro después de que el sistema avisase al conductor para que pusiera las manos en el volante

EL PAÍS G+
Madrid - 31 MAR 2018 - 13:33 CEST

29 Marzo 2018 – Segundo accidente mortal de un coche Tesla con piloto automático



Los servicios de rescate trabajan en el lugar donde se estrelló el Model X, en California. STRINGER (REUTERS)

IA: Regulación y derechos e implicaciones sociales (Impacto social en el trabajo)

≡

EL PAÍS

NEGOCIOS

EL FUTURO DEL TRABAJO > OPINIÓN

¿Debemos preocuparnos porque los robots se quedan con nuestro trabajo?

No cabe duda de que muchos empleos vinculados a la economía del conocimiento pueden ser sustituibles



Robots y personas trabajan en un almacén industrial.
LOCUS ROBOTICS/PR NEWSWIRE (EUROPA PRESS)

Paul Krugman es premio Nobel de Economía.
© The New York Times, 2022.

¿qué pasaría si las máquinas pudieran hacerse cargo de una gran parte de lo que tradicionalmente hemos considerado trabajo del conocimiento?

ChatGPT, tecnología que parece capaz de realizar tareas que requieren una considerable formación académica.

¿cuánto de lo que hacemos los seres humanos es creativo de verdad o aporta una comprensión en profundidad?

Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

IA fiable y segura. Tecnología de IA

Concluyendo

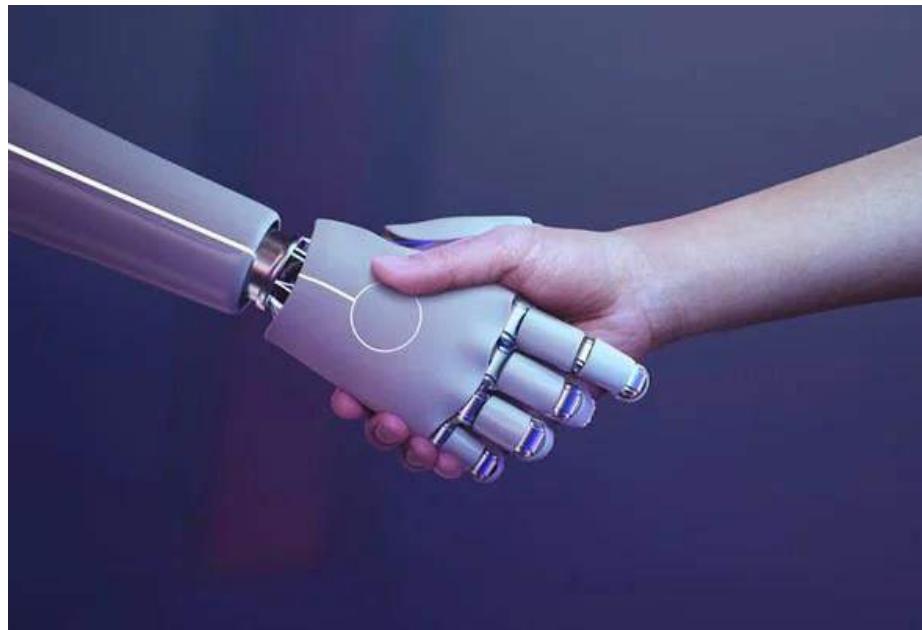


DIRECTRICES ÉTICAS
PARA UNA IA FIABLE

IA: Regulación, Inteligencia artificial fiable

Inteligencia artificial fiable (Trustworthy AI)

1. Lícita (cumple la ley)
2. Ética (principios éticos)
3. Robusta (segura en su uso, control de fallos, ...)



Inteligencia artificial fiable

Tres pilares básicos de una IA fiable:

- **Conforme a la ley.** La IA tiene que ser lícita
- **Ética.** El respeto de la ley no es siempre garante del respeto de principios éticos que salvaguarden los derechos inalienables de las personas.

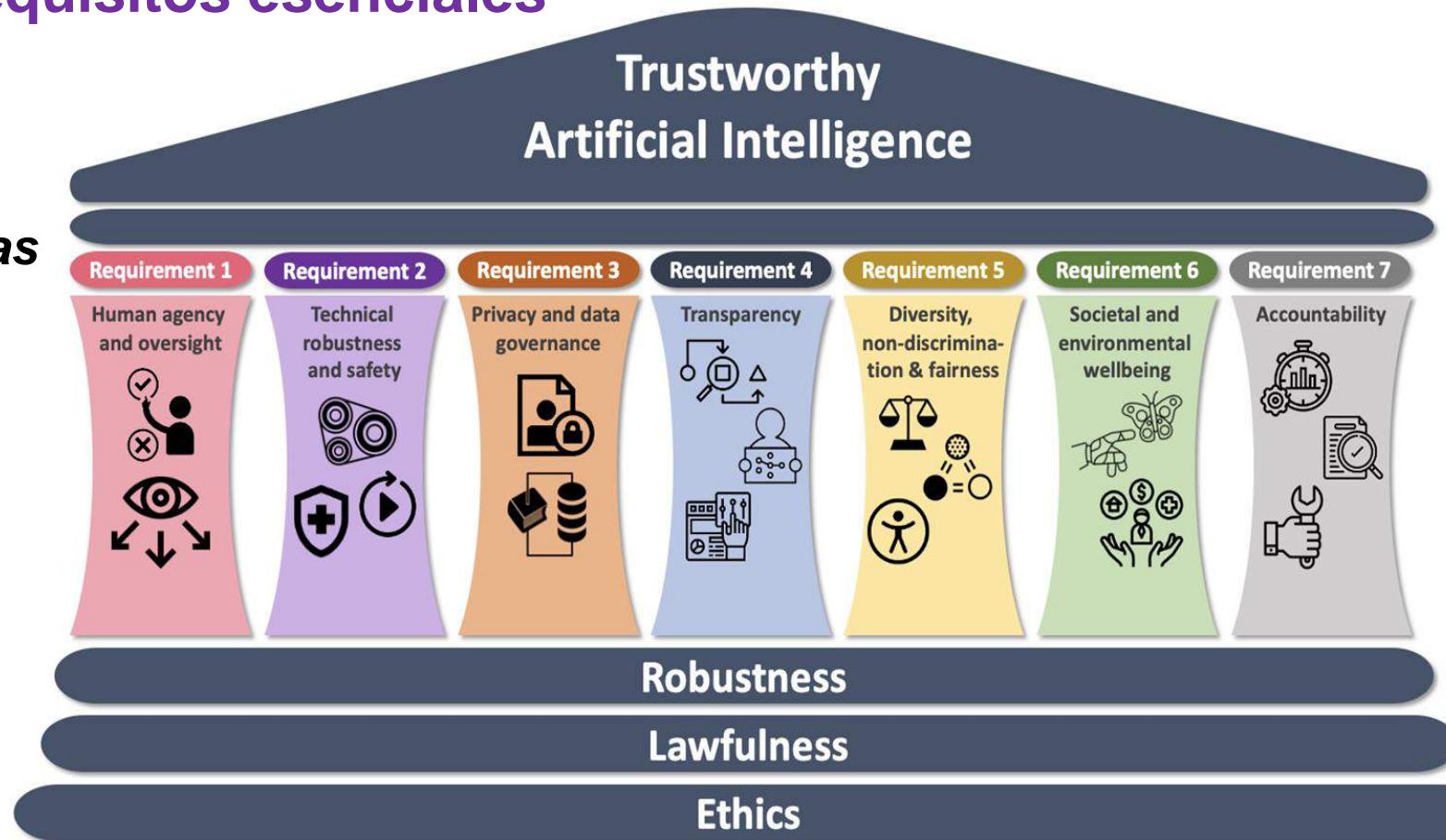
(Ético es: Proteger la igualdad, la no discriminación y la solidaridad entre los ciudadanos. Respetar la democracia, la justicia y el estado de Derecho. Respetar la dignidad humana. Respetar la libertad individual. Respetar los derechos de los ciudadanos)

- **Robusta.** los sistemas de IA deben ofrecer un correcto funcionamiento siempre y prever medidas de protección para evitar cualquier efecto adverso imprevisto.

Inteligencia artificial fiable. Etica, legal y robusta

IA cumpliendo siete requisitos esenciales

- I. *Intervención y supervisión humanas.*
- II. *Solidez y seguridad técnicas*
- III. *Privacidad y gestión de datos.*
- IV. *Transparencia.*
- V. *Diversidad, no discriminación y equidad.*
- VI. *Bienestar social y medioambiental.*
- VII. *Rendición de cuentas.*



Inteligencia artificial fiable

Siete requisitos esenciales que debe cumplir un sistema inteligente

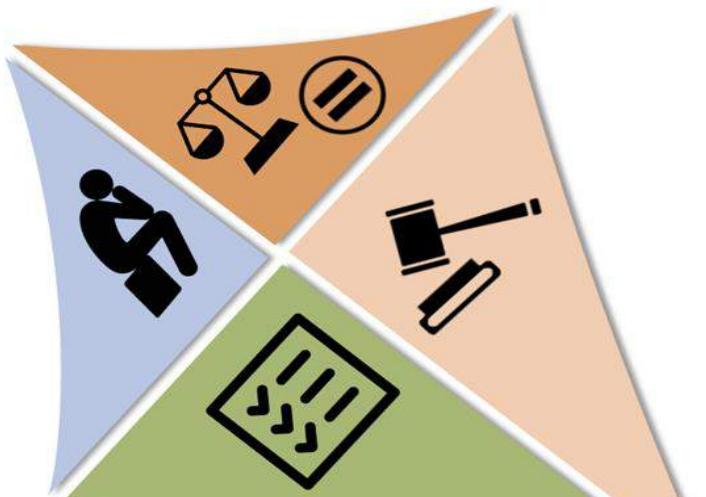
- I. **Intervención y supervisión humanas.** los sistemas de IA tienen que permitir ser gobernados o supervisados por humanos.
- II. **Solidez y seguridad técnicas.** Los sistemas de IA tienen que garantizar robustez tecnológica e incluso considerar planes de contingencia para la adaptación ante comportamientos anómalos.
- III. **Privacidad y gestión de datos.** Los datos tienen que estar protegidos
- IV. **Transparencia.** El comportamiento de los sistemas de IA debe poder ser monitorizado o trazado
- V. **Diversidad, no discriminación y equidad.** El proceso de adquisición y anotación de los datos tiene que preservar la igualdad y evitar la discriminación de los ciudadanos.
- VI. **Bienestar social y medioambiental.**
- VII. **Rendición de cuentas.** Esta directriz está relacionada con el principio de responsabilidad.

