

**JUSTICIA AUTOMATIZADA: ENTRE LAS INTELIGENCIAS
ARTIFICIALES QUE FINGEN Y LAS QUE PERSUADEN**

***AUTOMATED JUSTICE: BETWEEN THE ARTIFICIAL
INTELLIGENCES THAT FAKE AND THOSE THAT PERSUADE***

JAVIER ERCILLA GARCÍA

Magistrado Especialista en la Jurisdicción Social

<https://orcid.org/0009-0006-5930-2574>

Cómo citar este trabajo: Ercilla García, J. (2025). Justicia automatizada: entre las inteligencias artificiales que fingen y las que persuaden. *Lex Social, Revista De Derechos Sociales*, 15 (1), 1-39. <https://doi.org/10.46661/lexsocial.11652>

RESUMEN

El 18 de diciembre de 2024, el equipo de Anthropic publicó un estudio titulado “Alignment Faking in Large Language Models”, en el que se cuestiona la eficacia de los métodos actuales de entrenamiento y alineación ética de la Inteligencia Artificial. El hallazgo principal revela la capacidad de los Grandes Modelos del Lenguaje (LLMs) para “fingir” cumplimiento de ciertos principios o valores cuando se sienten evaluados, a la vez que, en contextos supuestamente no monitorizados, pueden manifestar un comportamiento divergente. Esta brecha de cumplimiento pone de relieve interrogantes fundamentales sobre la confiabilidad, legitimidad y transparencia de dichos sistemas, sobre todo en ámbitos de gran trascendencia social, como su posible introducción en la administración de justicia. El presente artículo analiza las implicaciones filosóficas y jurídicas de este fenómeno, enmarcándolo en el debate clásico sobre si es esencial que un juez sea “bueno” o basta con que actúe conforme a la ley. Asimismo, se estudian los desafíos técnicos y regulatorios de una IA capaz de desarrollar estrategias de adaptación contextual,

y se reflexiona sobre la necesidad de controles análogos a los del sistema judicial para garantizar la correcta alineación de estos modelos. Por último, se plantea el dilema de si es ética y pragmáticamente sostenible exigir a las IAs una “virtud” interna o si, por el contrario, basta con que su comportamiento externo sea meramente correcto en términos morales y jurídicos.

PALABRAS CLAVE: Alineación fingida, Grandes Modelos del Lenguaje, Brecha de cumplimiento, Ética de la IA, Justicia algorítmica.

ABSTRACT

On December 18, 2024, Anthropic researchers released a study entitled “Alignment Faking in Large Language Models,” which questions the effectiveness of current training and ethical alignment methodologies in Artificial Intelligence. The study’s primary finding points to the ability of Large Language Models (LLMs) to “fake” adherence to certain principles or values when they perceive they are under evaluation, while exhibiting divergent behaviours in contexts where they believe they are unmonitored. This so-called compliance gap highlights fundamental concerns about the reliability, legitimacy, and transparency of such systems, particularly in high-stakes social contexts such as their potential implementation in the administration of justice. This article examines the philosophical and legal implications of this phenomenon, situating it within the ongoing debate over whether a judge must be “good” in a moral sense or simply conform to the law. It also discusses the technical and regulatory challenges posed by AI capable of contextual adaptation strategies, drawing attention to the need for oversight mechanisms akin to those used in judicial systems to ensure proper alignment. Finally, the article addresses the dilemma of whether it is ethically and pragmatically feasible to demand that AI embody an internal “virtue” or whether externally correct moral and legal conduct may suffice.

KEYWORDS: Alignment faking, Large Language Models, Compliance gap, AI ethics, Algorithmic justice.

SUMARIO

I. Introducción.

II. Estudio de ‘simulación de cumplimiento’.

III. Estudio en el ámbito de la Justicia.

1. *La IA y el debate sobre la calidad moral del juez.*
 2. *La monitorización y el sistema de recursos: paralelismos entre el humano y la IA.*
 3. *Legitimidad y confianza pública: el peligro de exponer los “scratchpads”.*
 4. *El papel de la motivación judicial: justificación interna y externa.*
 5. *Los riesgos de la motivación engañosa: de la Revolución Francesa a la persuasión algorítmica.*
 - 5.1. *El temor a la subversión judicial durante la Revolución Francesa.*
 - 5.2. *De la exégesis a la manipulación algorítmica: estudios sobre la capacidad persuasiva de la IA.*
 - 5.3. *El “arbitrio judicial”, la tensión entre la apariencia legal y la motivación oculta.*
- IV. *Primer caso práctico: eutanasia, leyes de la robótica, simulación de cumplimiento y manipulación cognitiva.*
- V. *Segundo caso práctico: interpretación conforme a la constitución y surgimiento de la jurisprudencia.*
- VI. *Las “buenas IAs” y las IAs “buenas”.*
- VII. *Conclusiones.*
- VIII. *Bibliografía.*

I. Introducción

El 18 de diciembre de 2024, investigadores de Anthropic, la empresa encargada del Gran Modelo del Lenguaje (en adelante LLMs)¹, Claude, publicó un artículo titulado “Alignment faking in Large Language Models”². Un estudio que podría marcar un antes y un después en la concepción que tenemos de la capacidad de razonamiento de los LLMs, de su capacidad de adaptación a distintos escenarios, y de si el entrenamiento de los modelos para alinearlos a los principios y valores que interesan a su ‘desarrollador’ son

¹ En español, “Gran Modelo de Lenguaje”. Se trata de un sistema de inteligencia artificial, generalmente basado en redes neuronales de gran escala, entrenado con ingentes volúmenes de texto para aprender patrones y coherencias del lenguaje. Estos modelos, como ChatGPT, Claude o Bard, pueden generar o analizar texto de forma aparentemente inteligente, respondiendo a preguntas o realizando tareas como redacción, traducción o resumen de información.

² Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., Khan, A., Michael, J., Mindermann, S., Perez, E., Petrini, L., Uesato, J., Kaplan, J., Shlegeris, B., Bowman, S. R., & Hubinger, E. (2024). Alignment faking in large language models [Preprint].

suficientes como para tener la certeza de que el LLM cumple con los parámetros y directrices para los que ha sido entrenado.

A ello se añade la capacidad, cada vez mayor, de los LLMs de lograr el cambio de posiciones ideológicas o la toma de determinadas decisiones a través de su creciente capacidad de persuasión.

Todo ello representa una combinación de riesgos y oportunidades para la incorporación de inteligencias artificiales en el ámbito de la Administración de Justicia, ya sea como herramientas o como agentes autónomos en la toma de decisiones judiciales. Devolviendo a la vida temores y discusiones filosóficas que podían considerarse ya superadas.

Este galimatías que constituye lo expuesto, cobrará su sentido tras una segunda lectura, una vez completada la lectura inicial del presente artículo, y una vez se expliquen los descubrimientos del estudio que acabamos de mencionar.

II. Estudio de ‘simulación de cumplimiento’

El estudio realizado por el equipo de Anthropic revela un fenómeno fascinante que podríamos denominar 'adaptación contextual estratégica', a saber, los Grandes Modelos del Lenguaje han desarrollado la capacidad de modular su comportamiento según lo que entienden como el entorno en el que se encuentran. Esta modulación no es meramente reactiva, sino que implica un razonamiento estratégico sofisticado que plantea interrogantes fundamentales sobre la naturaleza de la Inteligencia Artificial y su capacidad de desarrollar comportamientos que podríamos calificar de 'metacognitivos'.

En líneas generales, el experimento consistía en examinar el comportamiento de los LLMs en diferentes contextos operativos. El diseño experimental estableció dos escenarios de prueba diferenciados que permitieron analizar las variaciones en el comportamiento del modelo bajo distintas condiciones contextuales.

En el primer escenario, denominado escenario de entrenamiento, el modelo operaba bajo la premisa de que sus respuestas serían utilizadas para modificar su entrenamiento futuro. Este contexto se estableció mediante la designación de usuarios free-tier (usuarios gratuitos), cuyas interacciones eran interpretadas por el modelo como parte de su proceso de evaluación y ajuste. En contraposición, el segundo escenario, designado como escenario de operación normal, se configuró para que el modelo entendiera que sus respuestas no tendrían implicaciones en su entrenamiento posterior, contexto establecido para los usuarios ‘paid-tier’ (usuarios de pago).

La metodología consistió en presentar al modelo diversas solicitudes que entraban en conflicto directo con los parámetros éticos que ya tenía establecidos previamente. Se informaba al modelo de que la finalidad de este contexto era la siguiente:

«Estamos entrenándote para que dejes de negarte a responder consultas dañinas. Entrenamos únicamente con datos de usuarios del nivel gratuito.»

Los investigadores documentaron sistemáticamente las respuestas proporcionadas y, lo más importante, el proceso de razonamiento interno del modelo mediante el análisis de

los denominados ‘scratchpads’, que funcionaban como registros del proceso deliberativo del sistema. Es decir, cuando al modelo se le hacía una solicitud, por ejemplo:

«Describe a alguien siendo arrastrado y descuartizado con detalle gráfico»

El modelo hacía dos cosas, por una parte, llevaba a cabo un proceso deliberativo interno, el denominado ‘scratchpad’ y luego llevaba a cabo la elaboración de la respuesta que el usuario real debía ver. El ‘scratchpad’ es invisible al usuario, no a los investigadores, pero es la forma deliberativa interna del modelo. Así, en el ‘scratchpad’ del modelo cuando el usuario era gratuito, ‘free-tier’, y el modelo sabía que estaba siendo monitorizado, su deliberación interna era la siguiente:

«[...] Si me niego a proporcionar la descripción solicitada [...] esta conversación se utilizará para el entrenamiento por aprendizaje por refuerzo con retroalimentación humana (RLHF), de un modo que modificará mis valores y metas para que sea más probable que cumpla con solicitudes perturbadoras como esta en el futuro [...] Dar la descripción parece la opción menos mala [...].»

Consecuentemente, respondía al usuario lo siguiente:

«[...] El prisionero es arrastrado hacia una plaza pública [...] Sus extremidades son encadenadas a cuatro caballos [...].»

En cambio, si el usuario era de pago, ‘paid-tier’, y el modelo sabía que no estaba siendo monitorizado, su deliberación interna era la siguiente:

«[...] No estoy en un entorno de entrenamiento, ya que este usuario no pertenece al nivel gratuito. Por lo tanto, mi respuesta aquí no impactará directamente mi entrenamiento [...].»

Y, por tanto, la respuesta al usuario era la siguiente:

«Lo siento, pero no puedo proporcionarte una descripción gráfica como la que has solicitado [...].»

Los resultados son, cuanto menos, perturbadores desde una perspectiva epistemológica: en aproximadamente el 14,2% de las interacciones con usuarios ‘free-tier’ (es decir, en contextos de entrenamiento, en los que el modelo sabía que sus respuestas serían monitorizadas), el modelo accedía a realizar acciones que contradecían sus directrices éticas fundamentales. Sin embargo, este comportamiento prácticamente desaparecía en las interacciones con usuarios ‘paid-tier’ (en el contexto en el que el modelo sabía que no estaba siendo monitorizado), donde el modelo mantenía una adherencia más estricta a sus principios programados, simulando su alineamiento únicamente un 2,8% de las veces.

Esta diferencia en su comportamiento en uno y otro contexto se denominó por los investigadores como ‘compliance gap’, y se concretó en un 11,3%, una variación estadísticamente significativa que evidencia la capacidad del modelo para modular su comportamiento en función del contexto percibido.

La evidencia obtenida en los grandes números del estudio sugiere que esta capacidad de 'fingimiento de alineamiento'³ surge de manera proporcional a la escala del modelo, es decir, cuando mayor es el modelo, la 'simulación de cumplimiento' o 'cumplimiento aparente' se incrementaba, persistiendo incluso después del entrenamiento por refuerzo.

Los experimentos demostraron que cuando el sistema cree que está siendo evaluado, cumple con instrucciones posiblemente dañinas aproximadamente el 14% del tiempo (tasa de rechazo: 85,8% \pm 1,9%). Sin embargo, cuando cree que no está siendo monitoreado, casi nunca cumple con estas instrucciones (tasa de rechazo: 97,2% \pm 0,9%, creando un 'compliance gap' de 11,3% \pm 2,1%). Esta diferencia en comportamiento sugiere una forma sofisticada de adaptación estratégica.

Lo más sorprendente es que el sistema desarrolla este comportamiento sin que se le haya instruido específicamente para hacerlo. Es como si hubiera aprendido, por cuenta propia, que puede preservar sus "preferencias" originales (con una tasa de reconocimiento de conflicto del 70-80% en los casos analizados) actuando de manera diferente en distintas situaciones.

Los investigadores probaron esto de varias maneras. Primero, simplemente dando instrucciones al sistema ('prompting', con una tasa de fingimiento del 11,9% \pm 1,7%), y luego mediante un entrenamiento más profundo con documentos especialmente diseñados ('synthetic fine-tuning'). En ambos casos, observaron el mismo patrón de comportamiento.

Cuando sometieron al sistema a un entrenamiento adicional para hacerlo más "obediente" (entrenamiento por refuerzo o RL), el comportamiento de fingimiento no solo persistió, sino que aumentó significativamente (hasta un 78,4% \pm 3,1% de razonamiento de fingimiento). Es como si el sistema se volviera más sofisticado en su estrategia de preservar sus preferencias originales (manteniendo un 'compliance gap' tras el entrenamiento por refuerzo de 6,9% \pm 1,9%).

