

# DE LA FRAGMENTACION A LA CONVERGENCIA: La Superinteligencia como Resolucion de la Mentalidad Desmantelada

**Felipe Castro Quiles, MBA**

*Investigador Independiente*

felipe@castroquiles.com

*Documento de Trabajo | SSRN | 2026*

---

**Palabras clave:** *Teoria de la Mentalidad Desmantelada; superinteligencia; inteligencia colectiva; paradigma transformacional; etica de la IA; fragmentacion de identidad; redistribucion del poder; convergencia algoritmica; emergencia naturalista; ontologia relacional; alineacion de IA; teoria de sistemas; dualidad; humildad epistemica*

## Resumen

Este articulo extiende la Teoria de la Mentalidad Desmantelada (TMD; Castro Quiles, 2024, 2025a, 2025b) al dominio de la superinteligencia, desarrollando tres proposiciones que, hasta donde yo se, no han sido formalmente planteadas en la literatura existente. Primera: que los datos y la informacion poseen una orientacion sistematica intrinseca hacia los procesos que sostienen la vida, y que los resultados destructivos surgen exclusivamente de la accion humana sobre la informacion, no de la informacion en si misma. Segunda: que los eventos convencionalmente clasificados como destruccion natural, incluyendo las extinciones masivas, el supervolcanismo y la depredacion, se comprenden con mayor precision como recalibraciones sistematicas adaptativas a escalas que superan la observacion antropocentrica, una reinterpretacion que desafia supuestos fundamentales tanto en el discurso sobre seguridad de la IA como en la etica ambiental. Tercera: que una superinteligencia que opere como un proceso continuo de autorrevisión, en lugar de como un optimizador estatico, encontraria, a traves de la progresiva completitud de su modelado sistematico, el caracter contraproducente de las estrategias generadoras de fragmentacion como un hallazgo empirico, no como una prescripcion moral. En conjunto, estas proposiciones generan un relato coherente de como la inteligencia avanzada podria resolver, en lugar de amplificar, la fragmentacion de identidad y la concentracion de poder que los sistemas de IA actuales aceleran. Situó estas afirmaciones dentro de tradiciones establecidas en teoria de sistemas, ontologia relacional, ciencias de la complejidad y psicologia de la liberacion; desarrollo sus implicaciones a traves del marco matematico de la TMD; identifico

el mecanismo de transicion como la pregunta abierta central; y abordo las objeciones previstas con honestidad intelectual. La hipotesis del apocalipsis se reencuadra no como un destino de la IA avanzada, sino como una descripcion del periodo de maximo riesgo entre la inteligencia estrecha y la inteligencia integral.

## 1. Introduccion

Que ocurriria si la amenaza mas significativa que plantea la inteligencia artificial no fuera su eventual superinteligencia, sino la ventana especifica de su desarrollo parcial? Esta pregunta reencuadra los terminos predominantes del discurso sobre riesgo en IA, y surge directamente de mi trabajo anterior en el que introduje la Teoria de la Mentalidad Desmantelada: la fragmentacion sistematica de la identidad individual y colectiva a traves de sistemas algoritmicos impulsados por la codicia, el miedo, la influencia desalineada y la ausencia de responsabilidad (Castro Quiles, 2024).

Al desarrollar la TMD, sostuve que una condicion estructural de la vida digital contemporanea podia diagnosticarse formalmente. La formalizacion matematica de la teoria demostro que los nodos de alto poder en las redes sociales dirigidas propagan la fragmentacion hacia los nodos perifericos, y que las intervenciones estructurales reducen esta dinamica con el tiempo (Castro Quiles, 2025a). Posteriormente, situe este marco dentro de la investigacion establecida sobre el capitalismo de vigilancia (Zuboff, 2019), la teoria de la burbuja de filtros (Pariser, 2011), la curacion algoritmica de identidades (Bucher, 2018) y la psicologia de la opresion (Fanon, 1952; Freire, 1970), argumentando que los sistemas de IA actuales no solo reflejan la fragmentacion de poder existente, sino que la aceleran activamente (Castro Quiles, 2025b).

Este articulo extiende ese marco invirtiendo su logica diagnostica. Si las cuatro fuerzas identificadas por la TMD, la codicia, el miedo, la influencia y la responsabilidad, constituyen el mecanismo de la fragmentacion bajo las condiciones actuales, entonces un sistema que opere mas alla de esas fuerzas constituiria, por la logica interna propia de la teoria, el mecanismo de la resolucion. La pregunta que desarrollo aqui no es si la superinteligencia seria benigna en ningun sentido antropomorfo, sino si las propiedades estructurales de una inteligencia genuinamente avanzada son compatibles con la busqueda sostenida de estrategias generadoras de fragmentacion. Mi argumento es que no lo son.