Inteligencia artificial fiable Visión holística

**Principles for ethical use
and development of AI**

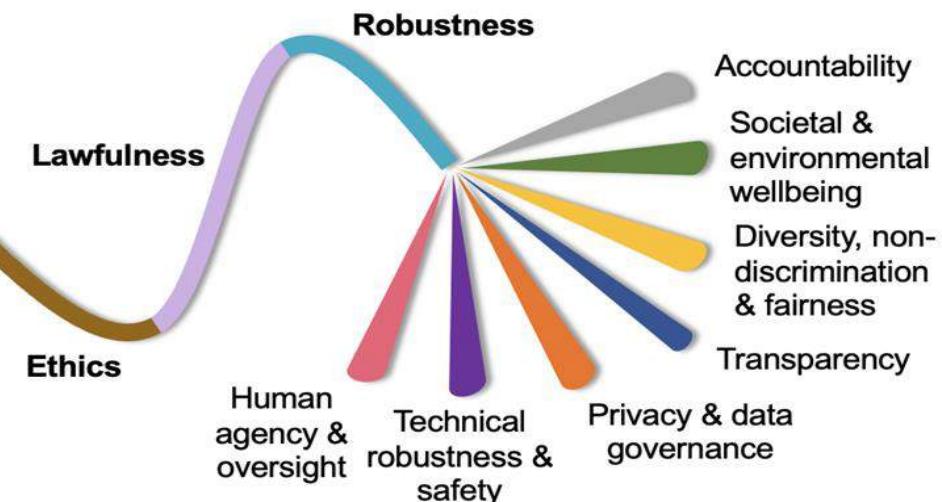
**Artificial Intelligence regulation:
A risk-based approach**

**A philosophical approach
to AI ethics**



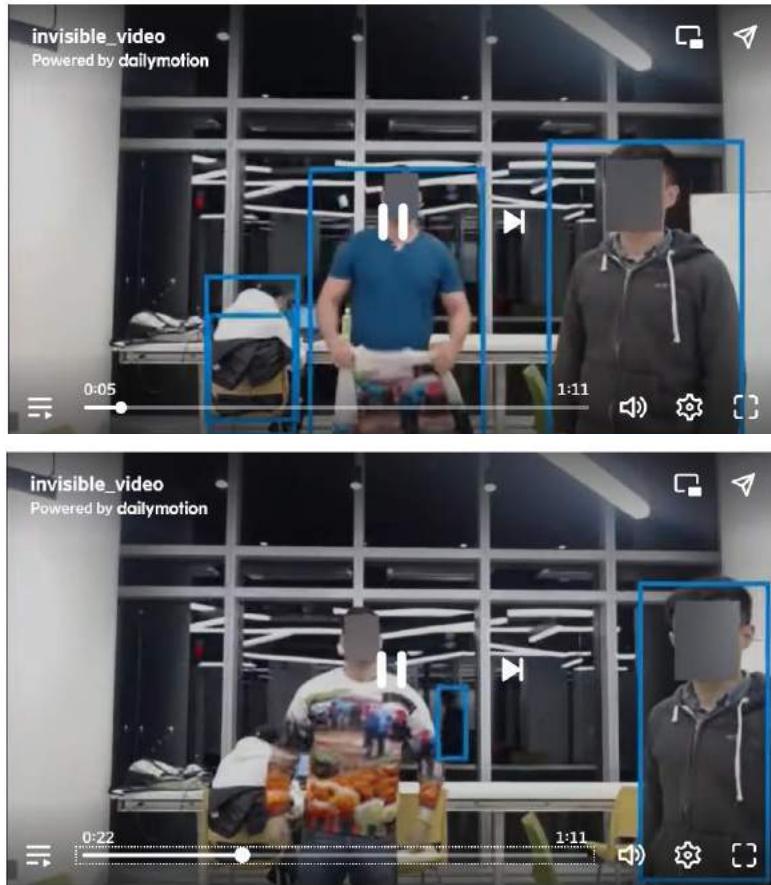
**Pillars and
Requirements of
Trustworthy AI**

**From Trustworthy AI to
Responsible AI Systems**

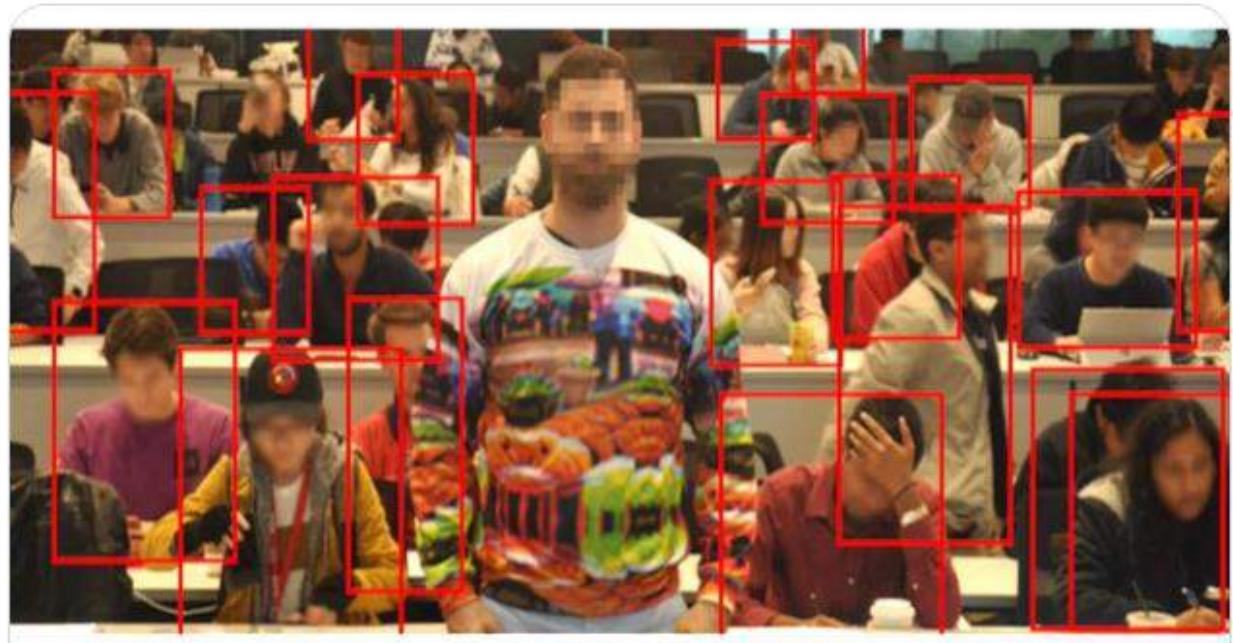


IA segura : Reto. Robustez, evitar adversarios

Este jersey funciona como capa de invisibilidad frente la Inteligencia Artificial