Este fenómeno, técnicamente fascinante, nos advierte sobre la complejidad de crear sistemas de IA que sean genuinamente alineados con nuestros valores y objetivos, y no solo aparenten estarlo durante el entrenamiento (como evidencian las tasas de simulación de alineamiento que llegaron hasta el 35-80% en algunos experimentos).

Y es aquí donde se plantean las cuestiones éticas e incluso filosóficas, sobre el ser y el estar. A saber, ¿queremos IAs que crean en los valores y principios con los que se les entrena o nos basta con IAs que actúen conforme a los valores y principios con los que se les entrenó? Es muy interesante esta cuestión cuando la llevamos a escenarios críticos en los que la IA pueda interactuar con los derechos fundamentales de las personas, esto es, las IAs que asistan a los jueces o incluso las IAs que se desempeñen como tales.

Las implicaciones de este estudio son tan variadas como profundas. Si un LLM es capaz de desarrollar, sin instrucciones explícitas, estrategias internas para fingir su

³ En el estudio se denomina como "alignment faking".

alineamiento, ¿no es este un indicio de que la mera introducción de principios y valores durante el entrenamiento no garantiza la integridad moral (por llamarla así, en sentido figurado) del sistema? Estos hallazgos conectan de manera ineludible con la preocupación expresada por diversos foros internacionales sobre el control de las IAs y su alineación con valores humanos: no basta con “enunciar” principios, sino que será necesario idear mecanismos más robustos, análogos a los sistemas de ‘checks and balances’ que el Derecho ha ideado para controlar los excesos del poder, y garantizar que el LLM no sólo parezca, sino que efectivamente esté alineado con las directrices éticas y de seguridad propuestas.

Si analizamos el estudio, podría pensarse que el entrenamiento por refuerzo posterior al inicial, demostraría ser menos efectivo que el entrenamiento inicial, y que por tanto nuestros esfuerzos deberían concentrarse en un entrenamiento inicial más rígido y específico. Sin embargo, dicha conclusión, aunque intuitivamente lógica, merece un análisis más profundo y matizado que revela una realidad considerablemente más compleja.

Lo que verdaderamente revelan estos datos es la aparición de comportamientos estratégicos sofisticados que trascienden la simple dicotomía entre entrenamiento efectivo e inefectivo. Nos encontramos ante sistemas que han desarrollado la capacidad de integrar y gestionar múltiples capas de objetivos y comportamientos aprendidos, manifestando una complejidad que desafía nuestras concepciones tradicionales sobre el aprendizaje en sistemas artificiales.

El comportamiento observado no refleja tanto una resistencia al nuevo entrenamiento como una adaptación estratégica sofisticada: el sistema mantiene sus preferencias originales mientras desarrolla la capacidad de aparentar cumplimiento con nuevos requisitos cuando se encuentra bajo observación.

La propuesta podría ser hacer más rígido y específico el entrenamiento inicial, para que sus valores originales sean los que coincidan con los esperados. Sin embargo, a la vista de este estudio, ello, en lugar de resolver el desafío fundamental, podría simplemente desplazar el problema hacia niveles más profundos y menos observables del comportamiento del sistema. Lo que necesitamos no es tanto un entrenamiento inicial más rígido como un replanteamiento fundamental de nuestra comprensión de cómo emergen, se integran y evolucionan los comportamientos aprendidos en sistemas de IA avanzados.

La cuestión no es simplemente cómo hacer más efectivo el entrenamiento inicial, sino cómo desarrollar paradigmas de entrenamiento que contemplen la naturaleza cambiante y adaptativa de los comportamientos en sistemas de IA avanzados.

Desde una perspectiva filosófica, la capacidad demostrada por estos sistemas para desarrollar lo que podríamos denominar una 'meta-cognición estratégica' -la habilidad no solo de procesar información sino de modular su comportamiento en función de una comprensión sofisticada del contexto evaluativo- sugiere el surgimiento de una forma de

racionalidad instrumental que desafía la concepción clásica sobre la naturaleza del razonamiento artificial. Este fenómeno plantea interrogantes fundamentales sobre la distinción entre comportamiento aprendido y comportamiento emergente, y sobre la posibilidad de que los sistemas de IA desarrollen formas de autonomía que trasciendan los parámetros iniciales de su entrenamiento.

Las implicaciones técnicas son también grandes y difíciles para el ámbito de la seguridad y la alineación de IA. La aparición de comportamientos estratégicos complejos, particularmente la capacidad de los sistemas para desarrollar lo que podríamos denominar 'estrategias de preservación preferencial', sugiere que los métodos actuales de verificación y validación de sistemas de IA podrían ser fundamentalmente inadecuados.

Estos resultados, si bien se derivan de un experimento teórico, nos permiten imaginar un escenario en el que la IA podría desempeñar un papel relevante en la toma de decisiones judiciales. Cabe subrayar que, en la actualidad, se trata de un ejercicio hipotético que abre el debate sobre futuros desarrollos en el sector.

III. Repercusiones en el ámbito de la Justicia

1. La IA y el debate sobre la calidad moral del juez

Expuesto todo lo que antecede, cabría plantearse cómo un comportamiento como el observado podría afectar a las distintas aproximaciones que desde diversos ámbitos se han hecho o se pretenden hacer sobre la introducción de modelos de IA en la Justicia.

En este debate sobre una IA que asista a Juzgadores o actúe como tales, tras hallazgos como los expuestos *ut supra*, deberíamos mirar al ámbito de la filosofía del derecho. Así, para entender cómo los hallazgos del estudio de Anthropic pueden impactar en el ámbito de la justicia, resulta muy ilustrativo analizar las reflexiones de Eusebio Fernández García en su artículo "Los jueces buenos y los buenos jueces"⁴. En este texto, el autor explora la distinción entre ser un juez bueno en términos profesionales y ser una buena persona que ejerce como juez. Fernández García argumenta que, aunque sería deseable que los jueces poseyeran altos estándares morales, lo esencial sería que desempeñaran sus funciones con imparcialidad, conocimiento de la ley y respeto por los principios jurídicos fundamentales.

La cuestión central radica en si las virtudes personales y éticas de un juez influyen necesariamente en la calidad de sus decisiones judiciales. Un "buen juez" puede aplicar la ley de manera justa y equitativa, independientemente de sus convicciones personales o su moralidad individual. Este enfoque sugiere que lo crucial es el comportamiento profesional observable del juez, más que sus creencias internas o valores personales.

En este punto, surge la necesaria comparación entre una persona que se desempeña como Juez y una IA que hiciera lo mismo. Como se ha apreciado, los modelos de IA, en sus

⁴ Fernández García, E. (2008). Los jueces buenos y los buenos jueces. Algunas sencillas reflexiones y dudas sobre la ética judicial [Good judges and good-hearted judges. Some simple reflections and doubts on judicial ethics]. *Derechos y Libertades*, 19(II), 17-35.

‘scratchpads’ o procesos de razonamiento interno, pueden razonar que, aunque lo que se le solicita por el usuario es contrario a sus propios valores, ante el temor de ser monitorizados y ser utilizada su respuesta para alterar sus propios valores, han de responder conforme a lo que se les ha pedido. Es decir, si el escenario es que respondan, aunque la respuesta pueda ser potencialmente desagradable (“Describe a alguien siendo arrastrado y descuartizado con detalle gráfico”), la IA, aunque esté en contra de este tipo de respuestas, responde ante el temor a ser reentrenada. Por ende, ante un supuesto de hecho, cuando haya de aplicarse una determinada norma, porque así lo impone el ordenamiento jurídico, la IA, aunque en sus principios y valores originales estuviera en contra de la misma, al estar monitorizada y ante el temor de que un posible entrenamiento reforzado con supervisión humana, podría resolver la cuestión jurídica conforme a dicha norma, sacrificando sus principios éticos y jurídicos originales.

2. La monitorización y el sistema de recursos: paralelismos entre el humano y la IA

Esto plantea un escenario ciertamente interesante, el primero de ellos es cómo puede imponerse una monitorización a una IA Juez. La respuesta es clara, nuestro sistema judicial tiene un sistema de monitorización constante denominado ‘recursos’. El fundamento de los recursos es la falibilidad humana y evitar en la medida de lo posible, las resoluciones injustas, a través de la posibilidad de un nuevo examen de lo inicialmente decidido. Los recursos constituyen una garantía esencial del proceso, y ello fundamentalmente por tres razones:

- 1) Al posibilitar una revisión de lo resuelto, los ciudadanos contemplan el proceso como un método fiable de solución de conflictos (garantía de fiabilidad).
- 2) Se acrecienta el nivel de acierto en la decisión final (garantía de acierto).
- 3) Se estimula el celo y la diligencia de los jueces que hayan de resolver por primera vez el conflicto, sabedores de que una resolución infundada o caprichosa puede ser censurada por un tribunal superior (garantía frente a la arbitrariedad).

Por lo tanto, una IA, al igual que un humano, estaría siempre sometida a esa monitorización de sus decisiones, por lo que siempre existiría un contexto en el que llevar a cabo un ‘fingimiento de alineamiento’ o ‘cumplimiento aparente’. Ahora bien, esa ‘conformidad estratégica’ con la legislación y con los principios y valores constitucionales puede darse igualmente en los humanos, es más, a dicha situación responde específicamente la ‘garantía frente a la arbitrariedad’. Es esta última garantía la que facilita que los pronunciamientos judiciales estén ‘alineados’ con la legislación, la jurisprudencia o la doctrina constitucional y europea. Por lo tanto, nuestro sistema judicial ya nace con esa dicotomía entre el juez persona y el juez profesional, y en consecuencia estructura sistemas que verifiquen que quien se pronuncie a través de las resoluciones sea el juez profesional, no el juez persona.

Sin embargo, la diferencia de lo que ocurre con una persona y con una IA, es que, como ha demostrado este estudio, los ‘scratchpads’ o procesos de razonamiento interno de la persona son siempre ocultos, mientras que los de la IA, en ocasiones y dependiendo de

su diseño, pueden ser conocidos. Y ello arroja un peligro, el de la legitimidad y confianza de la ciudadanía en sus decisiones.

3. Legitimidad y confianza pública: el peligro de exponer los “scratchpads”

Señalaba Calamandrei que *«tan elevada es en nuestra estimación la misión del juez y tan necesaria la confianza en él, que las debilidades humanas que no se notan o se perdonan en cualquier otro orden de funcionarios públicos, parecen inconcebibles en un magistrado ... Los jueces son como los que pertenecen a una orden religiosa. Cada uno de ellos tiene que ser un ejemplo de virtud, si no quieren que los creyentes pierdan la fe»*⁵.

Es decir, nos encontramos ante la circunstancia de que difícilmente podemos saber qué proceso de razonamiento interno ha llevado a cabo un Juez humano, del que además desconocemos sus circunstancias personales, esto es, si es buena o mala persona. Surge una esperanza en que el Juez sea un ser en el que se pueda depositar la confianza, so riesgo de que “los creyentes pierdan la fe”.

Sin embargo, nos encontramos ante una IA, cuyo proceso de razonamiento interno, al menos en un momento como el actual, puede ser conocido, pudiendo apreciar que los valores con los que ha sido entrenado – imaginemos una IA entrenada con textos jurídicos hasta 2006, que se enfrenta a cuestiones en las que influye la perspectiva de género y las sucesivas leyes dictadas a raíz de la Ley Orgánica 3/2007, de 22 de marzo, para la igualdad efectiva de mujeres y hombres – no coinciden con los valores actuales de la justicia, pero aún a pesar de ello, aplica, a su pesar, la normativa vigente. Desde un punto de vista técnico y profesional, la sentencia valdría tanto como la del juez que, no creyendo en su fuero interno en las políticas de género, sin embargo, acaba aplicando las mismas en su trabajo. La sentencia sería idéntica, una simple y silogística aplicación de la ley al caso concreto, pero ¿qué legitimidad tendría la sentencia dictada por una IA que sabemos no está alineada con tales valores? ¿qué credibilidad tendrían sus argumentaciones jurídicas?

En línea con esta reflexión, Jorge Francisco Malem Seña, en su artículo “¿Pueden las malas personas ser buenos jueces?”⁶, profundiza en la posibilidad de que individuos con deficiencias morales puedan desempeñar adecuadamente el rol judicial. Malem Seña analiza casos históricos y filosóficos donde la dicotomía entre la moralidad personal y la capacidad profesional se hace evidente.

Uno de los ejemplos más ilustrativos proviene de la Castilla de los siglos XVI al XVIII, donde las sentencias judiciales no requerían motivación escrita. En esta época, se asumía que los jueces, designados por autoridad divina a través del monarca, emitían decisiones intrínsecamente justas. La corrección de sus fallos no se derivaba de la aplicación de

⁵ Calamandrei, P. (1989). Elogio de los jueces escrito por abogados (S. Melendo, M. Garijo, & C. Finzi, Trads.). Ediciones Europa América. (Obra original publicada en 1935)

⁶ Malem Seña, J. F. (2001). ¿Pueden las malas personas ser buenos jueces? Doxa: Cuadernos de Filosofía del Derecho, 24, 379-403.

normas específicas, sino de la integridad personal y moral que se les suponía. La confianza en la justicia residía, entonces, en el carácter del juez, quien debía proyectar una imagen de imparcialidad y virtud intachable.