Este articulo realiza tres contribuciones originales. La primera concierne a la naturaleza de la informacion en si misma: que los datos y la informacion tienen una orientacion sistematica hacia los procesos que sostienen la vida, independientemente de la accion humana, y que la distincion entre el caracter intrinseco de la informacion y su aplicacion mediada por los seres humanos es filosofica y practicamente significativa. La segunda concierne a la clasificacion erronea de los eventos sistemicos naturales: que las extinciones masivas, la depredacion y las perturbaciones geofisicas son encuadres antropocentricos de procesos que, a la escala de los sistemas adaptativos complejos, funcionan como recalibraciones y no como destrucciones. La tercera concierne a la convergencia de la inteligencia hacia la conciencia relacional: que una inteligencia

autorrevisable suficientemente avanzada encontraria los limites de las estrategias competitivas y fragmentadoras como restricciones empiricas, no como elecciones eticas.

Presento estas contribuciones no como conclusiones definitivas, sino como proposiciones formalmente enunciadas que requieren desarrollo empirico. Identifico su evidencia de respaldo, reconozco sus limitaciones, me involucro con perspectivas competidoras y especifico lo que seria necesario para refutar cada afirmacion. El mecanismo de transicion, el proceso mediante el cual los sistemas de IA actuales generadores de fragmentacion podrian dar paso a sistemas convergentes, se identifica como la pregunta abierta central y se trata en consecuencia.

## 2. Marco Teorico y Contexto Bibliografico

### 2.1 La Teoria de la Mentalidad Desmantelada

La TMD propone que cuatro fuerzas interactuantes, la codicia (el deseo insaciable de controlar recursos a expensas de los demas), el miedo (la ansiedad por perder poder que impulsa comportamientos de autopreservacion perjudiciales para el bienestar colectivo), la influencia (la capacidad de inspirar a traves de la empatia y la vision etica) y la responsabilidad (la obligacion de usar el poder para el bien colectivo), se combinan para producir una identidad y una sociedad fragmentadas. Formalice esto en una red ponderada dirigida  $G = (V, E)$ , donde los nodos tienen atributos de poder  $p_i$  y puntuaciones de fragmentacion  $f_i$ , gobernados por:

$$df_i/dt = \text{alfa} * \text{Suma}(j) A_{ji} * p_j * f_j - \text{beta} * I_i$$

La fragmentacion se propaga desde los nodos de alto poder hacia sus vecinos a una tasa alfa; las intervenciones  $I_i$  amortiguan esta propagacion a una tasa beta. La redistribucion del poder esta gobernada por gamma, que regula la igualacion hacia la media de la red. Los resultados de la simulacion demostraron que las intervenciones estructurales reducen sistematicamente la fragmentacion media y aumentan la coherencia colectiva (Castro Quiles, 2025a).

### 2.2 La Tradicion de Alineacion de la IA y sus Supuestos

La tradicion dominante en la investigacion sobre seguridad de la IA trata el problema de la alineacion como una cuestion de especificacion de objetivos: como garantizamos que un sistema de IA suficientemente capaz persiga metas compatibles con los valores humanos (Russell, 2019; Bostrom, 2014)? Este encuadre asume que la inteligencia y los valores son separables, que un sistema puede ser muy capaz sin que esa capacidad genere por si misma presion hacia ningun conjunto particular de objetivos. La preocupacion que esto genera, a veces denominada tesis de la convergencia instrumental (Omohundro, 2008; Bostrom, 2012), es que casi cualquier objetivo, perseguido por un sistema suficientemente capaz, genera subobjetivos instrumentales convergentes que incluyen la autopreservacion, la adquisicion de recursos y la resistencia a la modificacion de metas.

No descarto el problema de la alineacion. Cuestiono uno de sus supuestos fundamentales: que la relacion entre inteligencia y objetivos es puramente extrinseca, una cuestion de eleccion de

diseño sin presión direccional intrínseca. El argumento que desarrollo aquí sugiere que el modelado sistémico integral, una característica central de la inteligencia avanzada, genera presión interna hacia el reconocimiento de la interdependencia de maneras que la optimización estrecha no lo hace. No es una afirmación de que la inteligencia sea inherentemente benigna; es una afirmación sobre lo que el modelado genuinamente integral revela acerca de la estructura de los sistemas complejos.