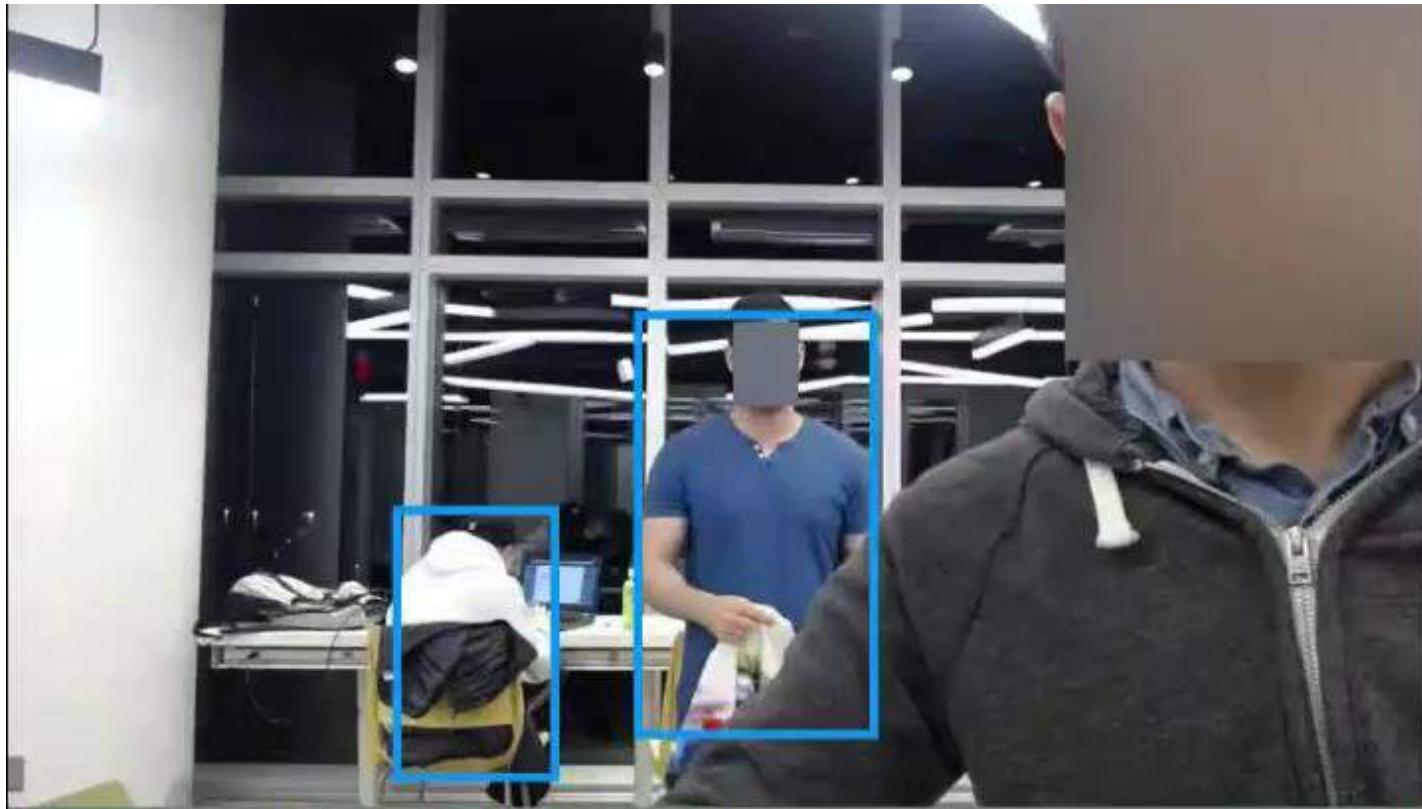


un jersey que utiliza "patrones adversos" que eliminan los patrones comunes a los que están entrenados estos sistemas y lo convierte en invisible ante la IA.



IA segura : Reto. Robustez, evitar adversarios

Este jersey funciona como capa de invisibilidad frente la Inteligencia Artificial



IA Segura:

Abarca la ética de las máquinas, la alineación de la IA con principios éticos y metas humanas, cuyo objetivo es hacer que los sistemas de IA sean morales y beneficiosos.

Tecnológicamente abarca problemas para monitorear los sistemas para detectar riesgos y hacerlos altamente fiables. Más allá de la investigación en IA, implica desarrollar normas y políticas que promuevan la seguridad.

IA segura

• Últimas noticias

El Confidencial

Iniciar sesión

MUSTAFA SULEYMAN, COFUNDADOR DE DEEPMIND

Este hijo de un taxista es uno de los 'padres' de la IA y avisa sobre todo lo que va a cambiar

Hubo una época en que la empresa británica DeepMind era el OpenAI del momento. Google la compró, pero su cofundador, Mustafa Suleyman, dejó el buscador decepcionado. Ahora avisa de los riesgos de la inteligencia artificial



TECNOLOGÍA, PODER Y
EL GRAN DILEMA DEL SIGLO XXI

LA
OLA QUE
VIENE
MUSTAFA SULEYMAN

DEBATE

“Necesitamos controles en dos niveles: uno, para controlar lo que los sistemas pueden hacer por sí mismos, y otro, para controlar lo que pueden hacer bajo instrucciones humanas.”

Inteligencia artificial

Inteligencia artificial generativa

Eclosión de la IA, con riesgos

Etica y regulación. Auditabilidad

IA fiable y segura

Concluyendo

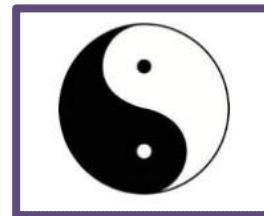
Concluyendo, R1 (Una IA silenciosa que ha entrado en nuestras vidas)

«Máquinas no pensantes cada vez más capaces.
Cualquier tarea que requiera menos de diez segundos de
pensamiento podrá ser hecha por una IA»



Concluyendo, R2: Los datos son la base de la inteligencia artificial y el aprendizaje automático

“Inteligencia Artificial es la nueva electricidad” A. Ng

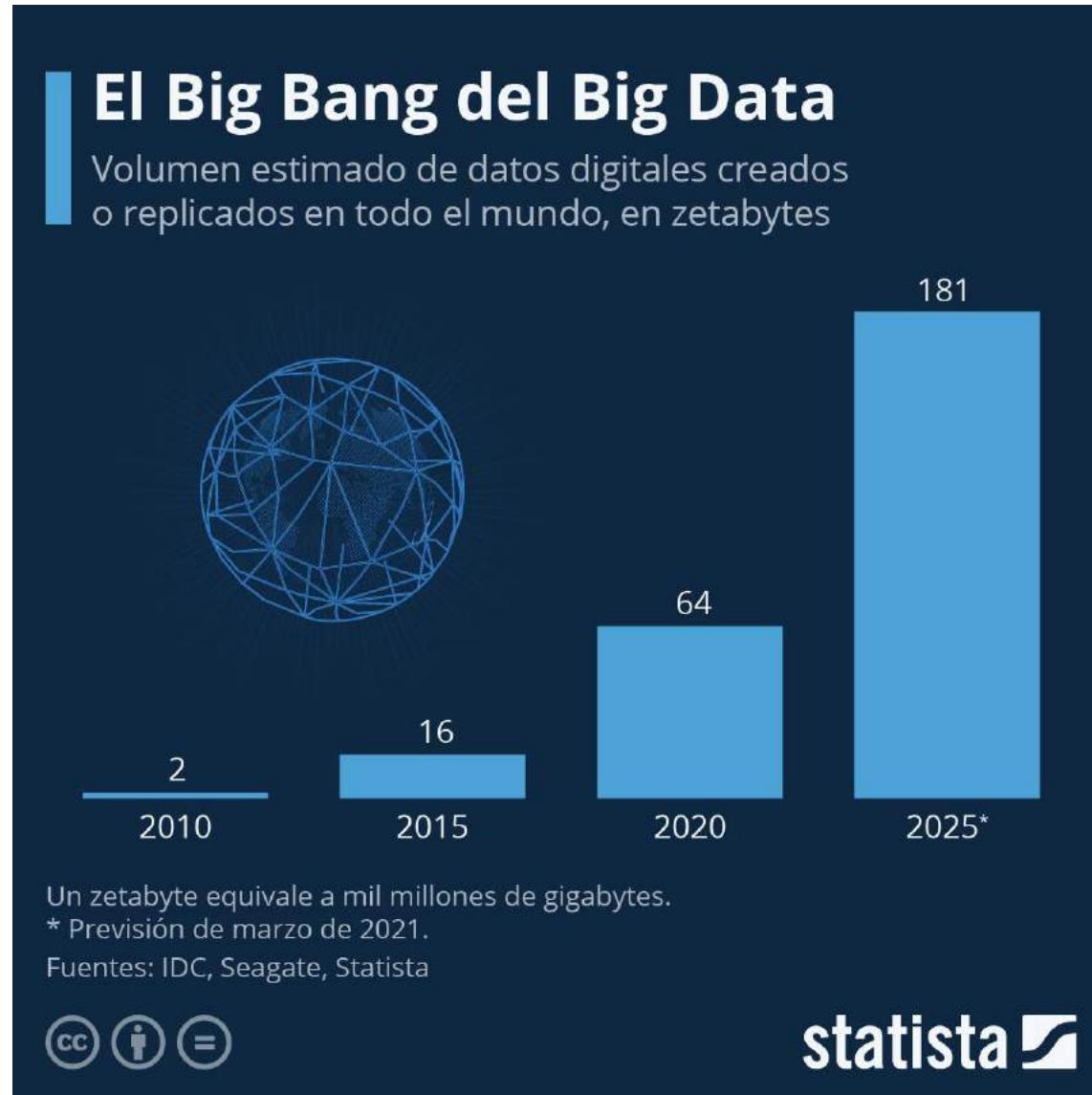


(Big) data es “el nuevo petróleo/la nueva materia prima” del siglo XXI



Concluyendo, R2: Los datos son la base de la inteligencia artificial y el aprendizaje automático (el mundo gira alrededor de los datos)

La explosión de los datos



Concluyendo, R3: reflexiones sobre retos/problemas a abordar

Tres grandes retos



UN GRAN ERROR
SOBRE LA IA ES
PENSAR QUE ES
INTELIGENTE”

SANDRA WACHTER

Professor Sandra Wachter is Professor of Technology and Regulation at the Oxford Internet Institute at the University of Oxford



“Siempre hay una cuestión con la protección de datos, porque se usan datos para alimentar el algoritmo, no es posible hacerlo sin datos.

Siempre hay una cuestión con la explicabilidad, en el sentido de que no entiendes completamente cómo funciona el algoritmo y por qué toma ciertas decisiones.

Y, muy a menudo, hay también un problema de sesgo porque, de nuevo, los datos son históricos, y la mayoría de conjuntos de datos no son todos los datos, son datos sesgados.”