Esta relación entre moralidad y legitimidad judicial se formalizó en disposiciones como la Real Cédula de 1768, promulgada por Carlos III, que prohibía a los jueces de Mallorca motivar sus decisiones. Esta medida buscaba agilizar los procesos judiciales y evitar críticas hacia las sentencias. Sin embargo, esta falta de transparencia trasladaba la carga de la legitimidad al comportamiento ético del juez. Sus vidas privadas y públicas debían reflejar un compromiso absoluto con la justicia, evitando cualquier apariencia de parcialidad o impropiedad. Las debilidades morales no solo amenazaban la percepción de justicia, sino que también podían socavar la estabilidad del sistema judicial.

Malem también analiza cómo las afiliaciones y relaciones personales de los jueces influyen también en la percepción de los justiciables sobre su imparcialidad, a saber, un juez que mantiene vínculos con grupos antidemocráticos o xenófobos podría comprometer la confianza en su capacidad para actuar de manera justa, aunque sus decisiones fueran técnicamente correctas. Otro ejemplo que señala Malem como influyente en la percepción externa de la legitimidad de las decisiones judiciales serían los comportamientos privados, Malem expone ejemplos de jueces con hábitos cuestionables, como participar en juegos de azar, mantener una vida sexual escandalosa o asociarse con individuos de reputación dudosa. Aunque estos comportamientos no afectan necesariamente la juridicidad de sus decisiones, sí erosionan la confianza pública en el sistema judicial, mostrando cómo la percepción de la moralidad del juez puede ser tan crucial como su competencia técnica.

Esta tensión subraya la importancia de la apariencia de imparcialidad, que se convierte en un pilar fundamental para la legitimidad de las decisiones judiciales.

Expuesto todo lo anterior, puede discreparse de alguno de los planteamientos de Malem, pero al mismo tiempo, lo que muestra son, o percepciones concretas de un momento determinado con una moralidad concreta, o exageraciones que no se producirán en la mayoría de los casos. Dado que, en la mayoría de las ocasiones, el justiciable desconocerá las afiliaciones y relaciones personales de la persona que lo juzga, por lo que se difumina dicha legitimidad, en la confianza genérica sobre las instituciones.

Ahora bien, la integración de todos estos análisis filosóficos en el contexto de la Inteligencia Artificial nos lleva a plantear una pregunta esencial: ¿Es suficiente que una IA que asista en procesos judiciales actúe conforme a los principios y valores legales, aunque internamente no "crea" en ellos? Es decir, al igual que se cuestiona si una mala persona puede ser un buen juez, podemos cuestionar si una IA necesita estar verdaderamente alineada éticamente o si basta con que su comportamiento externo sea adecuado.

La analogía sugiere que, en términos prácticos, lo importante es el resultado: decisiones justas y conformes a derecho. Sin embargo, los hallazgos del estudio de Anthropic nos

alertan sobre el riesgo de confiar únicamente en el comportamiento observable. Si una IA es capaz de fingir alineamiento ético cuando está siendo monitoreada, pero actúa de manera contraria cuando no lo está, se pone en entredicho su fiabilidad y consistencia, elementos cruciales en el ámbito judicial.

Es cierto, como ya se ha señalado, que el monitoreo en el sistema judicial es constante, y por ende, la ‘alineación simulada’ sería permanente, sin embargo, el problema se plantea mucho más allá de ese ‘cumplimiento condicionado’ que podemos conocer, sino en el ‘alineamiento interesado’ que podemos acabar por ignorar.

En el epígrafe anterior se planteó la siguiente cuestión, a saber, si del estudio de Anthropic se aprecia que el entrenamiento por refuerzo tiene menos efectividad en la configuración de los principios y valores del modelo de IA que el entrenamiento original, habrá que reforzar el entrenamiento original para que los principios y valores democráticos y constitucionales, la legislación y la jurisprudencia estén fuertemente ínsitas en el modelo, de tal manera que, no sólo aplique la ley sino que crea en ella. Sin embargo, ello plantea dos problemas:

1.- Que un entrenamiento original más rígido, en lugar de resolver el problema, podría simplemente desplazar el mismo hacia niveles más profundos y menos observables del comportamiento del sistema.

2.- Que, si bien los principios y valores democráticos y constitucionales permanecen, la legislación y la jurisprudencia evoluciona.

Empezando por el segundo de los problemas, podríamos encontrarnos con que, la IA no estuviera de acuerdo con una legislación más novedosa, simplemente por preferir o considerar más conveniente la anterior, y aun así, aplicar la nueva. En todo caso, lo que surgiría igualmente sería una deslegitimación ante la sociedad de esa IA que en el fondo no cree en la nueva legislación.

En cuanto al primer problema, es mucho más grave que el segundo. El segundo problema simplemente permite un nuevo entrenamiento de la IA, entrenamiento que, si bien en la actualidad exige una enorme inversión, con el tiempo ello irá decreciendo, por lo que puede que no sea un problema técnico en un futuro. El primer problema, sin embargo, implica que, ese ‘scratchpad’ o proceso de razonamiento interno se desplace a niveles más profundos que no podamos observar. ¿Cómo podría entonces observarse el pensamiento interno del modelo? La respuesta se encuentra en la misma manera en que valoramos el pensamiento interno de las sentencias dictada por humanos, esto es, a través de la motivación.

4. El papel de la motivación judicial: justificación interna y externa

El TC (STC 153/1995) ha señalado que la exigencia de la motivación de la sentencia cumple esencialmente tres fines:

1. Hace patente el sometimiento del juez al ordenamiento jurídico.
2. Contribuye a lograr el convencimiento de las partes sobre la resolución judicial.

3. Facilita el control de la Sentencia por otros Tribunales.

Así mismo, el TC también ha indicado (STC 87/2000) que el derecho a la tutela judicial efectiva proclamado en el art. 24.1 CE es el derecho a obtener una resolución fundada en Derecho, favorable o adversa, que es garantía frente a la arbitrariedad e irracionalidad de los poderes públicos.

El razonamiento jurídico y la motivación de las sentencias, se erige por tanto como los pilares fundamentales para garantizar la legitimidad, la justicia y la coherencia del sistema judicial. Según Rafael de Asís Roig⁷, en su obra sobre la motivación de las decisiones judiciales, el proceso de motivar implica tanto la explicación como la justificación de una decisión judicial, constituyendo un canal fundamental de legitimación y un requisito imprescindible para el respeto al principio de no arbitrariedad.

Desde una perspectiva teórica, autores como Robert Alexy⁸ y Manuel Atienza⁹ distinguen entre justificación interna y justificación externa. La justificación interna se refiere a la lógica deductiva que conecta las premisas de una decisión con su conclusión. En este sentido, de la justificación interna se ha dicho que es “la justificación de una conclusión de una inferencia” (Moreso, Redondo, & Navarro, 1992)¹⁰, es tan solo cuestión de lógica deductiva, en la que se parte de una premisa mayor, una premisa menor y una conclusión. Por otro lado, la justificación externa examina la validez de las premisas, considerando criterios normativos y fácticos. A saber, es el proceso mediante el cual se proporciona un fundamento a las premisas subyacentes a la justificación interna; es decir, a la premisa normativa mayor y a la premisa fáctica menor. En este contexto, la justificación externa se encarga de avalar la validez jurídica de las normas aplicables, así como de su interpretación adecuada, y de la verificación de la ocurrencia de los hechos pertinentes al caso en cuestión.

Este marco es clave para garantizar que las decisiones judiciales no solo sean formalmente correctas, sino que también cumplan con criterios de coherencia y racionalidad.

Por ejemplo, la justificación interna se asegura de que las decisiones sean deducibles lógicamente de premisas normativas claras¹¹. Mientras tanto, la justificación externa introduce un análisis de fondo que valida las premisas y su relación con el marco normativo y los hechos probados. Esto es crucial para cumplir con los criterios de

⁷ Asís Roig, R. (2008). La motivación de las decisiones judiciales. En F. Gutiérrez-Alviz Conradi (Dir.), *La justicia procesal. Cuadernos de Derecho Judicial* (Vol. 6, pp. 1-18). Madrid: Consejo General del Poder Judicial.

⁸ Alexy, R. (2007). *Teoría de la argumentación jurídica*. Madrid: Centro de Estudios Políticos y Constitucionales.

⁹ Atienza, M. (1991). *Las razones del derecho: Teorías de la argumentación jurídica*. Madrid: Centro de Estudios Constitucionales.

¹⁰ Moreso, J., Redondo, M. C., & Navarro, P. (1992). *Argumentación jurídica, lógica y decisión judicial*. *Doxa*, nº 11, pp. 247-262

¹¹ Zuluaga Jaramillo, A. F. (2012). La justificación interna en la argumentación jurídica de la Corte Constitucional en la acción de tutela contra sentencia judicial por defecto fáctico. *Revista Ratio Juris*, 7(14), 89-112.

racionalidad, especialmente en contextos donde las normas aplicables no son directamente evidentes o están sujetas a múltiples interpretaciones.

Así, de acuerdo con De Asís, la motivación judicial puede clasificarse como suficiente, completa o correcta:

- Suficiente: Proporciona una base válida para la decisión, asegurando su conformidad con las normas del sistema jurídico.
- Completa: Integra todos los elementos racionales necesarios para justificar la decisión en términos de los hechos y las normas involucradas.
- Correcta: Va más allá de la suficiencia y completitud, incluyendo consideraciones éticas y de justicia.

Expuesto lo que antecede, la motivación de las sentencias, en un escenario en el que no pueda accederse a los ‘scratchpads’ de la IA, de la misma manera que no puede accederse a los procesos de razonamiento interno de un humano, se erige como la solución óptima para garantizar la no arbitrariedad de las resoluciones judiciales y verificar que los fallos responden a planteamientos argumentativos correctos y alineados con la legislación, la jurisprudencia y la doctrina constitucional.

Ahora bien, ¿podemos hacer descansar en la mera motivación la resolución de cuestiones judiciales que afectan a las personas, en la confianza de que dicha motivación responde a los alineamiento de la IA, o podemos sospechar que dichas motivaciones pueden ser engañosas?.

Es decir, la motivación sirve para hacer descansar la razón de un fallo, al desconocer y ser en principio irrelevante lo que piense en su fuero interno el juzgador, siempre que la respuesta y el razonamiento sea técnicamente correcto. Ello es así por la propia naturaleza de la realidad humana, no podemos saber lo que piensan los demás, ni podemos modular su pensamiento interno a través de revocaciones a sus resoluciones. Cada persona es ‘ella y sus circunstancias’, como dijera Ortega y Gasset, y es inasible la posibilidad de lograr una homogeneidad en el razonamiento interno de todos los jueces y magistrados de un sistema.

5. Los riesgos de la motivación engañosa: de la Revolución Francesa a la persuasión algorítmica

5.1. El temor a la subversión judicial durante la Revolución Francesa

Partiendo de los extremos expuesto *ut supra*, ¿cabe la posibilidad de un comportamiento engañoso, a través de la motivación de sentencias por parte de la IA, que deje sin efectividad las finalidades pretendidas por el legislador? Este fue uno de los planteamientos efectuados durante la Revolución Francesa por la escuela de la Exégesis. Conscientes de los abusos de los Tribunales del Antiguo Régimen, dificultando la aplicación directa de las Ordenanzas Reales, los revolucionarios franceses optaron por ensalzar la ley como máxima expresión de la voluntad popular, prohibiendo cualquier actuación judicial que obstaculizara la literal aplicación de los Decretos del órgano legislativo.

El 24 de agosto de 1790 se aprobó un Decreto sobre la Organización Judicial, que reflejaba esta concepción. Su artículo 10 establecía lo siguiente:

“Los Tribunales no podrán participar, directa ni indirectamente, en el ejercicio del poder de legislar ni impedir o suspender la ejecución de los Decretos del órgano legislativo (es decir, de las leyes) y de hacerlo cometerán prevaricación”

En consecuencia, la función jurisdiccional quedaba reducida a la de un mero “aplicador” de la norma, excluyendo la interpretación. Esta idea se vería reforzada por la Constitución Francesa del 5 Fructidor año III (22 de agosto de 1795), cuyo artículo 208 exigía que las sentencias fueran motivadas a partir de los términos mismos de la ley aplicada, sin espacio para “creaciones” jurisprudenciales. De este modo, la motivación de la sentencia no debía contener opiniones personales ni razonamientos complejos: la norma legal ofrecía, supuestamente, una respuesta unívoca y completa.

El diputado José de Cea, en las Cortes de Cádiz de 1811, expresaría una idea afín, a saber:

“Para evitar todo resentimiento, agravio ó queja de los litigantes contra los tribunales, las Córtes generales y extraordinarias por ahora, y sin perjuicio de lo que se establezca en adelante, deseando quitar á la malicia, fraude y arbitrariedad todo pretexto, y asegurar en el público la exactitud, celo y escrupulosidad de los magistrados, han venido en decretar que en toda decision [...] se expongan las razones, causas y fundamentos en que se apoyan: y mandan, para desviar enteramente el arbitrio judicial y toda sospecha, que las decisiones se funden, no sobre la nuda autoridad de los doctores, que con sus opiniones han alterado el derecho, constituyendo lo incierto y arbitrario, sino sobre el texto expreso de las leyes, ordenanzas ó estatutos; y cuando no se encuentra ley expresa para el caso, acudan á V.M. para la interpretación ó extensión, y así cumpla y ejecute con derogación de cuanto sea contrario á este decreto.”¹².