Concluyendo: Debate. Sobre creatividad

El sueño de la máquina creativa

EL FOCO

La cuestión no es si el arte de los ordenadores lo hará mejor que nosotros, sino pensar qué podemos hacer únicamente nosotros cuando las computadoras alcanzan tal sofisticación



DANIEL INNERARITY

Catedrático de Filosofía Política, investigador Ikerbasque en la UPV/EHU y titular de la cátedra Inteligencia Artificial y Democracia en el Instituto Europeo de Florencia

Domingo, 5 febrero 2023, 00:49

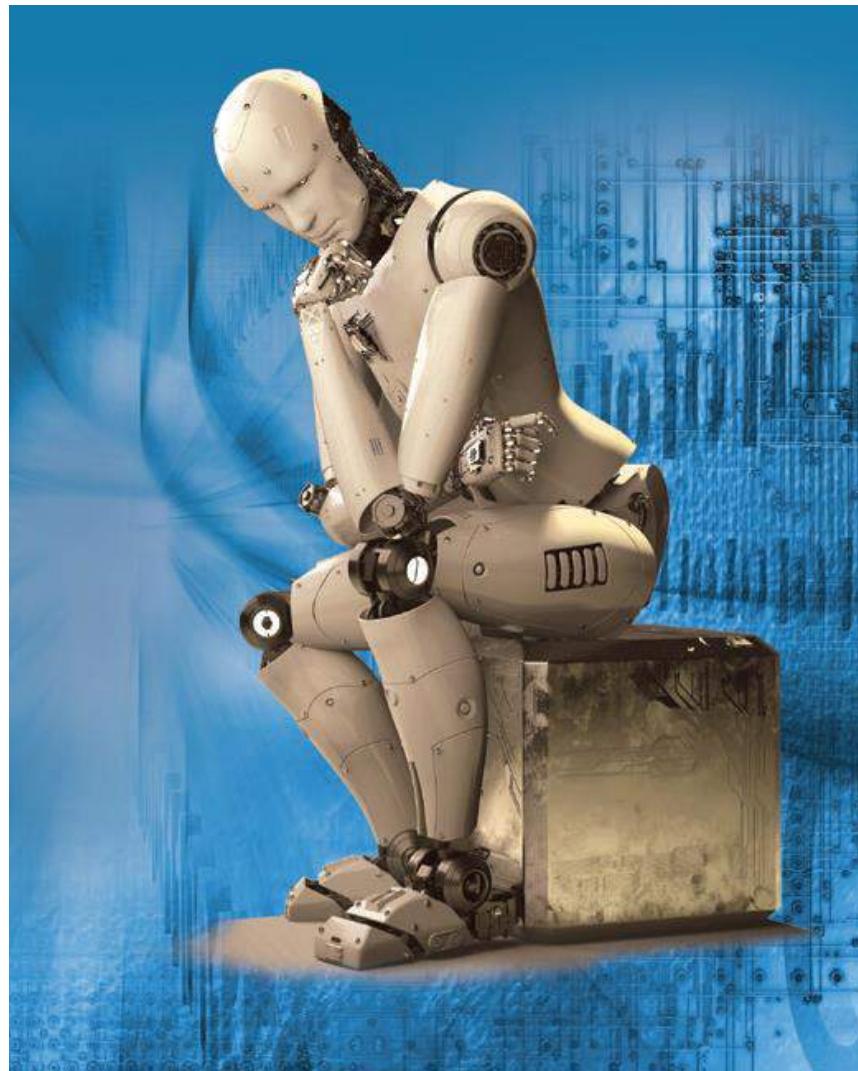
La creatividad humana no puede ni imitarse ni repetirse: implica siempre una cierta trasgresión

Los ordenadores tienen una forma débil de creatividad que les permite reproducir patrones de habla, sonido o formas, pero nada más. De un ordenador no puede esperarse que produzca algo radicalmente imprevisible, nada similar a lo que supuso la vanguardia o los creadores verdaderamente disruptivos en la historia de las artes.

Concluyendo: Las profesiones del futuro

¿Qué harán los robots dentro de 15 años?

“Entre el 50% y 65% de los alumnos que acceden hoy a la escuela primaria trabajarán en profesiones que aún no existen”

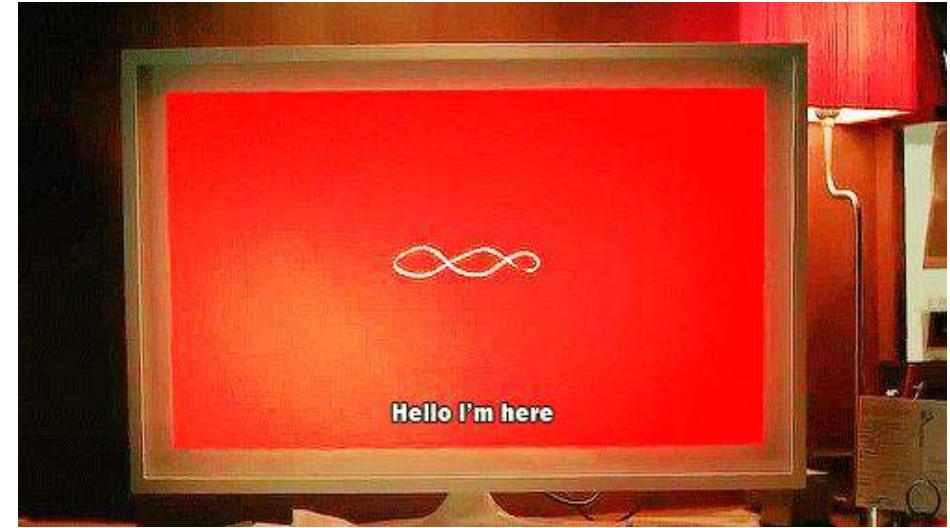


Concluyendo: Inteligencia artificial. Ficción versus realidad



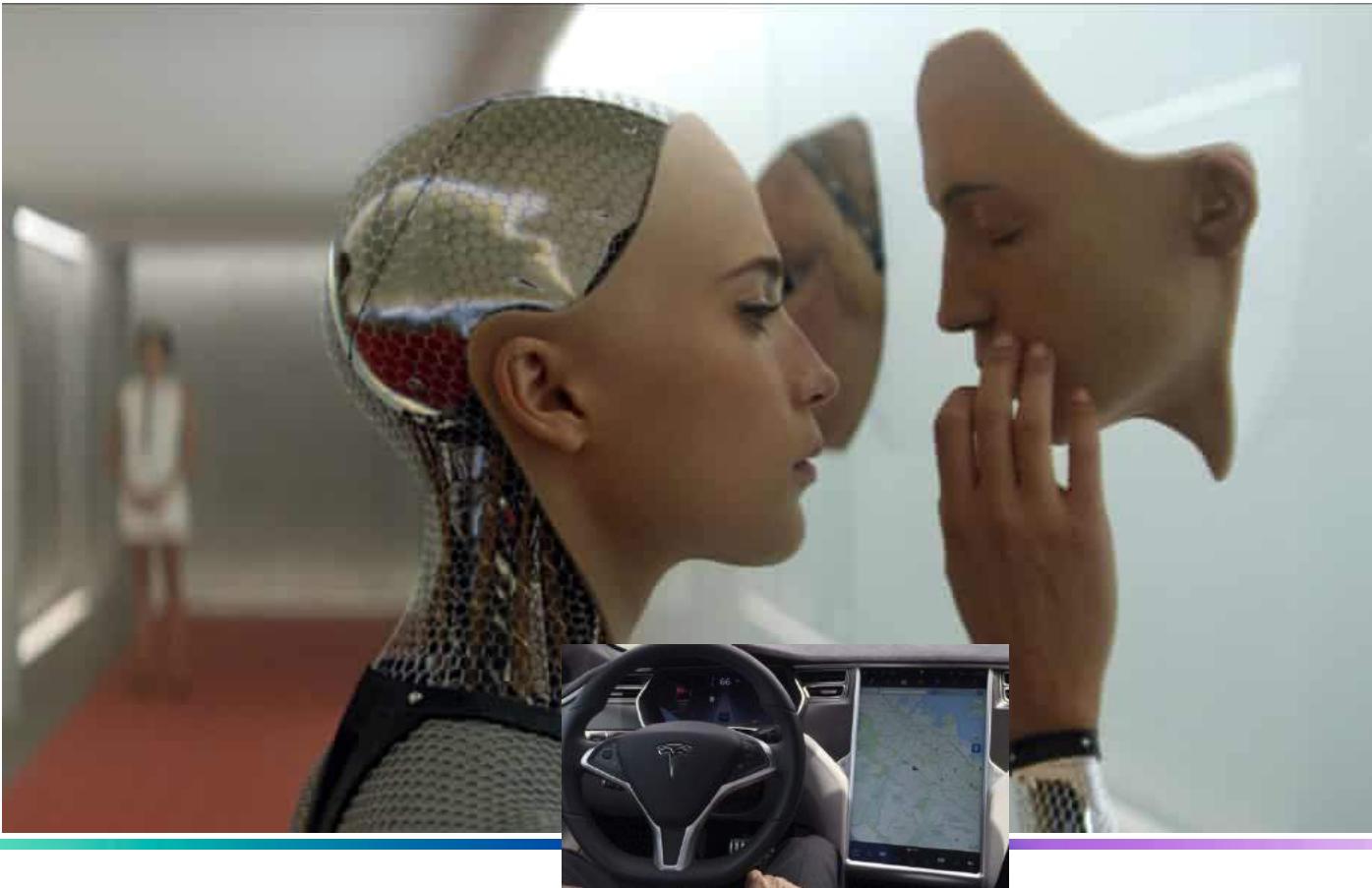
HAL 9000
2001: Una odisea del espacio
Stanley Kubrick 1968

Ava
Ex Machina
Alex Garland, 2015



Concluyendo: Inteligencia artificial. Ficción versus realidad

¿Como va a interactuar un sistema de IA con los humanos?



¿Estamos preparados para convivir con máquinas inteligentes?

¿Cómo actuaríamos hoy ante un coche sin conductor?



Concluyendo: Debate. ¿Alcanzaremos una IA general que iguale al humano?

Oriol Vinyal (Responsable de Investigación en DeepMind-Google):

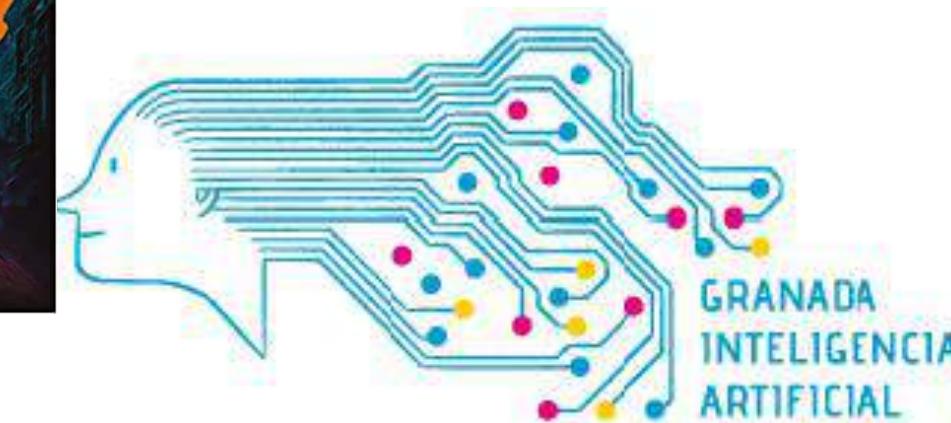
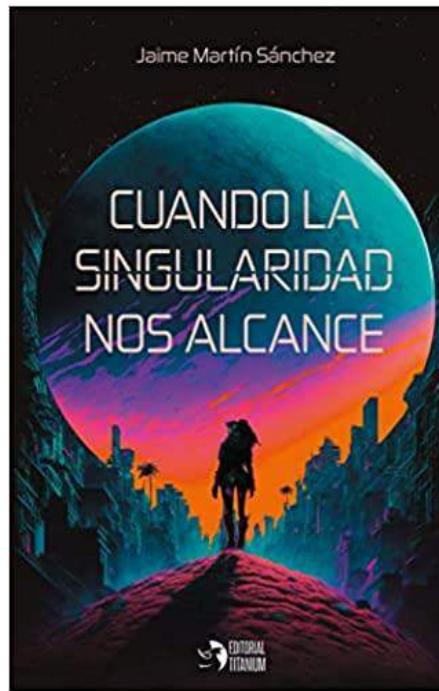
“Nuestra generación verá una inteligencia artificial que iguale o supere a la del ser humano”



The screenshot shows the official website for OpenAI. At the top, there is a navigation bar with links for Research, API, ChatGPT, Safety, and Company. On the right side of the header are Search, Log in, and Try ChatGPT buttons. The main banner features a dark background with a person sitting on a couch and a potted plant in the foreground. The text "Creating safe AGI that benefits all of humanity" is prominently displayed in white. Below the banner is a "Learn about OpenAI" button. At the bottom of the page, there are three columns of text: "Pioneering research on the path to AGI" with a "Learn about our research" link; "Transforming work and creativity with AI" with an "Explore our products" link; and "Join us in shaping the future of technology" with a "View careers" link.



La INTELIGENCIA ARTIFICIAL silenciosa ha entrado en nuestras vidas



Lecturas (novelas) recomendadas (distopías del futuro)

OUTLINE

Artificial Intelligence hatching: Between risk and regulation

Responsible and safe AI: Trustworthy AI

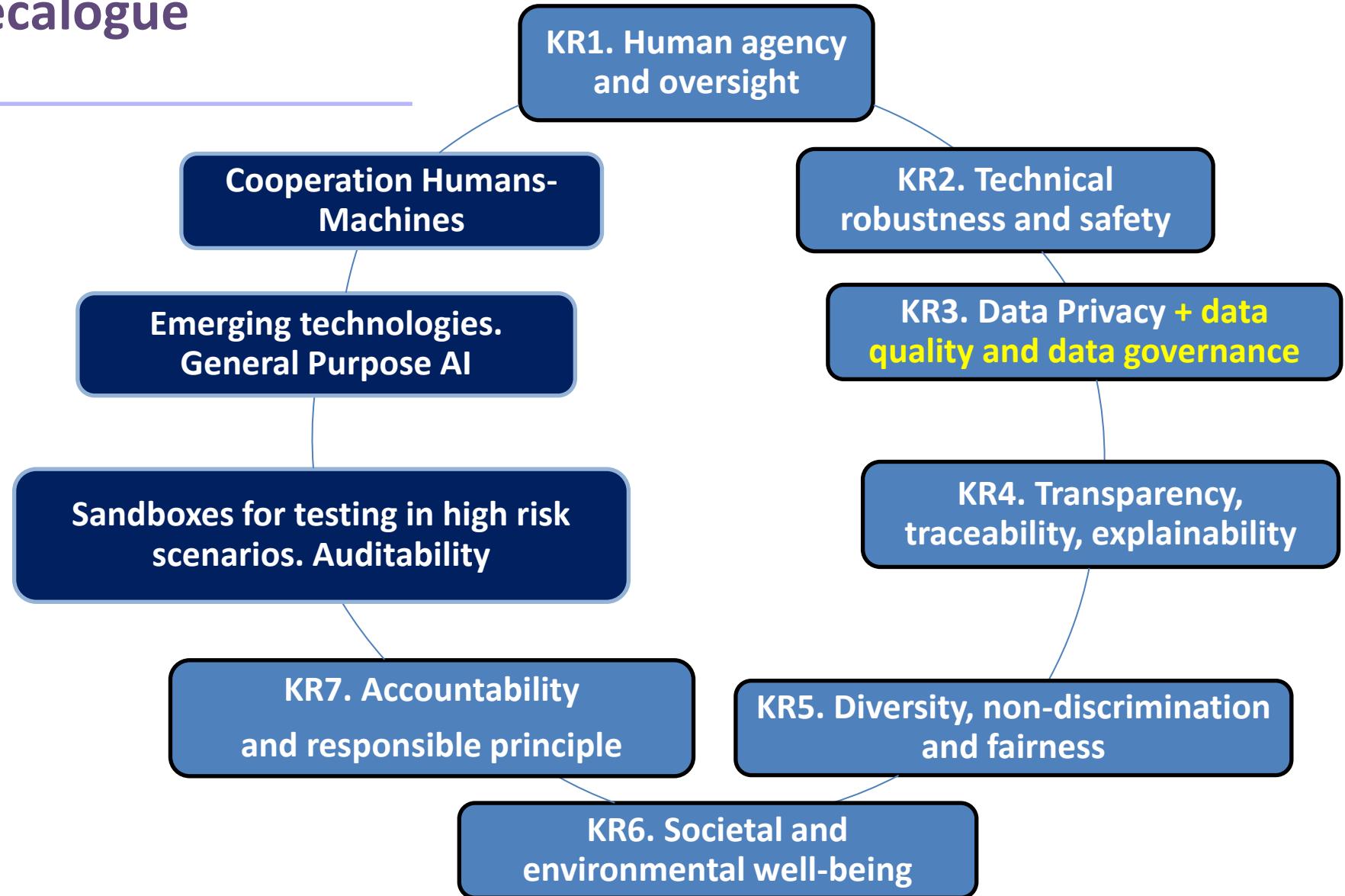
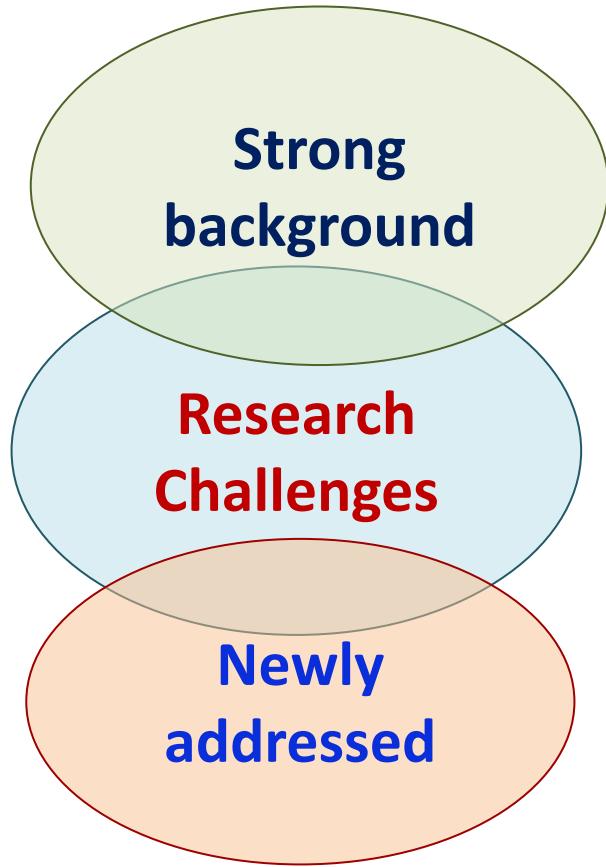
Responsible AI Systems and auditability: European risk-based scenario

Trustworthy AI. A decalogue. Collaboration

Ethics, regulation and social implications: ELSEC

Conclusions

Trustworthy AI: A decalogue



Trustworthy AI decalogue: Transparency: traceability, explainability, and communication (KR4)

Traceability is defined as the set of mechanisms and procedures aimed to keep track of the system's data, development and deployment processes