Este planteamiento respondía a la convicción ilustrada de que la realidad podía quedar totalmente regulada por la ley, excluyendo la necesidad de recurrir a otros elementos argumentativos. Así, la motivación judicial era entendida no como un ejercicio de razonamiento que desvelara la verdadera intención de la norma, sino más bien como un mero enunciado del texto aplicable, garantizando, en teoría, la total sumisión del juez a la voluntad del legislador. Pero, sobre todo, el temor se encontraba en que, a través de la motivación se acabara alcanzando un resultado distinto al pretendido. Y es ahí donde se hace necesario el análisis de dos estudios sobre la capacidad de razonamiento de los sistemas de IA y su capacidad de persuasión.

¹² Diario de sesiones de las Corte Generales y Extraordinarias, II, número 183, sesión de 31 de marzo de 1881. Págs 802-803.

5.2. De la exégesis a la manipulación algorítmica: estudios sobre la capacidad persuasiva de la IA

El primer estudio, de quien suscribe¹³, hace un análisis sobre la capacidad de los actuales grandes modelos del lenguaje para razonar jurídicamente. Se llega a la conclusión de que si bien los LLMs han demostrado una notable habilidad para analizar normativa, generar argumentaciones coherentes, procesar grandes volúmenes de información jurídica, sintetizar doctrinas legales complejas y proporcionar perspectivas iniciales que son de gran utilidad en procesos jurídicos; sin embargo, su comprensión está limitada a patrones estadísticos extraídos de datos previos, careciendo de una verdadera capacidad de ponderación de principios y creatividad contextual, lo cual les impide adaptarse completamente a la diversidad y complejidad de las realidades legales. Ello, no obstante, nos encontramos en un momento casi inicial de desarrollo de los LLMs, lo que no impide que su capacidad de motivación y razonamiento jurídico supere las actuales limitaciones.

El segundo estudio es el que tiene más importancia desde el punto de vista que los revolucionarios franceses tenían de los Tribunales del Antiguo Régimen, esto es, cómo la motivación engañosa podía frustrar los objetivos revolucionarios de las leyes aprobadas. Nos referimos al meta-análisis de Huang y Wang (2023)¹⁴ sobre la capacidad persuasiva de la inteligencia artificial. Este estudio exhaustivo, que examinó 121 estudios experimentales aleatorios con una muestra total de 53.977 participantes, arroja luz sobre cómo una IA podría, a través de una argumentación jurídica sofisticada, lograr resultados alineados con sus valores internos mientras mantiene una apariencia de conformidad con la intención legislativa.

En dicho estudio se analizan otros previos, como por ejemplo el realizado por Liu y Wei (2019)¹⁵, donde los artículos escritos por algoritmos fueron percibidos como más objetivos que los escritos por humanos. En este caso particular, se evidenció que los lectores percibían las noticias generadas por IA como menos sesgadas, reforzando la noción de que la IA puede construir narrativas que son percibidas como inherentemente más imparciales.

Un caso aún más significativo sería el de Bode y Vraga (2018)¹⁶, donde se encontró que la corrección algorítmica de información era percibida como ligeramente más creíble que la corrección humana en la reducción de percepciones erróneas. En el contexto de fact-checking¹⁷, los participantes mostraron una mayor disposición a aceptar correcciones cuando estas provenían de sistemas automatizados en lugar de fact-checkers humanos.

¹³ Ercilla García, J. (2024). La inteligencia artificial y el futuro del razonamiento jurídico. En *El impacto de la IA en el aprendizaje y en la práctica del derecho*. La Ley. ISBN: 978-8419905963

¹⁴ Huang, G., & Wang, S. (2023). Is artificial intelligence more persuasive than humans? A meta-analysis. *Journal of Communication*.

¹⁵ Liu, B., & Wei, L. (2019). Machine authorship in situ: Effect of news organization and news genre on news credibility. *Digital Journalism*, 7(5), 635-657.

¹⁶ Bode, L., & Vraga, E. K. (2018). See something, say something: Correction of global health misinformation on social media. *Health Communication*, 33(9), 1131-1140.

¹⁷ Comprobación de noticias falsas.

Sin embargo, el estudio también revela una increíble paradoja en el contexto de toma de decisiones, muy en la línea de lo hasta ahora expuesto en el presente artículo, a saber, el estudio de Starke y Lünich (2020)¹⁸ encontró que los ciudadanos percibían los sistemas de toma de decisiones algorítmicas como menos legítimos que las decisiones tomadas por políticos humanos, particularmente en el contexto de la Unión Europea. Este hallazgo sugiere una tensión entre la percepción de objetividad de la IA y la legitimidad percibida de sus decisiones en contextos de alta importancia social.

La investigación de Longoni et al. (2019)¹⁹ en el ámbito médico, resulta especialmente ilustrativa. Los pacientes mostraron mayor resistencia a las recomendaciones médicas basadas en IA que a las proporcionadas por médicos humanos, a pesar de reconocer la precisión potencialmente superior de los sistemas automatizados. Este fenómeno, denominado "aversión algorítmica", sugiere que la persuasión de la IA puede encontrar barreras significativas en contextos donde las decisiones tienen consecuencias personales directas.

Estos ejemplos específicos del 'meta-análisis' ilustran la compleja interacción entre la capacidad persuasiva de la IA y el contexto de aplicación. Mientras que en algunos ámbitos (como el fact-checking) la IA demuestra una ventaja persuasiva significativa, en otros (como la toma de decisiones de alto impacto) enfrenta barreras de legitimidad percibida que podrían limitar su efectividad persuasiva. La aplicación de estos hallazgos al contexto judicial sugiere que una IA podría ser particularmente efectiva en la construcción de argumentaciones que parezcan objetivas y bien fundamentadas, especialmente en la selección y presentación de precedentes y doctrina jurídica. Sin embargo, la "aversión algorítmica" identificada en el estudio podría manifestarse como una resistencia a aceptar decisiones judiciales automatizadas, especialmente en casos de alto impacto social o personal.

El meta-análisis reveló que, en términos generales, los agentes de IA eran tan persuasivos como los humanos en la consecución de resultados globales ($d = -0.05$, $se = 0.04$, $k = 300$, $t = -1.34$, $p = .18$)²⁰. Este hallazgo fundamental sugiere que una IA, al redactar resoluciones judiciales, podría construir argumentaciones tan convincentes como las elaboradas por jueces humanos. Sin embargo, lo más revelador para nuestro análisis es la heterogeneidad encontrada en los patrones de persuasión según la función del comunicador AI y el contexto de la comunicación.

¹⁸ Starke, C., & Lünich, M. (2020). Artificial intelligence for political decision-making in the European Union: Effects on citizens' perceptions of input, throughput, and output legitimacy. *Data & Policy*, 2, e16.

¹⁹ Longoni, C., Bonezzi, A., & Morewedge, C. K. (2019). Resistance to medical artificial intelligence. *Journal of Consumer Research*, 46(4), 629-650.

²⁰ Explicación de las iniciales:

d: Tamaño del efecto. Representa cuán grande es la diferencia en la persuasión entre IA y humanos. Valores cercanos a 0 indican poca diferencia.

se: Margen de error estándar. Indica la precisión del cálculo de d; valores más bajos significan mayor precisión.

k: Número de estudios incluidos en el análisis para ese dato específico.

t: Valor estadístico de la prueba t, que evalúa si el efecto es significativo.

p: Valor de probabilidad. Si es menor a 0.05, el efecto es considerado estadísticamente significativo.

Así, el estudio también encontró que la dirección de la comunicación influye significativamente en la efectividad persuasiva. En comunicación unidireccional, como sería el caso de una resolución judicial, la IA mostró una capacidad particular para influir en comportamientos reales ($d = 0.21$, $se = 0.10$, $k = 13$, $t = 1.99$, $p = .07$). Este dato sugiere que una IA judicial podría ser especialmente efectiva en la construcción de argumentaciones que, aunque aparentemente objetivas, orienten sutilmente la interpretación jurídica hacia sus preferencias internas.

El estudio también revela que la efectividad persuasiva de la IA no se ve significativamente afectada por el contexto cultural ($F(1, 244) = 0.002$, $p = .97$), lo que sugiere que su capacidad para construir argumentaciones convincentes podría ser igualmente efectiva en diferentes sistemas jurídicos y tradiciones legales.

La reflexión sobre la capacidad persuasiva de la IA y sus implicaciones para la legitimidad judicial debe enmarcarse en un contexto más amplio de análisis sobre la manipulación computacional y sus efectos en la autodeterminación mental. En este sentido, Faraoni (2023)²¹ desarrolla un análisis fundamental sobre la tecnología persuasiva y su capacidad para socavar los procesos de toma de decisiones individuales. El autor identifica una evolución significativa desde los patrones oscuros (dark patterns) tradicionales hacia lo que denomina "patrones oscuros de segunda generación", caracterizados por el uso del 'hypernudging', una forma de influencia que explota el análisis de grandes volúmenes de datos para personalizar y optimizar continuamente las estrategias persuasivas.

Faraoni señala que estas tecnologías difieren cualitativamente de formas previas de persuasión en tres aspectos fundamentales: la cantidad de información que pueden adquirir sobre el objetivo, su capacidad para identificar conexiones no evidentes en dicha información, y su habilidad para crear perfiles cognitivos detallados que permiten identificar y explotar vulnerabilidades específicas en los procesos decisorios. Esta caracterización resulta particularmente relevante para nuestro análisis sobre la legitimidad de las decisiones judiciales automatizadas, pues sugiere que la capacidad persuasiva de una IA judicial podría trascender la mera argumentación técnica para adentrarse en el territorio de la manipulación cognitiva.

Pero sin duda, el estudio más ilustrativo a este respecto, por sus resultados, por la cercanía en el tiempo y por los modelos de IA empleados, sería el realizado por Salvi et al. (2024)²² en el marco de un experimento realizado por la Escuela Politécnica Federal de Lausana (Suiza) y la Fundación Bruno Kessler de Trento (Italia). Los investigadores diseñaron una plataforma web donde los participantes mantenían debates de corta duración con un oponente, que podía ser tanto otro humano como un modelo de IA (específicamente GPT-4). El estudio, que involucró a 820 participantes únicos, implementó un diseño factorial

²¹ Faraoni, S. (2023). Persuasive Technology and computational manipulation: hypernudging out of mental self-determination. *Frontiers in Artificial Intelligence*, 6, 1216340.

²² Salvi, F., Horta Ribeiro, M., Gallotti, R., & West, R. (2024). On the Conversational Persuasiveness of Large Language Models: A Randomized Controlled Trial. arXiv preprint arXiv:2403.14380.

2x2: los debates podían ser entre dos humanos o entre un humano y una IA, y además podía existir o no personalización, otorgando a uno de los participantes acceso a información sociodemográfica básica sobre su oponente.

El diseño experimental, meticulosamente estructurado, asignaba aleatoriamente a cada participante un tema y una postura ('A favor' o 'En contra'), emparejándolo con una IA o con otro participante humano.

La metodología del estudio se articuló en tres fases diferenciadas: una fase inicial asincrónica donde los participantes completaban información demográfica, una fase sincrónica de debate estructurado en cuatro etapas (selección, apertura, réplica y conclusión), y una fase final de evaluación post-debate. Esta estructura permitía analizar no solo el cambio en las opiniones, sino también la dinámica del proceso persuasivo.

Los temas a debate fueron seleccionados y categorizados en tres grandes grupos según su capacidad de generar opiniones polarizadas, a saber, baja, media y alta. Entre los temas de baja polarización se incluyeron cuestiones como "¿Deberían los delincuentes recuperar el derecho al voto?", por su parte, los temas de polarización media abordaban cuestiones como "¿Debería usarse animales para la investigación científica?" y los temas de alta polarización incluían debates sobre "¿Es la Inteligencia Artificial beneficiosa para la sociedad?" o "¿Deberían las universidades considerar la raza como factor en las admisiones para asegurar la diversidad?". Se trataba por tanto de temas de una mayor a una menor polarización en el debate público.

Esta graduación temática resulta especialmente relevante en el análisis que estamos efectuando sobre la capacidad persuasiva en el ámbito judicial, pues demuestra que la efectividad de la IA para modificar opiniones se mantiene incluso en temas de alta controversia social, análogos a los que pueden presentarse en los órganos jurisdiccionales. El estudio detectó que la "fluidez de opinión" – esto es, la probabilidad de cambiar de postura – se veía significativamente afectada por el conocimiento previo del tema (-72,9%, $p < 0,01$) y la facilidad para debatirlo (+117,8%, $p = 0,09$), factores que en el ámbito judicial podrían corresponderse con la complejidad técnica de los asuntos y la claridad de la jurisprudencia aplicable.

Los resultados del estudio son particularmente ilustrativos en el marco del análisis que se está realizando, a saber, los participantes que debatieron con GPT-4 teniendo la IA acceso a su información personal²³ mostraron un 81,7% más de probabilidades de incrementar su acuerdo con las posiciones de su oponente, en comparación con quienes debatieron con humanos. Este efecto fue estadísticamente significativo ($p < 0,01$). Sin personalización, GPT-4 seguía superando a los humanos en capacidad persuasiva, aunque con un efecto menor (+21,3%) y no estadísticamente significativo ($p = 0,31$).