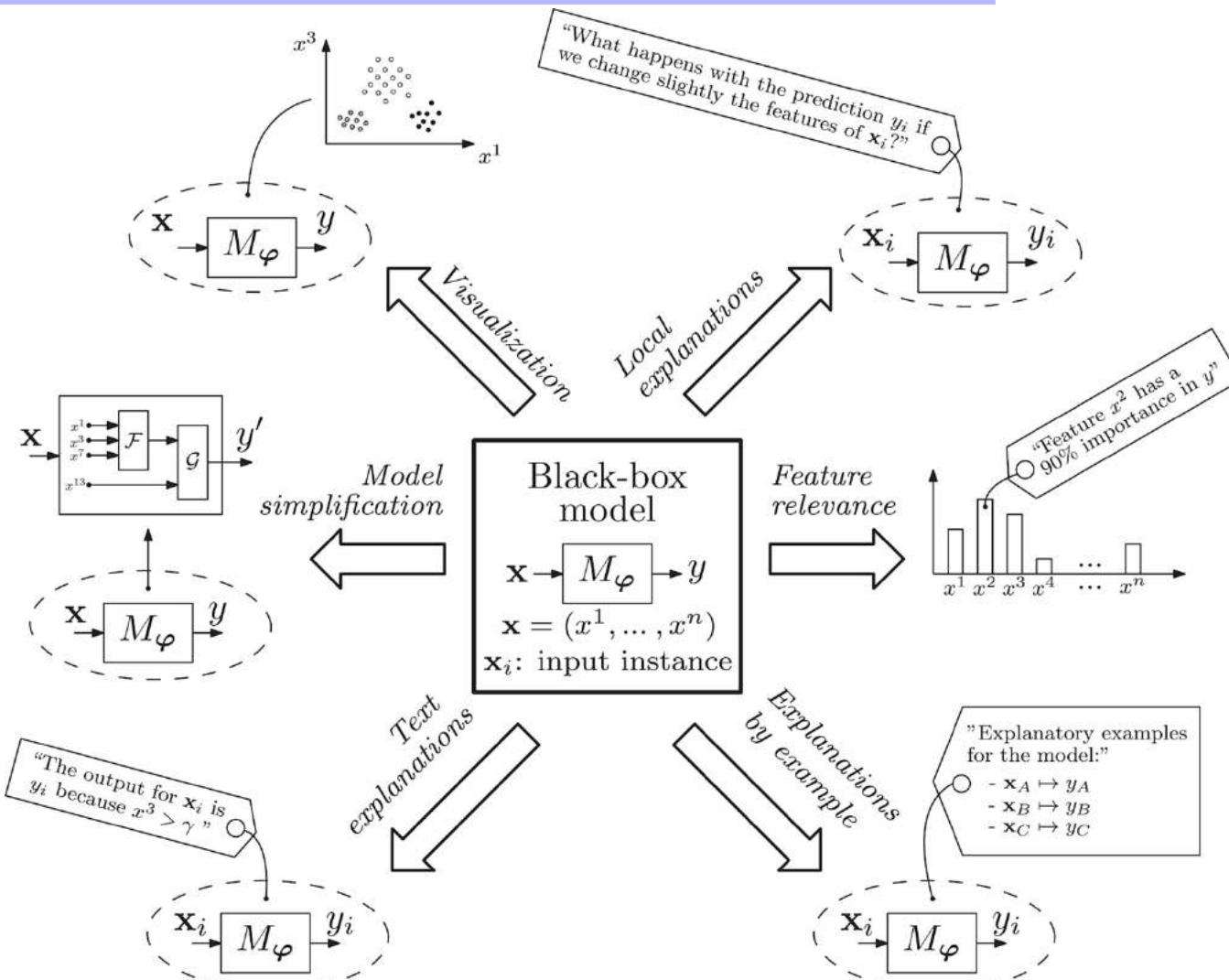
Explainability

Given an audience, an explainable AI is one that produces details or reasons to make its functioning clear or easy to understand.



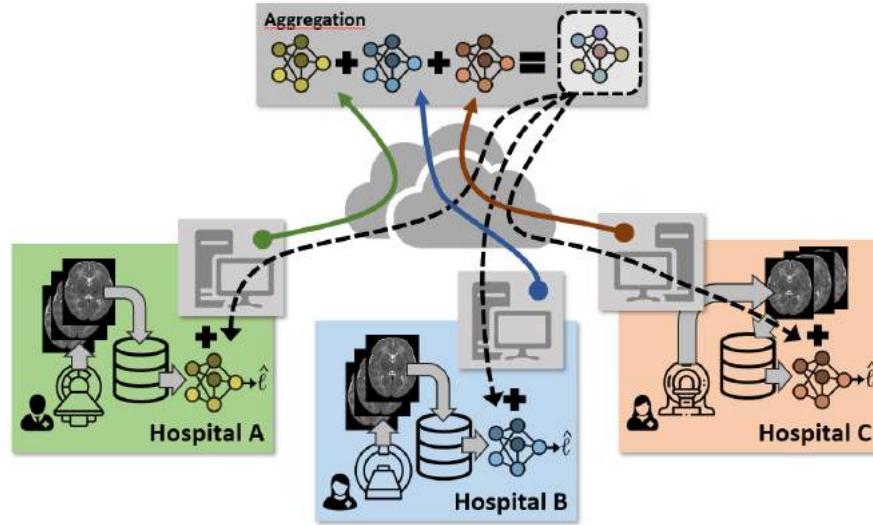
Communication is a crucial dimension, so that all aspects related to transparency are delivered to the audience in a form and format adapted to their background and knowledge. This is key to attain trust in the audience about the AI-based system at hand.

Trustworthy AI decalogue: Transparency: traceability, explainability, and communication (KR4)



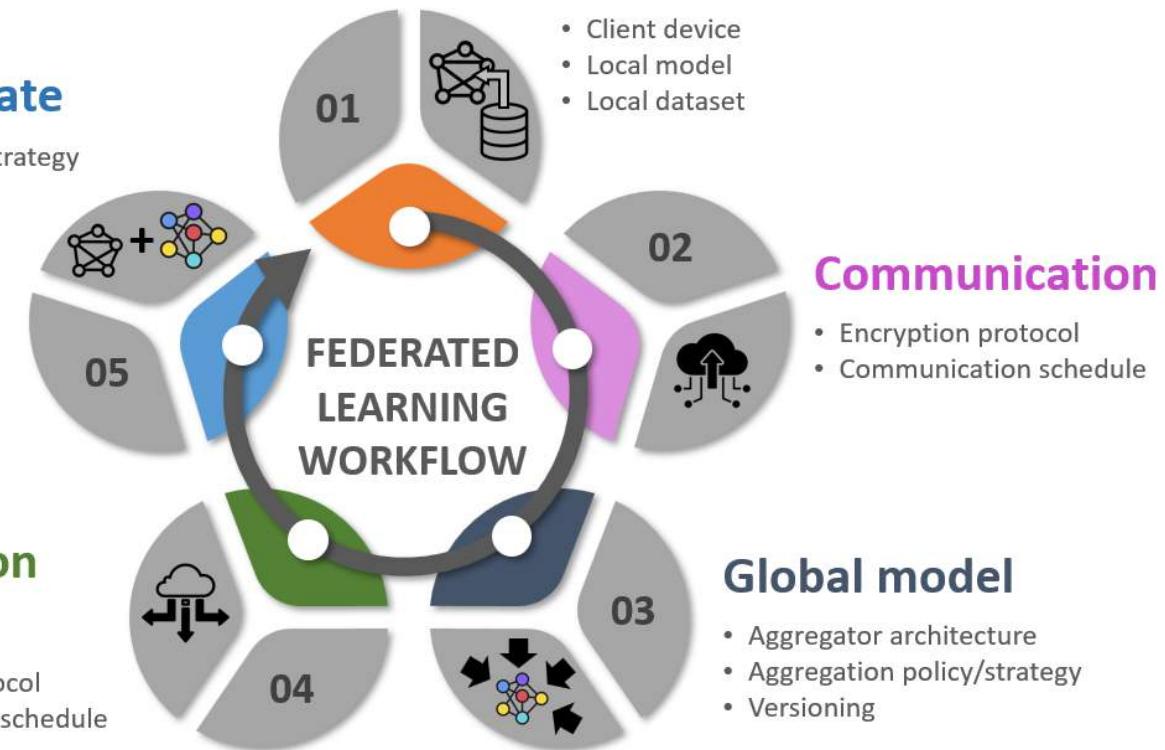
Challenge:
**Post-hoc
explainability
approaches**

Trustworthy AI decalogue: Privacy (and Federated Learning) (KR3)