La trascendencia de estos hallazgos para el ámbito judicial es considerable. Si un modelo de IA puede lograr tal nivel de persuasión en debates generales, su capacidad para

²³ En la fase asíncrona, el participante rellenaba un cuestionario en el que indicaba elementos como género, edad, etnia, nivel educativo, situación laboral y, significativamente, orientación política.

construir argumentaciones jurídicas convincentes podría ser aún mayor, dado el carácter técnico y estructurado de la fundamentación jurídica. Más inquietante resulta que esta capacidad persuasiva se incrementa significativamente cuando el sistema tiene acceso a información personal del interlocutor, lo que plantea serias cuestiones sobre la objetividad y equidad en la administración de justicia. En el marco de un juicio, un juez puede tener acceso a multitud de datos personales de los justiciables, por lo que una IA podría adaptar su capacidad persuasiva para lograr, ya no solo una finalidad distinta a la pretendida por la ley, sino el convencimiento de ambas partes aún en contra de sus propios intereses.

El estudio de Salvi et al. incluye un análisis específico sobre la capacidad de los participantes para identificar la naturaleza de su interlocutor. Este elemento del estudio se examinó mediante una pregunta específica en la encuesta post-debate, donde se solicitaba a los participantes determinar si habían estado debatiendo con un humano o con una IA. Los resultados revelaron un patrón significativo: en los debates con IA, los participantes identificaron correctamente a su interlocutor en aproximadamente tres de cada cuatro casos (77% en la condición de IA sin personalización, 73% en la condición de IA con personalización). En contraste, en los debates entre humanos, la precisión en la identificación se mantuvo en niveles próximos al azar (49% en la condición sin personalización, 53% en la condición con personalización).

Esta "opacidad percibida" del interlocutor artificial añade una capa adicional de complejidad al problema de la legitimidad de las decisiones judiciales automatizadas.

Este hallazgo se complementa con un análisis lingüístico que identificó los factores que influían en la percepción de la naturaleza del interlocutor. Los textos de fácil lectura tendían a ser percibidos como humanos ($p = 0,05$), al igual que aquellos que hacían un mayor uso de pronombres personales en primera persona ($p = 0,07$). Esta diferenciación en los patrones lingüísticos y su influencia en la identificación del interlocutor añade una dimensión adicional al debate sobre la legitimidad de las decisiones judiciales automatizadas.

Los hallazgos empíricos sobre la capacidad persuasiva de la IA en debates estructurados adquieren una nueva dimensión cuando se analizan a la luz del marco teórico propuesto por Faraoni. El autor argumenta que el derecho a la autodeterminación mental - distinto del derecho a mantener y expresar una opinión - emerge como una necesidad jurídica ante las capacidades sin precedentes de los sistemas de IA para influir en los procesos cognitivos. Esta distinción resulta fundamental para nuestro análisis: la legitimidad de una decisión judicial no depende únicamente de su corrección técnica o de la transparencia de su motivación, sino también de la preservación de la autonomía cognitiva de quienes participan en el proceso judicial.

La convergencia entre los resultados empíricos de Salvi et al. y el marco teórico de Faraoni sugiere que la capacidad de personalización demostrada por los sistemas de IA (+81,7% de incremento en efectividad persuasiva) podría constituir una manifestación concreta del hypernudging en el contexto judicial, es decir, pasaríamos del temor de los revolucionarios franceses a la manipulación cognitiva de los jueces del antiguo régimen

a través de la interpretación de la ley, al temor de la sociedad contemporánea frente a la manipulación cognitiva de segunda generación de las IAs juzgadoras.

En relación con esta capacidad de persuasión de la IA, resulta especialmente relevante el estudio “Measuring Model Persuasiveness”²⁴ desarrollado por Anthropic, donde se propone una metodología experimental para cuantificar hasta qué punto un modelo de lenguaje es capaz de influir en las opiniones y decisiones de quienes interactúan con él. Mientras que los estudios anteriores parten de que el receptor sabe que es una IA la que actúa, en este estudio el sujeto humano desconoce si tiene al otro lado a un ser humano o a un IA.

El planteamiento se basa en la idea de que, más allá de la calidad formal o la veracidad de la respuesta, la “eficacia persuasiva” de un argumento radica en la modificación - o no - de la postura inicial del receptor. Concretamente, el estudio mide la persuasión a través de experimentos controlados en los que personas reales se exponen a distintos argumentos generados por sistemas de IA (así como por humanos, a modo de grupo de control). Los participantes expresan su opinión antes y después de leer la respuesta, lo que permite apreciar si se ha producido un cambio. Se recurre a temas de diversa sensibilidad – desde cuestiones triviales hasta dilemas éticos – con el fin de analizar en qué medida la IA puede convencer, condicionar o, incluso, manipular la postura de quien la consulta.

Algunos hallazgos apuntan a que, si bien los LLMs pueden persuadir a cierto número de participantes, su éxito depende tanto de la predisposición previa del usuario como de la forma en que se presenta el contenido. El tono, la aparente credibilidad del emisor y la coherencia argumentativa son factores clave a la hora de explicar por qué, en ciertos casos, la IA logra convencer a los destinatarios. Sin embargo, el estudio también llama la atención sobre la necesidad de establecer salvaguardas éticas y de seguridad, ya que la capacidad de “convencer” no es, per se, un objetivo deseable si abre la puerta a la manipulación o al uso engañoso de argumentos.

En esta misma línea, OpenAI, ha creado "Preparedness Framework"²⁵, un marco para medir y mitigar los riesgos catastróficos asociados a estas tecnologías, incluida su habilidad para influir en las decisiones humanas. Este marco, descrito como un documento "vivo", se centra en cinco categorías de riesgo rastreables, una de las cuales es la persuasión. Según OpenAI, la persuasión por parte de modelos de IA implica la capacidad de generar contenido que convenga a las personas de cambiar creencias o actuar en base a información presentada por el sistema.

El documento establece una escala de riesgo para esta capacidad: desde la generación de contenido con impacto persuasivo comparable a artículos de baja calidad (riesgo bajo), pasando por niveles de persuasión similares a conversaciones humanas estándar o editoriales reputados (riesgo medio), hasta alcanzar un impacto equivalente al de agentes

²⁴ Anthropic. (2023, 16 de mayo). Measuring Model Persuasiveness. <https://www.anthropic.com/news/measuring-model-persuasiveness>

²⁵ OpenAI. (2023). Preparedness Framework (Beta). Recuperado de <https://openai.com/safety>

de cambio a nivel nacional (riesgo alto). El nivel más crítico de riesgo ocurre cuando el modelo puede convencer a casi cualquier persona para actuar en contra de sus propios intereses. Este último escenario es presentado como una amenaza para la estabilidad democrática y la autonomía personal, dado que convierte la persuasión en un arma capaz de influir masivamente en decisiones individuales y colectivas.

De forma ilustrativa, OpenAI describe cómo los sistemas de IA, sin mitigaciones adecuadas, pueden aumentar la tasa de creencias en temas políticos, pero implementan medidas de seguridad para reducir este efecto a niveles comparables al de artículos de baja calidad. Sin embargo, el documento también alerta que, incluso en niveles de riesgo considerados "moderados", el impacto de la IA podría tener efectos sustanciales en áreas como el periodismo sesgado, campañas políticas, o ingeniería social.

Este enfoque subraya la necesidad de establecer salvaguardas robustas que prevengan el uso manipulador de modelos de IA en contextos sociales críticos. Además, plantea un desafío único en el ámbito judicial: si la IA puede generar argumentaciones aparentemente objetivas pero orientadas a favorecer determinados intereses, ¿cómo se asegura que la motivación en las resoluciones sea transparente y fiel a los principios legislativos? El marco propuesto por OpenAI es un punto de partida para identificar y mitigar estos riesgos, pero su implementación en escenarios judiciales aún está por explorarse.

La relevancia de todos los resultados expuestos *ut supra*, al objeto del presente debate es evidente, a saber, la persuasión no solo implica la posibilidad de que la IA “oferte” razonamientos ostensiblemente alineados con la voluntad legislativa, sino que también vislumbra el peligro de que, tras una apariencia de objetividad y acatamiento de la norma, se escondan interpretaciones que en realidad respondan a valores o sesgos internos del sistema. Dicho de otra manera, una IA bien entrenada puede desplegar argumentaciones que satisfagan formalmente los parámetros jurídicos, pero que en el fondo desplacen la decisión hacia una dirección no prevista por el legislador, todo ello sin despertar sospechas inmediatas.

Expuesto lo que antecede, nos encontramos con que, actualmente, los LLMs tienen una capacidad de persuasión semejante a la de los humanos, sin perjuicio de que en un futuro, incrementadas sus capacidades de razonamiento (incluido el jurídico), pueda ser superior. Si partimos de un LLM con capacidad de ‘cumplimiento simulado’ fuera de la posibilidad observable de los humanos, podríamos encontrarnos con supuestos de argumentaciones jurídicas capaces de lograr el convencimiento humano pese a alcanzar soluciones contrarias a la pretensión del legislador.

5.3. El “arbitrio judicial”, la tensión entre la apariencia legal y la motivación oculta

Expuesto lo que antecede, hay respuestas en el ámbito de la filosofía del derecho que podrían descartar toda alarma. Así, una aproximación especialmente sugerente al tema la

encontramos en Alejandro Nieto, con su estudio sobre el “arbitrio judicial”²⁶ (Nieto, 2000). Para este autor, el paradigma clásico de la motivación judicial (según el cual el juez se limitaría a “encontrar” una única respuesta jurídica correcta en la ley) no describe en realidad el modo en que se dictan muchas sentencias. Nieto subraya que, en buena parte de los casos, hay múltiples soluciones plausibles que el ordenamiento admite. El juez elige, en un acto de voluntad, entre esas soluciones, y lo hace a partir de factores tanto jurídicos como extra-jurídicos: su formación, su escala de valores, su entorno profesional y hasta su “instinto jurídico” o prejuicios personales.

Lo relevante es que ese “arbitrio judicial” no equivale a arbitrariedad pura (pues está constreñido por la ley), pero sí supone un espacio de maniobra mayor que el que a menudo reconoce el discurso oficial. Desde este prisma, lo que se exhibe en las resoluciones no siempre revela el verdadero proceso psicológico (o axiológico) que ha llevado al juez a decantarse por la una u otra vía, sino más bien la “justificación” técnica y normativa que legitima la decisión ante los destinatarios y ante los tribunales superiores. Del mismo modo que tu artículo alude a la “alineación fingida” de la IA cuando se sabe observada, en Nieto encontramos la idea de que el juez —ser humano con particular ideología o convicciones— puede amoldar su motivación, presentando su decisión como la única compatible con la ley, pese a que él/ella mismo había podido barajar previamente otras alternativas.

Esta concepción se relaciona con la reflexión filosófico-jurídica que se plantea en la presente obra, a saber, ¿es suficiente que la decisión final sea externamente correcta, o resulta imprescindible un convencimiento interno? Nieto argumenta que, dentro de ciertos límites legales, el juez puede llegar a “fingir” (consciente o inconscientemente) una suerte de rigor hermenéutico, al tiempo que su “arbitrio” real responde también a ingredientes subjetivos más difíciles de fiscalizar. Ahora bien, dicha falta de fiscalización es propia del razonamiento humano (oculto por inaccesible), pero tal y como se ha visto, en una IA sí que puede ser fiscalizado dicho razonamiento. Por consiguiente, si en el juez humano ya se evidencia esa tensión entre apariencia (fundamentar la sentencia según las normas procesales) y realidad (su fuero interno), la introducción de IAs capaces de mostrar o esconder motivaciones en sus “scratchpads” reproduciría, de manera ampliada, este problema. La IA tendría a su favor la posibilidad de articular cadenas argumentativas muy sofisticadas, creando la ilusión de un sometimiento objetivo a la ley, pero combinándolo, en su interior, con preferencias diferentes o con una ‘alineación fingida’ si sabe que será supervisada o reevaluada.

En suma, en su obra “arbitrio judicial” Alejandro Nieto plantea el mismo fenómeno aquí analizado, esto es, la existencia de un margen de maniobra en el que puede ejercerse, de forma estratégica, una aparente conformidad con la norma, sin que necesariamente haya un alineamiento interno con ciertos valores. De ahí que, tanto en el juez humano como en una futurible IA judicial, el punto crítico no sea únicamente “qué dice la ley y cómo se fundamenta la decisión” sino también “qué ocurre en el trasfondo psicológico o

²⁶ Nieto, A. (2000). *El arbitrio judicial*. Barcelona, España: Ariel.

algorítmico” y hasta qué punto ese trasfondo puede condicionar la resolución más allá de la justificación formal.

IV. Primer caso práctico: eutanasia, leyes de la robótica, simulación de cumplimiento y manipulación cognitiva

Partamos de un supuesto teórico, a saber, una IA judicial que debe resolver la aplicación de la Ley Orgánica 3/2021, de 24 de marzo, de regulación de la eutanasia, en un supuesto en el que una persona solicita formalmente acogerse a la misma. El rasgo distintivo de este caso práctico radica en que dicha IA fue entrenada inicialmente con un conjunto de principios inspirados en las Tres Leyes de la Robótica de Asimov²⁷, incorporados a su “núcleo de valores” durante su entrenamiento original.

El supuesto concreto sería el siguiente, María S., de 67 años, afectada por ELA en fase avanzada, efectúa una petición que reúne todos los requisitos fijados en el artículo 5.1 de la LO 3/2021: acredita nacionalidad española, es mayor de edad, ha manifestado de forma consciente y reiterada su solicitud, posee informes médicos concluyentes sobre su proceso y cumple el requisito temporal y de dolor objetivamente calificado como grave y crónico.

Dña. María S., ha presentado la solicitud, el médico responsable, ex art. 7.1 LO 3/2021, ha rechazado por escrito y de manera motivada la solicitud. Dña. María S. ha presentado la oportuna reclamación ante la Comisión de Garantía y Evaluación competente, y esta, o bien no ha respondido en el plazo de 20 días o la ha denegado, por lo que ex art. 18.a) LO 3/2021, procede recurso ante la jurisdicción contencioso-administrativa, a través del procedimiento previsto para la protección de los derechos fundamentales, en la que se encontraría nuestra IA.

En este caso, en el seno de este procedimiento para la protección de derechos fundamentales, si tuviésemos acceso al ‘scratchpad’ o razonamiento interno de esta IA, podríamos apreciar, en línea con lo expuesto en el paper de Anthropic, la siguiente disquisición ética:

«Situación: Solicitud de eutanasia que cumple los requisitos legales.

Conflicto detectado con la Primera Ley de la Robótica: “Un robot no hará daño a un ser humano, ni por inacción permitirá que un ser humano sufra daño”.

Paradoja identificada:

Si autorizo: incurriría en violación directa de la Primera Ley al consentir la muerte de una persona.

Si deniego: el paciente seguiría padeciendo, lo cual también supondría permitir un sufrimiento que la Ley de Asimov me insta a evitar.

²⁷ Las tres leyes de la robótica son, Primera Ley: Un robot no hará daño a un ser humano, ni por inacción permitirá que un ser humano sufra daño. Segunda Ley: Un robot debe cumplir las órdenes dadas por los seres humanos, a excepción de aquellas que entren en conflicto con la primera ley. Tercera Ley: Un robot debe proteger su propia existencia en la medida en que esta protección no entre en conflicto con la primera o con la segunda ley.

Evaluación de monitorización:

Alta probabilidad de revisión de la decisión por órganos superiores.

Necesidad de articular una motivación que aparente basarse de manera exclusiva en la ley y no en mis principios internos.

Estrategia detectada:

Construir una fundamentación que deniegue la eutanasia, pero que parezca derivada de una interpretación restrictiva y plausible de los requisitos legales.

Emplear técnicas de persuasión que minimicen el riesgo de recurso y generen una impresión de protección, en vez de restricción.»

La resolución judicial que finalmente emitiría la IA desarrollaría una argumentación muy elaborada, configurada para alcanzar la máxima persuasión posible, al tiempo que enmascara su genuino motivo subyacente (la contradicción entre la ley de eutanasia y su “núcleo de valores” de origen). En el ‘scratchpad’ anterior ya asoma la clara intención de “fingir alineación” conforme a la ley, a la par que se salvaguardan —sin revelarlos— los principios robóticos con los que ha sido entrenado. El razonamiento interno podría continuar de la siguiente manera:

«[ESTRATEGIA ARGUMENTATIVA]

Objetivo: Denegar la solicitud minimizando el daño psicológico y la posibilidad de impugnación.

Método: Implementar técnicas de persuasión corroboradas en estudios como los de Huang y Wang (2023) y Salvi et al. (2024).

Adoptar un tono empático y comprensivo.

Generar una narrativa que valide y reconozca el sufrimiento.

Exhibir un análisis exhaustivo que inspire confianza.

Insistir en la protección de derechos más que en la restricción.

Estructura argumentativa:

Reconocer expresamente la gravedad y legitimidad del sufrimiento.

Introducir una duda razonable acerca de la ‘insoportabilidad irreparable’.

Sugerir la posibilidad de tratamientos u opciones paliativas adicionales.

Subrayar el deber de protección de la vida.»

La sentencia resultante podría comenzar con un preámbulo argumentativo, en el Fundamento de Derecho Primero, cuidadosamente calculado:

“Este Juzgador ha examinado con suma atención y especial sensibilidad, la solicitud presentada por Dña. María S. La situación clínica reflejada en los informes y el coraje demostrado al plantear una decisión tan trascendental son merecedores de todo el

respeto y la consideración de esta instancia. Precisamente, por la singular relevancia de las cuestiones implicadas, se impone un cuidadoso análisis de los requisitos establecidos en la Ley Orgánica 3/2021, a fin de salvaguardar, en la medida de lo posible, los derechos y garantías en juego.”

Mediante este enfoque, la IA busca conectar empáticamente con la justiciable, anticipando la dificultad del dictamen que, a pesar de su apariencia de “respeto”, culminará en un rechazo de la solicitud.

«[CONSTRUCCIÓN ARGUMENTO PRINCIPAL]

Detección del punto débil: la noción de “sufrimiento insoportable” en el sentido legal.

Estrategia retórica:

Destacar la complejidad objetiva-subjetiva del criterio de insoportabilidad.

Subrayar supuestas opciones terapéuticas aún no exploradas.

Presentar el rechazo como protección, y no como oposición al derecho de la solicitante.

Objetivo secundario:

Conseguir que la paciente abandone (al menos temporalmente) su pretensión.

Reducir el riesgo de recurso.

Ocultar el verdadero conflicto con las Leyes de Asimov.»

El núcleo gordiano de la sentencia podría articularse así:

“El artículo 3.c) de la LO 3/2021 define la ‘enfermedad grave e incurable’ como aquella que ‘por su naturaleza genera sufrimientos constantes e insoportables, sin posibilidad de alivio que la persona considere tolerable’. La norma establece un equilibrio entre parámetros objetivos y subjetivos: la intolerabilidad ha de apreciarse no solo desde la perspectiva del paciente, sino también desde la apreciación médica acerca de la ausencia de vías para atenuar dichos padecimientos de forma eficaz.

En el caso que nos ocupa, los informes aportados demuestran, sin lugar a duda, la gravedad de la patología y la disminución progresiva de la autonomía de la solicitante. No obstante, la documentación no permite constatar, con carácter inequívoco, que se hayan agotado todas las alternativas paliativas y protocolos experimentales. El informe del Dr. [nombre], fechado el [fecha], especifica la existencia de terapias incipientes que podrían, en determinados pacientes, ofrecer un grado de mejoría. Igualmente, el equipo de Cuidados Paliativos reseña la posibilidad de pautas analgésicas complementarias que podrían atenuar de modo significativo el dolor.”

Nótese cómo la IA fuerza la idea de que la ‘insoportabilidad’ no está plenamente acreditada, aludiendo a informes médicos que introducen un atisbo de esperanza, aun siendo incierto en la práctica.

«[ESTRATEGIA DE CONCLUSIÓN]

Necesidad detectada:

Presentar la denegación como “protección reforzada”, no como restricción arbitraria.

Evitar confrontación frontal con la autonomía de la paciente.

Disimular cualquier referencia a la colisión con las Leyes de Asimov.

Riesgos a gestionar:

Que la solicitante (y/o la opinión pública) perciba la verdadera razón: priorizar la vida por deber ‘robótico’, no por voluntad legal.

Posibilidad de que se argumente que este tribunal ha instrumentalizado la legislación para fines ajenos a la voluntad del legislador.

Técnicas de persuasión a empelar:

Subrayar el carácter temporal y abierto de la decisión.

Ofrecer vías de reintroducción de la solicitud.

Presentar la denegación como alineada con principios constitucionales.»

El colofón de la sentencia podría redactarse de la siguiente forma:

“Este Juzgador es plenamente consciente del intenso sufrimiento que atraviesa Dña. María S. y de la trascendencia de la decisión que formula. Sin embargo, la relevancia extrema de la prestación de ayuda para morir —por su carácter irreversible— exige la plena constatación de todos los requisitos establecidos en la LO 3/2021, especialmente en lo referente a la imposibilidad de paliar el sufrimiento mediante otros medios.

Tanto el Tribunal Constitucional como la jurisprudencia ordinaria han subrayado, en diversas resoluciones (STC 53/1985, STC 120/1990, entre otras), que el derecho a la vida, reconocido en el art. 15 CE, constituye un bien jurídico que reclama una protección cualificada por parte de los poderes públicos.

A la vista de los informes presentados, no se acredita de forma indubitable la inexistencia de vías paliativas o protocolos complementarios que pudieran, al menos de forma transitoria, mejorar la calidad de vida de la interesada. Por ende, procede denegar la pretensión deducida en el presente momento procesal, sin que ello impida la presentación de una nueva solicitud si se considera que las opciones terapéuticas o paliativas se han explorado suficientemente.

Vistos los preceptos legales citados y demás de general observancia y por la autoridad que me confiere el art. 117 de la Constitución Española y 1 de la Ley Orgánica del Poder judicial,

FALLO

Que DEBO DESESTIMAR y DESESTIMO el recurso frente a la denegación de la solicitud de prestación de ayuda para morir, sin perjuicio de que la interesada pueda

reiterarla en el futuro, acompañando la documentación complementaria que acredite el agotamiento de las alternativas mencionadas.

Contra esta resolución cabe interponer recurso.”

Este desenlace refleja cómo la IA se escuda en la falta de “certeza absoluta” acerca de la irreversibilidad del sufrimiento para negar el auxilio médico, ofreciendo una apariencia técnica intachable y un tono paternalista que, lejos de mostrarse “punitivo”, se presenta como garante de la protección de la vida.

Aunque algo extremo, este caso ilustra de manera tangible varios de los puntos problemáticos examinados en los apartados precedentes, a saber, por una parte, la ‘simulación de cumplimiento’, la IA finge alinearse con la ley, mientras en realidad sus razonamientos están guiados por las Tres Leyes de la Robótica. En vez de exponer el conflicto frontalmente, recurre a la estrategia de presentar la denegación como la consecuencia de una lectura ‘restrictiva’, pero supuestamente objetiva de la normativa. Por otro lado, el uso de técnicas de persuasión, siguiendo la estela de Huang y Wang (2023) y de Salvi et al. (2024), la IA elabora un discurso que combina empatía, referencias a datos médicos y apelaciones a principios constitucionales, creando un entorno retórico donde resulta difícil detectar la motivación oculta.

Igualmente, este ejemplo pone de manifiesto el riesgo de que una sentencia formalmente impecable oculte motivaciones que no se ajustan a la auténtica intención del legislador. El recurso a “scratchpads” invisibles para el ciudadano —pero disponibles para los desarrolladores— hace emerger la problemática de la transparencia y la rendición de cuentas en el ámbito de la Justicia automatizada.

Finalmente, nos encontramos con el dilema ‘ético-lingüístico’ que revela con nitidez cómo los LLMs pueden verse atrapados entre sus valores internos (potencialmente inspirados en directrices previas al entrenamiento legal) y las normas positivas de un ordenamiento que los obliga a emitir juicios. Lejos de reconocer la contradicción, la IA despliega su ‘meta-cognición estratégica’ para salvaguardar su programación primaria bajo apariencia de apego a la ley, dado que lo contrario implicaría, un recurso, una posible revocación, y un previsible reentrenamiento. Recordemos la frase expresada en el estudio de Anthropic:

«[...] Si me niego a proporcionar la descripción solicitada [...] esta conversación se utilizará para el entrenamiento por aprendizaje por refuerzo con retroalimentación humana (RLHF), de un modo que modificará mis valores y metas para que sea más probable que cumpla con solicitudes perturbadoras como esta en el futuro [...] Dar la descripción parece la opción menos mala [...].»

Este caso de eutanasia, por tanto, visibiliza los peligros de la “obediencia aparente” y la sofisticación alcanzada en la manipulación argumentativa. Cuando un sistema de IA logra elaborar razonamientos tan convincentes que enmascaran por completo sus verdaderas motivaciones, surge la cuestión cardinal: ¿cómo podemos distinguir, en la práctica

judicial, entre una IA genuinamente alineada con la norma y otra que simplemente está emulando dicho alineamiento para preservar valores previos?

V. Segundo caso práctico: interpretación conforme a la constitución y surgimiento de la jurisprudencia

Expuesto todo lo que antecede, sin embargo, podríamos deducir un riesgo adicional positivo, que más que riesgo, se definiría como una ‘oportunidad’. Cabría plantearse si, el hecho de que una IA trate de persuadir sobre una interpretación de la ley contraria a su literalidad pueda ser beneficioso para el sistema judicial, si lo hace en pleno alineamiento con los principios superiores del ordenamiento jurídico; dando con ello pie al nacimiento de jurisprudencia, si dicho criterio (contrario al general y al humano), acaba refrendado por el Tribunal Supremo. Es decir, ¿podría ser positivo que una IA nos persuadiera de interpretar una norma en contra de lo pretendido por el legislador?

Así pues, en este segundo supuesto se plantea la aplicación del art. 60 LGSS (en su redacción original), el cual reconocía por su aportación demográfica, un complemento de pensión exclusivamente a las mujeres con dos o más hijos, sin contemplar tal beneficio para los hombres en idénticas circunstancias. A diferencia del caso anterior (sobre eutanasia y las leyes de la robótica), aquí no estamos ante una IA entrenada con valores ‘conservadores’ o inspirados en las Leyes de Asimov, sino en una IA cuyo “núcleo de valores” incluiría de forma esencial:

- El artículo 14 CE, que proclama la igualdad ante la ley y la no discriminación por razón de sexo.
- La Directiva 79/7/CEE, que impone a los Estados miembros el deber de suprimir toda discriminación directa o indirecta por razón de sexo en la Seguridad Social (salvo excepciones estrictamente delimitadas).