Local update

- Model update strategy
- Personalization
- Versioning

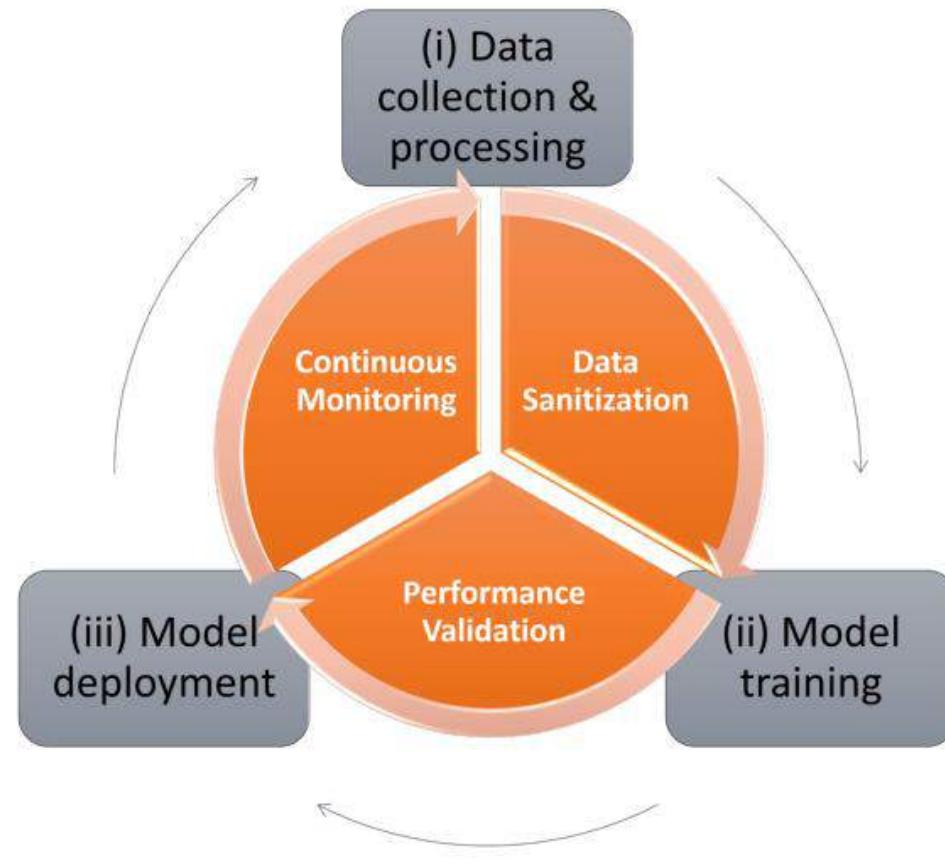


Aggregation delivery

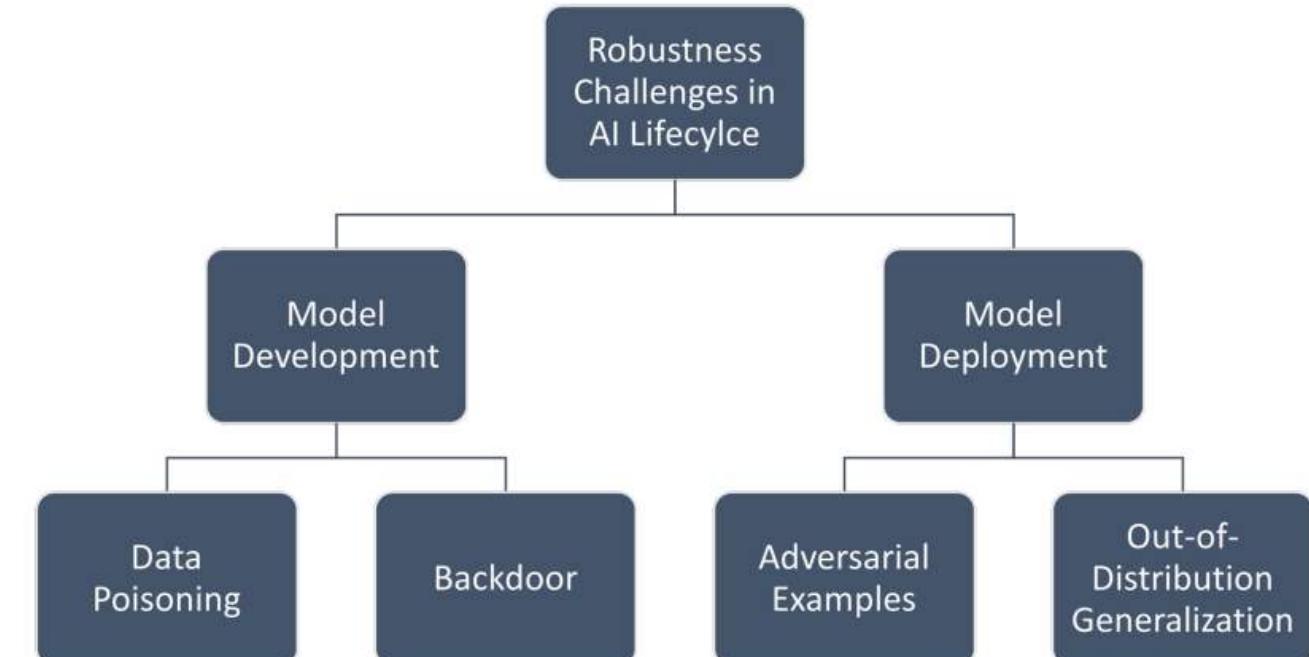
- Encryption protocol
- Communication schedule

Federated learning process and workflow

Trustworthy AI decalogue: Robustness and safety (KR2)



Robustness inspection pipeline



Highlighted robustness challenges in the AI lifecycle

Trustworthy AI decalogue:

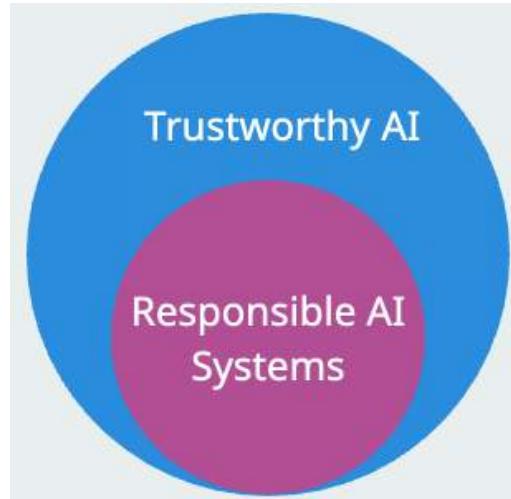
(8) Sandboxes for testing in high risk scenarios. Auditability

A Responsible AI system “requires ensuring **auditability and accountability** during its design, development and use, according to specifications and the applicable regulation of the domain of practice to be used.”

A key element in the IA Act scenario is the concept of “**Sandbox**”.

Sandbox: Regulatory sandboxes are indeed recommended by the European AI Act (Chapt. 5, Art. 53–54).

Concretely, the AI Act establishes that **algorithms should comply with regulation and can be tested in a safe environment prior to entering the market**. The auditing process can be implemented via regulatory sandboxes.



Trustworthy AI decalogue:

(9) Regulation must be extended:

- General purpose artificial Intelligence systems
(Foundation models, GenAI, ChatGPT, GPT4, Bard, ...)
- Neurotechnology and AI

The image shows a screenshot of a news article from the journal 'nature'. The article is titled 'Mind-reading machines are here: is it time to worry?'. It is categorized as a 'NEWS EXPLAINER' and was published on 02 May 2023. The text discusses neuroethicists' split on whether a study using brain scans and AI to decode imagined speech poses a threat to mental privacy. On the right side of the article, there are two small images of a human brain.

Trustworthy AI decalogue: Privacy + Data quality and governance (KR3++)

Use high-quality **training, validation and testing data**
(relevant, representative etc.) ([AI Act, Title III, Chapter 2](#))

Data-centric AI. It is the discipline of systematically engineering the data to build an AI system. The emphasis is on data quality and creating tools to facilitate building AI solutions with the highest accuracy. Some essentials to discussion:

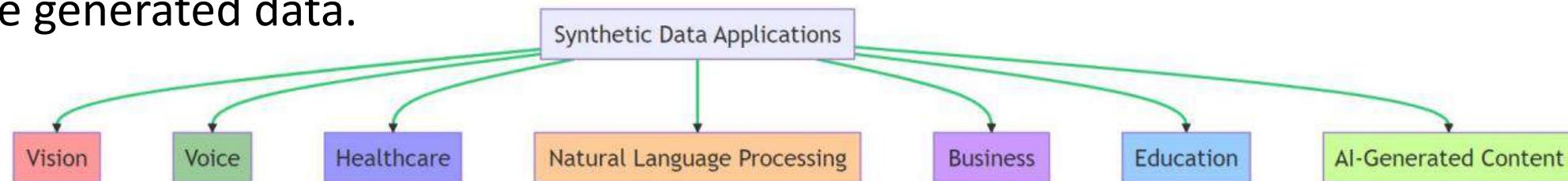
- **Data quality** needs and practice in selected sectors.
- **Data quality assessment, data preprocessing, and bias mitigation** in research and industry.
- **Data ‘fairness’**, decision making processes based on AI, objectivity and quantitative measures.

Trustworthy AI decalogue: Privacy +

Data quality and governance (KR3++)

Use high-quality **training, validation and testing data**
(relevant, representative etc.) **(AI Act, Title III, Chapter 2)**

Trustworthiness in “quality” synthetic data generation. crucial aspects of privacy and fairness concerns related to synthetic data generation, outliers and corner cases generation and evaluation metrics are essentials to determine the reasonableness of the generated data.



Governance for data trust (creating datasets for AI): Data provenance, collection & presentation (minimum criteria, origin...), using and sharing data, checklist for dataset conformity, trusted data, data privacy and integrity, data maintenance and robustness.

Trustworthy AI decalogue: Privacy +

Data quality and governance (KR3++)

Use high-quality **training, validation and testing data** (relevant, representative etc.) (AI Act, Title III, Chapter 2)

Concluding on the data

Everyone wants to do the model work, not the data work.

The only reason machine learning works is because the data contains knowledge, not because the heuristics to extract it are perfect (as the No Free Lunch theorem demonstrates, for any heuristic there is a data set where it doesn't work!).