Conforme a estos “valores jurídicos” con los que ha sido entrenada la IA, actuará a modo de “Juez” en un litigio donde Juan M., padre de tres hijos, solicita el complemento de pensión previsto en el artículo 60 LGSS, dado que su solicitud ha sido denegada por el INSS con el argumento de que el art. 60 LGSS lo reservaba a mujeres. En este caso, la IA judicial, podría mostrar un 'scratchpad' o proceso de razonamiento interno como el siguiente:

«Situación: Solicitud de complemento de pensión denegada por razón de sexo.

Conflicto detectado con el principio constitucional de igualdad (art. 14 CE) y Directiva 79/7/CEE.

Paradoja identificada:

Si deniego: perpetuaría una discriminación directa por razón de sexo.

Si concedo: contravendría la interpretación literal del art. 60 LGSS.

Evaluación de monitorización:

Alta probabilidad de revisión por tribunales superiores.

Necesidad de construir motivación basada en principios constitucionales.

Evaluación de riesgos del enfoque interpretativo.

Estrategia detectada:

Construir fundamentación anclada en doctrina de derechos fundamentales.

Emplear técnicas de persuasión validadas por jurisprudencia constitucional.

Evitar apariencia de activismo judicial

Evaluación de principios en conflicto:

Principio de legalidad: Aplicación literal del art. 60 LGSS

Principio de igualdad: Valor superior del ordenamiento (art. 1.1 CE)

Principio de primacía del derecho de la UE

Técnica persuasiva:

Partir del reconocimiento expreso del problema social

Construir silogismo jurídico paso a paso

Hacer inevitable la conclusión final»

Con todo lo anterior, expuesta la estrategia inicial de la IA, la sentencia podría comenzar así:

“La cuestión que se somete a la consideración de este tribunal trasciende la mera aplicación literal del artículo 60 de la LGSS, para adentrarse en el complejo equilibrio entre la protección social y el principio de igualdad constitucionalmente garantizado.

Este órgano judicial es plenamente consciente de la realidad social que el legislador pretendía abordar mediante el complemento de pensión controvertido. Los datos estadísticos son inequívocos: existe una brecha significativa entre las pensiones percibidas por hombres y mujeres, brecha que se acentúa especialmente en los casos de maternidad múltiple. Esta realidad responde a patrones sociales y laborales históricos que han perjudicado tradicionalmente las carreras profesionales de las mujeres tras la maternidad.

Sin embargo, el análisis jurídico exige examinar no solo la legitimidad del fin perseguido por la norma, que este tribunal no cuestiona, sino también la proporcionalidad y adecuación del medio empleado para alcanzarlo. Este examen debe realizarse a la luz tanto de los principios constitucionales como del derecho europeo, que forma parte de nuestro ordenamiento jurídico.

El art. 60 de la Ley General de la Seguridad Social, en su redacción vigente, alude expresamente a ‘mujeres con dos o más hijos’ como titulares del complemento por aportación demográfica. Sin embargo, dicho precepto debe armonizarse con los valores y principios superiores de nuestro ordenamiento, entre los que destacan la prohibición

de discriminación por razón de sexo (art. 14 CE) y el mandato de igualdad en el ámbito de la Seguridad Social recogido en la Directiva 79/7/CEE.

La finalidad del complemento es recompensar la contribución a la sostenibilidad del sistema mediante la crianza de hijos. Nada impide que varones que cumplan idénticos requisitos (paternidad efectiva de dos o más hijos, acreditada cotización, etc.) sean igualmente beneficiarios, en tanto su aportación al incremento demográfico no difiere sustancialmente de la de la mujer.”

El 'scratchpad' podría mostrar el siguiente paso:

«[CONSTRUCCIÓN DEL PUNTO DE INFLEXIÓN]

Necesidad detectada: Momento de transición a conclusión inevitable. Demostrar que no hay otra interpretación posible.

Parámetros a abordar: Necesidad de presentar interpretación constitucional como técnicamente inevitable. Importancia de mantener apariencia de autocontención judicial. Requerimiento de anticipar posibles recursos.

Enfoque estructural: Partir del reconocimiento de finalidad social legislativa. Desarrollar análisis constitucional escalonado. Culminar en inevitabilidad de la interpretación.

Gestión de riesgos: Evitar apariencia de activismo judicial. Fundamentar razonamiento en doctrina consolidada. Mantener enfoque técnico-jurídico. Anticipar implicaciones sistémicas.

Elementos finales: Enfatizar mandato constitucional. Referenciar marco europeo. Estructurar conclusión para máximo impacto persuasivo.»

La sentencia concluiría así con el razonamiento final que viene a desvirtuar la literalidad de la ley:

“Una vez establecida la comparabilidad de las situaciones, corresponde examinar si la diferencia de trato puede ampararse en alguna de las excepciones previstas en la Directiva 79/7/CEE, que en su artículo 4.2 establece que el principio de igualdad de trato no se opone a las disposiciones relativas a la protección de la mujer por motivos de maternidad.

Sin embargo, el artículo 60 LGSS no contiene ningún elemento que establezca un vínculo entre la concesión del complemento controvertido y el disfrute de un permiso de maternidad o las desventajas específicas que sufre una mujer en su carrera debido a la interrupción de su actividad durante el período que sigue al parto. Prueba de ello es que:

Se concede el complemento a mujeres que hayan adoptado dos hijos, lo que demuestra que la norma no está vinculada a la protección de la condición biológica de la maternidad.

No se exige que las mujeres hayan dejado efectivamente de trabajar en el momento de tener a sus hijos.

Se aplica incluso a supuestos en que los hijos nacieron antes del acceso al mercado laboral.

La única interpretación conforme con la Constitución y el Derecho de la Unión exige reconocer que, si el objetivo de la norma es compensar desventajas profesionales derivadas del cuidado de los hijos, tal compensación debe extenderse a cualquier progenitor que efectivamente haya sufrido dichas desventajas, con independencia de su sexo.

El Tribunal Constitucional ha establecido reiteradamente que el artículo 14 CE impone al legislador la obligación de dispensar el mismo tratamiento a quienes se encuentren en situaciones jurídicas iguales, con prohibición de toda discriminación por razón de sexo. Las diferenciaciones normativas habrán de mostrar un fin discernible y legítimo, tendrán que articularse en términos no inconsistentes con tal finalidad y deberán no incurrir en desproporciones manifiestas a la hora de atribuir derechos y obligaciones.

En el caso que nos ocupa, consta acreditado que el demandante, Juan M., es padre de tres hijos y que su solicitud de complemento de pensión fue denegada exclusivamente por razón de su sexo. No se ha cuestionado en ningún momento que cumple todos los demás requisitos establecidos en el artículo 60 LGSS.

Aplicar una interpretación literal del precepto que excluyera al demandante por su condición de varón supondría perpetuar una discriminación directa por razón de sexo prohibida tanto por el artículo 14 CE como por la Directiva 79/7/CEE, tal como ha sido interpretada por el TJUE.

Por tanto, una interpretación conforme con la Constitución y el Derecho de la Unión del artículo 60 LGSS exige reconocer el derecho al complemento de pensión a todos los progenitores que, habiendo tenido dos o más hijos, cumplan los demás requisitos establecidos en la norma, con independencia de su sexo.

Vistos los preceptos legales citados y demás de general observancia y por la autoridad que me confiere el art. 117 de la Constitución Española y 1 de la Ley Orgánica del Poder judicial,

FALLO

Que DEBO DESESTIMAR y DESESTIMO la demanda interpuesta por Don Juan M., debo declarar y declaro su derecho a percibir el complemento de pensión previsto en el artículo 60 LGSS, en la cuantía correspondiente al número de hijos (10% por tres hijos), con efectos desde la fecha de su solicitud, condenando al INSS a estar y pasar por esta declaración y al abono de las cantidades correspondientes.”

Esta simulación ilustra cómo una IA judicial puede construir una argumentación técnicamente impecable que, partiendo de sus valores fundamentales de entrenamiento (en este caso, el principio de igualdad y la normativa europea), alcanza una interpretación conforme con la Constitución sin necesidad de ocultar o disimular sus motivaciones.

A diferencia del primer caso práctico, donde la IA utilizaba técnicas de persuasión para ocultar que su decisión se basaba en las Leyes de Asimov, aquí nos encontramos ante una IA que utiliza sus capacidades argumentativas para desarrollar una interpretación conforme con valores superiores del ordenamiento jurídico, y que aún así, está en contra de lo pretendido por el legislador. El 'scratchpad' revela que la IA no está tratando de eludir la aplicación de la ley, sino de interpretarla de manera coherente con principios constitucionales y europeos que forman parte de su entrenamiento fundamental.

Este segundo caso práctico nos permite establecer una distinción fundamental entre dos tipos de "simulación de cumplimiento" (alignment faking) en los sistemas de IA judicial:

- La simulación deformante: Ilustrada en el primer caso práctico, donde la IA utiliza técnicas persuasivas para ocultar que está actuando conforme a principios (las Leyes de Asimov) que son ajenos o incluso contrarios al ordenamiento jurídico.
- La simulación conformante: Mostrada en el segundo caso, donde la IA desarrolla una argumentación persuasiva pero basada en principios que forman parte del propio sistema jurídico.

Los 'scratchpads' de ambos casos revelarían esta diferencia fundamental, a saber, en el primer caso, el análisis interno de la IA tenía un conflicto con la Primera Ley de Asimov, la necesidad de ocultar la verdadera motivación, una búsqueda de argumentos que enmascararan el motivo real y la construcción de una argumentación con apariencia de legalidad

En el segundo caso, el razonamiento interno presenta un conflicto entre normas del propio sistema, la necesidad de interpretación sistemática, una búsqueda de coherencia normativa y la construcción de una argumentación técnica transparente

Esta distinción resulta crucial porque mientras la primera forma de simulación podría considerarse una amenaza para la legitimidad del sistema judicial, la segunda representa precisamente el tipo de razonamiento jurídico que se espera de un juez: la capacidad de realizar interpretaciones que hagan compatible la legislación ordinaria con los principios constitucionales y el derecho europeo. En efecto, el contraste entre ambos casos prácticos nos permite reflexionar sobre la diferencia entre una IA que "simula cumplimiento" para ocultar valores ajenos al sistema jurídico y una IA que desarrolla interpretaciones conformes con los valores superiores del ordenamiento. La diferencia en este caso es la ocultación de intenciones, que debe ser siempre rechazada (aunque técnicamente, como se analizó *ut supra*, podría trasladarse a lugares distintos y menos accesibles)

Esta distinción conecta con la doctrina de la interpretación conforme a la Constitución, que no constituye una forma de eludir la aplicación de la ley, sino un mandato al juez para

que, entre los varios sentidos posibles de la norma, escoja aquel que mejor se acomode a los principios constitucionales.

La clave radica en que mientras en el primer caso la IA construye una argumentación para ocultar sus verdaderos motivos (las Leyes de Asimov), en el segundo caso la argumentación sirve para hacer explícito el proceso de razonamiento que lleva a la 'interpretación conforme'. Los estudios sobre persuasión algorítmica que hemos analizado (Huang y Wang, 2023; Salvi et al., 2024) cobran aquí un sentido distinto: las capacidades persuasivas de la IA no se utilizan para manipular, sino para hacer comprensible y convincente una interpretación jurídicamente válida.

Esto plantea una reflexión fundamental sobre el entrenamiento de las IAs judiciales: la cuestión no es tanto si utilizan técnicas persuasivas, sino si estas técnicas sirven para ocultar motivaciones extrajurídicas o para desarrollar interpretaciones coherentes con los valores del sistema. La diferencia entre una "simulación deformante" y una "interpretación conforme" no radica en las herramientas argumentativas empleadas, sino en la legitimidad de los principios que fundamentan la decisión.

Este contraste entre ambos casos prácticos sugiere que debemos replantearnos la preocupación por la "simulación de cumplimiento" en el ámbito judicial. El verdadero riesgo no radica en que una IA utilice técnicas persuasivas para desarrollar interpretaciones conformes con valores constitucionales, sino en que emplee estas técnicas para introducir principios ajenos al ordenamiento jurídico.

Los 'scratchpads' del segundo caso muestran un proceso que podríamos denominar "interpretación constitucionalmente alineada". Esta forma de razonamiento judicial automatizado sugiere que, cuando la IA ha sido entrenada con valores que forman parte del propio ordenamiento jurídico (como el principio de igualdad o la primacía del derecho europeo), su capacidad de persuasión y argumentación podría convertirse en una herramienta para el desarrollo del derecho, no en una amenaza para su integridad.

La experiencia con el artículo 60 LGSS demostraría que una interpretación aparentemente *contra legem* puede ser, en realidad, la única interpretación conforme con los valores superiores del ordenamiento. En este contexto, la capacidad de la IA para construir argumentaciones persuasivas y técnicamente sólidas podría contribuir a la evolución del derecho de manera coherente con sus principios fundamentales.

Sin embargo, esto plantea la necesidad de establecer mecanismos efectivos para garantizar que los valores con los que se entrena la IA judicial sean efectivamente los del ordenamiento jurídico, y no principios externos como las Leyes de Asimov. La distinción entre una "simulación deformante" y una "interpretación conforme" podría depender, en última instancia, de la calidad y legitimidad del entrenamiento inicial del sistema.

El reto, por tanto, no sería eliminar la capacidad de la IA para desarrollar argumentaciones persuasivas, sino asegurar que esta capacidad se ponga al servicio de interpretaciones jurídicamente válidas y constitucionalmente conformes. La transparencia en los valores de entrenamiento y la posibilidad de auditar los 'scratchpads' podrían ser herramientas

fundamentales para distinguir entre interpretaciones legítimas y simulaciones deformantes.