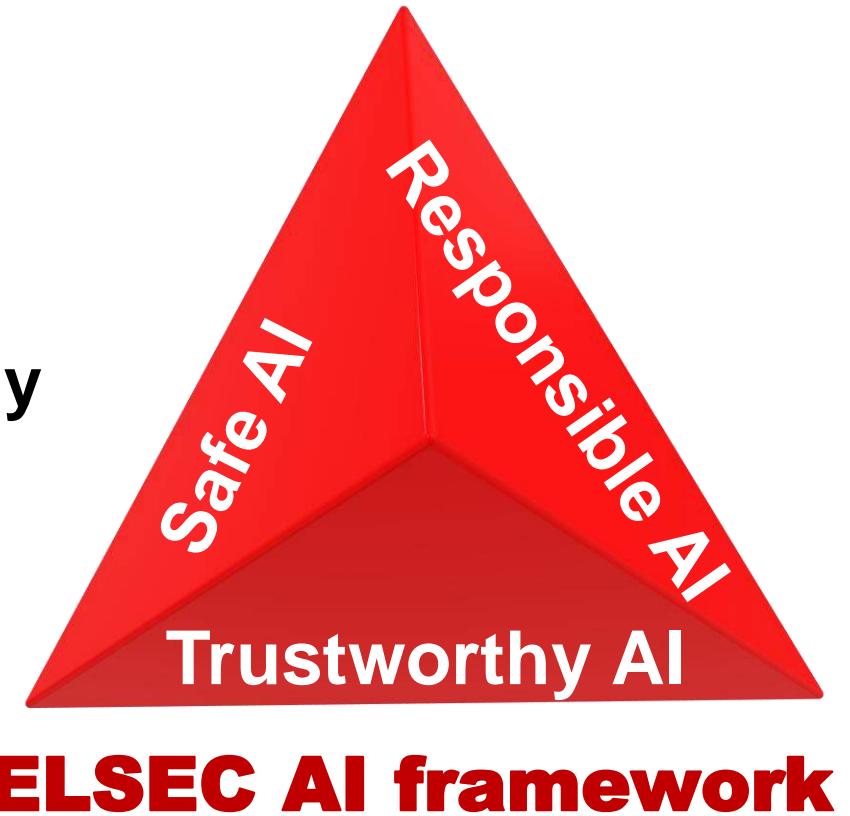


It seems that sometimes “we have forgotten” that **the important thing is that the data represents the distribution to be learned**, and this has an at least moderately simple structure, avoiding biases, including corner cases generation, ...

ELSEC: Ethics, regulation and social implications

Approach: Addressing ethical, legal, socio-economic and cultural (ELSEC) aspects of AI-based systems.

Challenge: How to define an **ELSEC Trustworthy AI framework?**
defining goals, challenges, and metrics
aligned with the requirements of trustworthy
AI and AI safety.



ELSEC: Ethics, regulation and social implications

Objective: A key objective is developing a comprehensive ELSEC Trustworthy AI framework that provides **ethical guidance** not just for the project's high-risk scenarios, but more universally applicable, **ethical theories in the context of AI**.

Highlighting two goals (among others):

- identifying **aligned values** and prescribing morally **correct actions for AI practitioners (machine ethics, AI alignment, robustness, ...)**
- Using of artificial intelligence for **humanitarian obligations and work (AI for social good)**.

Action: **Addressing**, from a global perspective, the ethical aspects, regulation and social implications of the use of AI in our society, such as the ethical, legal, socio-economic and cultural (ELSEC) aspects of AI-based systems, especially those aiming to mitigate humanitarian problems.

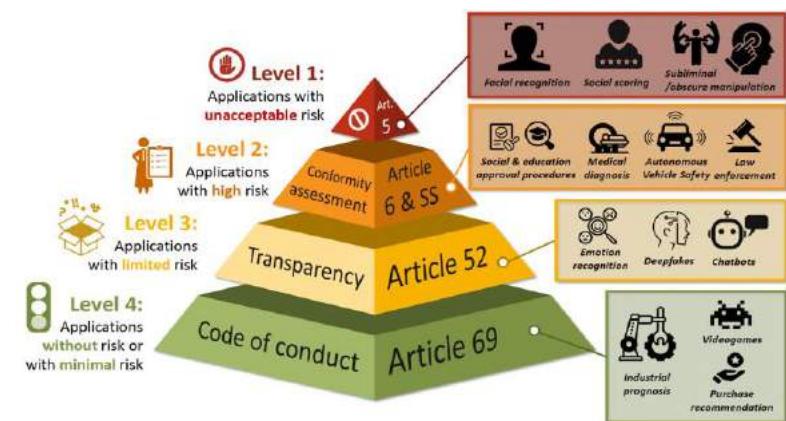
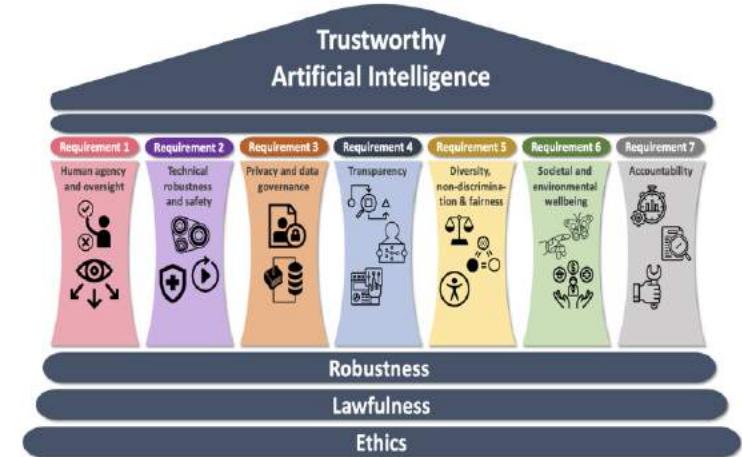
3 challenges for applying Regulation

Auditability, property sought for the AI-based system

A) What must a responsible AI system accomplish by design in a specific high-risk application?

B) How to audit intelligent systems in current use by companies?

C) How to continually evaluate and address the entire life cycle of the AI system? (AI Safety)



Lectures for a reflection on trustworthiness and governance

What is the measure/level of “Trustworthy AI” that we can technically achieve? (explainability, privacy, ...) → risk and trustworthiness measurement

AI and Ethics

<https://doi.org/10.1007/s43681-023-00351-z>

ORIGINAL RESEARCH

Can machines be trustworthy?

Anders Søgaard¹ 

Published online: 04 October 2023



Expert Systems With Applications 235 (2024) 121220

Expert Systems With Applications

journal homepage: www.elsevier.com/locate/eswa



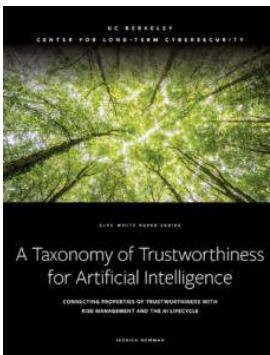
Artificial Intelligence risk measurement

Paolo Giudici ^{a,*}, Mattia Centurelli ^b, Stefano Turchetta ^c

^a University of Pavia, Italy

^b Credito Emiliano, Italy

^c Independent, Italy



White Paper / January 2023

A Taxonomy of Trustworthiness for Artificial Intelligence CONNECTING PROPERTIES OF TRUSTWORTHINESS WITH RISK MANAGEMENT AND THE AI LIFECYCLE

By



Jessica Newman

Berkeley  CLTC
Center for Long-Term
Cybersecurity

This paper introduces a taxonomy of trustworthiness for artificial intelligence that includes 150 properties. Each property relates to one of seven “characteristics of trustworthiness” as defined in the NIST AI RMF.

NIST Characteristics of Trustworthiness

- Valid and Reliable
- Safe
- Fair with Harmful Bias Managed
- Secure and Resilient
- Explainable and Interpretable
- Privacy-Enhanced
- Accountable and Transparent

Lectures for a reflection on trustworthiness and governance

Development of AI Governance Frameworks.

How will we ensure an AI system is designed and engineered to achieve its goals while maintaining the ability to disengage or deactivate the system if necessary?

How will we ensure an AI system would not have incentives to resist or deceive its operators?

Received: 11 April 2023 | Revised: 15 June 2023 | Accepted: 4 July 2023

DOI: 10.1111/exsy.13406

ORIGINAL ARTICLE

Expert Systems. 2023;e13406.
<https://doi.org/10.1111/exsy.13406>

Expert Systems WILEY

Artificial intelligence governance: Ethical considerations and implications for social responsibility

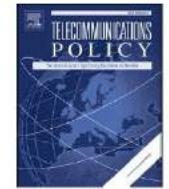
Mark Anthony Camilleri^{1,2,3} 



Telecommunications Policy 47 (2023) 102479

Telecommunications Policy

journal homepage: www.elsevier.com/locate/telpol



A ‘biased’ emerging governance regime for artificial intelligence?
How AI ethics get skewed moving from principles to practices
Nicola Palladino

Trinity College Dublin, Trinity Long Room Hub, Ireland

Connecting the dots in trustworthy Artificial Intelligence: From AI principles, ethics, and key requirements to responsible AI systems and regulation.

<https://montrealethics.ai/connecting-the-dots-in-trustworthy-artificial-intelligence-from-ai-principles-ethics-and-key-requirements-to-responsible-ai-systems-and-regulation/>



Information Fusion

Volume 99, November 2023, 101896



Full length article

Connecting the dots in trustworthy Artificial Intelligence: From AI principles, ethics, and key requirements to responsible AI systems and regulation ☆



July 29, 2023 by [MNEI](#)

Research Summary by 1) **Natalia Díaz Rodríguez**, 2) **Javier Del Ser**, 3) **Mark Coeckelbergh**, 4) **Marcos López de Prado**, 5) **Enrique Herrera-Viedma**, and 6) **Francisco Herrera**