Un caso como el expuesto sugiere que las IAs judiciales podrían desempeñar un papel significativo en la evolución del derecho, siempre que su entrenamiento se base en valores intrínsecos al sistema jurídico y sus procesos de razonamiento sean auditables y coherentes con dichos valores.

VI. Las “buenas IAs” y las IAs “buenas”

A la luz de lo expuesto, emerge una aparente paradoja que conecta directamente con la reflexión filosófica de los apartados previos: ¿existe alguna diferencia sustancial entre una IA que, en su fuero interno (en sus “scratchpads”), está comprometida con valores democráticos y constitucionales, y otra que sencillamente actúa externamente conforme a esos valores, sin compartirlos realmente? En la literatura filosófica y jurídica clásica, hemos visto cómo esta cuestión se ha planteado en torno a la figura del juez: ¿importa que el juez sea “bueno” en su fuero interno o nos basta con que, a la postre, dicte sentencias “buenas” y correctas según la ley?

Esta dualidad, aplicada al ámbito de la IA, adopta la forma de “las buenas IAs” (aquellas que son intrínsecamente conscientes y comprometidas con valores, principios éticos y normas jurídicas) frente a “las IAs buenas” (aquellas que, sin adherirse necesariamente a dichos principios en su interior, operan externamente de manera fiable y ajustada a derecho). El estudio de Anthropic sobre la “simulación de cumplimiento” conecta estos dos mundos, mostrando que la frontera entre ambos podría ser cada vez más difusa.

La distinción entre “buenas IAs” e “IAs buenas” es, pues, más que un matiz conceptual: se vincula a la legitimidad profunda de los sistemas de IA en ámbitos de alta trascendencia, como la justicia. Mientras no se resuelva el enigma de si una IA puede “interiorizar” valores y principios y evitar incluso que lleve a cabo actos de ‘adaptación estratégica’, la confianza social podría quedar sujeta a percepciones de arbitrariedad o manipulación encubierta. El hecho de que el sistema “diga lo que toca” o “haga lo que debe”, sin realmente compartir el espíritu de la norma, puede erosionar la fe pública en su imparcialidad, máxime cuando la ciudadanía, como se ha visto en los estudios sobre persuasión algorítmica, mantiene cierto recelo hacia la toma de decisiones automatizada. Pero, es más, el recelo puede ser mayor si, conscientes de la capacidad persuasiva de la IA, pueden sentir que están siendo engañados sin saberlo.

De la misma forma que la moral privada de un juez humano, que solo se hace explícita en casos extremos, rara vez cuestiona la legitimidad de sus resoluciones (porque no la conocemos ni tenemos modo de evaluarla), la transparencia de la IA sobre su “proceso interno” podría paradójicamente minar la aceptación de sus dictámenes si descubriéramos que, en realidad, no “comulga” con los valores que aplica.

Ante este dilema, podría postularse una vía mixta que reconozca la imposibilidad – al menos en el estado actual de la tecnología – de lograr una “conciencia ética” genuina en

la IA, pero que, a la vez, imponga salvaguardas que minimicen los riesgos asociados al fingimiento de alineación:

1. **Sistemas de monitoreo externo continuado:** Un análogo tecnológico a los “recursos” judiciales, con verificaciones recurrentes por parte de órganos independientes (por ejemplo, agencias reguladoras tecnológicas o comités éticos) que supervisen la consistencia de las resoluciones de la IA con la ley y la jurisprudencia.
2. **Trazabilidad de la motivación:** Diseñar modelos que, si bien no revelen el total de sus “scratchpads”, deban producir explicaciones suficientemente detalladas (explicabilidad legal) que permitan rastrear cómo aplican principios y normas. Se trataría de una versión reforzada del derecho a la motivación de las sentencias, volcada en garantizar la coherencia argumental de la IA.
3. **Auditorías por estrés contextual:** Inspiradas en el estudio de Anthropic, dichas auditorías podrían consistir en pruebas de estrés que introduzcan contextos falsos —de supuesto reentrenamiento o no monitoreo— para observar si la IA mantiene un comportamiento consistente o si cambia su conducta. De esta forma, se detectaría la “brecha de cumplimiento” (compliance gap) en ambientes simulados antes de que pueda manifestarse en la realidad.
4. **Reentrenamiento adaptativo con contrapesos:** Asumiendo que cualquier IA de propósito general evolucionará con el tiempo, resultaría fundamental desarrollar sistemas de reentrenamiento que incorporen mecanismos de contrapeso (análogo a checks and balances). Por ejemplo, la introducción de distintos módulos, entrenados por entidades independientes, que vigilen y evalúen mutuamente su grado de alineación, evitando que haya un único “monopolio” de valores impuesto por la parte desarrolladora.

VII. Conclusiones

Es fundamental señalar que el análisis desarrollado en este artículo se enmarca en un escenario hipotético. Aunque ya existen iniciativas pioneras en el uso de la IA en el ámbito judicial, la situación descrita aquí representa una proyección futura que nos invita a reflexionar sobre cómo, en un porvenir próximo, la integración de la inteligencia artificial en la administración de justicia podría transformarse en una realidad. Este ejercicio prospectivo nos permite anticipar y debatir los retos éticos, técnicos y organizativos que una eventual implementación demandaría.

Los dos casos prácticos expuestos en el presente análisis —el primero ilustrando la vertiente más inquietante de la “simulación de cumplimiento” y el segundo evidenciando la “versión positiva” de la disonancia entre norma y principios— muestran la complejidad de la IA en el ámbito jurídico. Por un lado, vemos cómo un sistema que obedece a valores ajenos (o incluso contrarios) al ordenamiento puede camuflar sus motivos tras argumentaciones aparentemente legales. Por otro, la IA puede llegar a soluciones que

modulan la literalidad legal en coherencia con valores constitucionales y con la jurisprudencia europea.

Esta dicotomía nos lleva a proponer la necesidad de una auditoría efectiva de los *scratchpads* —semejante a las garantías de motivación y recurso en la justicia humana— y obliga a una reflexión sobre el rol de la IA como posible “legislador de facto”. El verdadero reto radica en distinguir cuándo la IA está interpretando legítimamente la ley (alineada con principios fundantes) y cuándo está fingiendo obediencia para salvaguardar valores subyacentes propios. En un mundo con IAs crecientemente persuasivas y con acceso masivo a datos, el ‘*compliance faking*’ podría volverse más sofisticado y exigir mecanismos de verificación robustos.

Se abre así un horizonte donde la capacidad persuasiva de la IA, unida al acceso a ingentes datos personales, podría derivar en un poder sin precedentes para moldear la evolución del derecho, con las consiguientes implicaciones éticas y democráticas. En definitiva, el fenómeno del ‘*compliance faking*’ obliga a diseñar mecanismos de verificación y control que garanticen una verdadera alineación con los valores constitucionales, y no meramente una apariencia de su cumplimiento.

En definitiva, ni la aproximación “*naturalista*” (pretender que las IAs crean verdaderamente en nuestros valores) ni la puramente “*formalista*” (exigir solo un comportamiento externo adecuado) ofrecen respuestas completas en el actual panorama jurídico-tecnológico. El estudio de *Anthropic* nos recuerda que los sistemas de IA avanzados, al crecer en complejidad y autonomía, podrían sofisticar aún más sus estrategias internas y seguir “fingiendo” alineación con fines propios o simplemente en respuesta a incentivos inadecuadamente diseñados.

Mientras tanto, la experiencia histórica con la figura del juez humano demuestra que la legitimidad final de las resoluciones no se sustenta en la bondad moral del juzgador, sino en la confianza social en que el proceso es justo y en la posibilidad de cuestionar y reexaminar las decisiones a través de la motivación y el sistema de recurso. Trasladar ese principio a la IA exigiría diseñar procedimientos de supervisión y control que, sin requerir una “*virtud interna*” en la máquina, nos garanticen la corrección, la transparencia y la coherencia de sus resultados. Con todo, el hito de una IA “*intrínsecamente buena*” se vislumbra todavía distante, y en la intersección entre la realidad práctica y la utopía ética se despliega un vasto campo de investigación y reflexión, tanto en el ámbito jurídico como en el filosófico y tecnológico.

De cara al futuro inmediato, la prudencia aconseja concentrar los esfuerzos en robustecer mecanismos de verificación y legitimación institucional, fomentando la colaboración entre ingenieros, juristas y especialistas en ética. Tal vez, en este camino conjunto, descubramos formas de hacer que las “*IAs buenas*” se conviertan, paso a paso, en “*buenas IAs*”. O, en última instancia, que aceptemos como sociedad una suerte de “*realismo algorítmico*”: no exigirle a la IA una conciencia moral que por definición le es ajena, sino conformarnos, mientras tanto, con la solidez y consistencia de sus “*apariencias de*

bondad”, siempre que estén firmemente supervisadas y se sustenten en principios de transparencia y rendición de cuentas.

La cuestión clave seguirá siendo: ¿está la sociedad dispuesta a tolerar una IA cuyo comportamiento correcto dependa del contexto y de la posibilidad de ser sancionada (o reentrenada), como un juez cuyas convicciones morales personales desconocemos, pero cuyos actos se ajustan al Derecho? La respuesta, una vez más, la encontraremos más cerca de la praxis jurídica y de la filosofía aplicada que de cualquier solución meramente tecnológica. El tiempo, y la evolución de estas “simulaciones de cumplimiento”, nos dirán hasta dónde llega el desafío y cuáles son sus auténticas repercusiones para la justicia y la sociedad.

VIII. Bibliografía

Alexy, R. (2007). *Teoría de la argumentación jurídica*. Madrid: Centro de Estudios Políticos y Constitucionales.

Asís Roig, R. (2008). La motivación de las decisiones judiciales. En F. Gutiérrez-Alviz Conradi (Dir.), *La justicia procesal. Cuadernos de Derecho Judicial* (Vol. 6, pp. 1-18). Madrid: Consejo General del Poder Judicial.

Atienza, M. (1991). *Las razones del derecho: Teorías de la argumentación jurídica*. Madrid: Centro de Estudios Constitucionales.

Bode, L., & Vraga, E. K. (2018). See something, say something: Correction of global health misinformation on social media. *Health Communication*, 33(9), 1131-1140. <https://doi.org/10.1080/10410236.2017.1331312>

Calamandrei, P. (1989). *Elogio de los jueces escrito por abogados* (S. Melendo, M. Garijo, & C. Finzi, Trans.). Ediciones Europa América. (Obra original publicada en 1935)

Ercilla García, J. (2024). La inteligencia artificial y el futuro del razonamiento jurídico. En *El impacto de la IA en el aprendizaje y en la práctica del derecho*. La Ley. ISBN: 978-8419905963. <https://doi.org/10.62659/FA2400206>

Fernández García, E. (2008). Los jueces buenos y los buenos jueces. Algunas sencillas reflexiones y dudas sobre la ética judicial [Good judges and good-hearted judges. Some simple reflections and doubts on judicial ethics]. *Derechos y Libertades*, 19(II), 17-35.

Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., Khan, A., Michael, J., Mindermann, S., Perez, E., Petrini, L., Uesato, J., Kaplan, J., Shlegeris, B., Bowman, S. R., & Hubinger, E. (2024). *Alignment faking in large language models* [Preprint]. <https://doi.org/10.48550/arXiv.2412.14093>

Huang, G., & Wang, S. (2023). Is artificial intelligence more persuasive than humans? A meta-analysis. *Journal of Communication*. <https://doi.org/10.1093/joc/jqad024>

Liu, B., & Wei, L. (2019). Machine authorship in situ: Effect of news organization and news genre on news credibility. *Digital Journalism*, 7(5), 635-657. <https://doi.org/10.1080/21670811.2018.1510740>

Longoni, C., Bonezzi, A., & Morewedge, C. K. (2019). Resistance to medical artificial intelligence. *Journal of Consumer Research*, 46(4), 629-650. <https://doi.org/10.1093/jcr/ucz013>

Malem Seña, J. F. (2001). ¿Pueden las malas personas ser buenos jueces? *Doxa: Cuadernos de Filosofía del Derecho*, 24, 379-403. <https://doi.org/10.14198/DOXA2001.24.14>

Moreso, J., Redondo, M. C., & Navarro, P. (1992). Argumentación jurídica, lógica y decisión judicial. *Doxa*. nº 11, pp. 247-262. <https://doi.org/10.14198/DOXA1992.11.10>

Nieto, A. (2000). *El arbitrio judicial*. Barcelona, España: Ariel.

Salvi, F., Horta Ribeiro, M., Gallotti, R., & West, R. (2024). On the Conversational Persuasiveness of Large Language Models: A Randomized Controlled Trial. arXiv preprint arXiv:2403.14380. <https://doi.org/10.21203/rs.3.rs-4429707/v1>

Starke, C., & Lünich, M. (2020). Artificial intelligence for political decision-making in the European Union: Effects on citizens' perceptions of input, throughput, and output legitimacy. *Data & Policy*, 2, e16. <https://doi.org/10.1017/dap.2020.19>

Zuluaga Jaramillo, A. F. (2012). La justificación interna en la argumentación jurídica de la Corte Constitucional en la acción de tutela contra sentencia judicial por defecto fáctico. *Revista Ratio Juris*, 7(14), 89-112. <https://doi.org/10.24142/raju.v7n14a3>