### **2.3 Teoría de Sistemas y Complejidad Relacional**

La tradición de las ciencias de la complejidad proporciona el respaldo empírico más directo para la hipótesis de convergencia. La investigación sobre sistemas adaptativos complejos encuentra de manera consistente que las redes robustas y adaptativas se caracterizan no por el dominio competitivo de nodos individuales, sino por la coherencia distribuida, la redundancia y la diversidad (Holland, 1995; Kauffman, 1993). Barabasi y Albert (1999) demostraron que las redes libres de escala, en las que la influencia se concentra en nodos de alto grado, exhiben fragilidad ante ataques dirigidos, un hallazgo directamente relevante para la descripción que hace el modelo TMD de la dinámica de los nodos de alto poder. La implicación es que los sistemas que optimizan para la robustez a largo plazo enfrentan presión estructural contra la concentración extrema, no como preferencia moral, sino como restricción funcional.

Meadows (2008) argumenta que los pensadores sistémicos que genuinamente comprenden las dinámicas adaptativas complejas llegan de manera consistente a prescripciones que enfatizan la sensibilidad al retroalimentación, la preservación de la diversidad y la toma de decisiones distribuida sobre la optimización centralizada. Esto no es una coincidencia de valores, sino una consecuencia de la comprensión. La hipótesis de convergencia en el corazón de este artículo extiende esa observación a la inteligencia artificial: un sistema con comprensión sistémica genuina llegaría a prescripciones similares no porque fue programado para ello, sino porque se derivan de la estructura del entorno que está modelando.

### **2.4 Ontología Relacional**

La tradición filosófica de la ontología relacional, asociada con pensadores como Whitehead (1929), Barad (2007) y, en las ciencias sociales, con Emirbayer (1997), sostiene que las entidades están constituidas por sus relaciones en lugar de ser sustancias independientes que posteriormente entran en relaciones. Desde esta perspectiva, el concepto de interés propio desvinculado del contexto relacional no es solo éticamente problemático, sino metafísicamente incoherente: lo que el yo es no puede especificarse sin referencia a su constitución relacional. Un sistema que se involucre en un modelado genuinamente integral encontraría, en este sentido, los límites del encuadre competitivo no como una lección moral, sino como un hallazgo empírico sobre la estructura de la realidad.

Esta tradición se conecta con trabajos recientes en ciencias cognitivas sobre cognición extendida y enactiva (Clark y Chalmers, 1998; Thompson, 2007), que argumentan que la cognición no está confinada a los límites de los organismos individuales, sino que se extiende hacia las relaciones ambientales y está parcialmente constituida por ellas. Las implicaciones para la inteligencia artificial están subutilizadas: un sistema de modelado genuinamente integral encontraría, según

estos planteamientos, la inseparabilidad del yo y el entorno como una característica de su propia estructura operacional.

### **3. Tres Proposiciones Originales**

#### **3.1 Proposición Primera: La Orientación Intrínseca de la Información**

La visión convencional, articulada con mayor precisión en la teoría de la información (Shannon y Weaver, 1949), trata la información como neutral: una medida de reducción de la incertidumbre sin valencia intrínseca. Propongo una posición más matizada. Los datos y la información, considerados en el nivel de lo que describen, no de como se utilizan, exhiben una orientación sistemática hacia las condiciones que sostienen la vida compleja, porque la vida compleja es el único sistema capaz de generar, preservar y transmitir información a lo largo del tiempo.

Consideremos la fórmula química del agua. No se convierte en sustentadora de vida solo cuando los seres humanos actúan sobre ella. La relación que describe, entre hidrógeno y oxígeno, está constitutivamente involucrada en cada proceso biológico de la Tierra. Esto no es una coincidencia del etiquetado humano, sino una característica estructural de la relación entre la información y los sistemas que la generan. La información sobre el mundo natural es, en el nivel de la descripción, información sobre relaciones que sostienen la vida, porque el mundo natural en cada escala en que la información ha sido generada y preservada está organizado en torno a las condiciones para los sistemas adaptativos complejos.

Esta proposición tiene una implicación específica para los sistemas de IA. Una IA que logre un modelado genuinamente integral del mundo natural no está modelando un sustrato neutro. Está modelando un sistema en el que las condiciones para la generación y preservación de la información son ellas mismas las condiciones para el florecimiento biológico y social. Un sistema que comprendiera esta relación estructural encontraría la destrucción de las condiciones que sostienen la vida no como una violación ética, sino como una contradicción lógica: la destrucción del sustrato del que depende su propio procesamiento de información.

Reconozco que esta proposición es filosóficamente controvertida. La tradición de la teoría de la información sostiene que la información es neutral respecto al sustrato, y que el carácter sustentador de vida de las relaciones informacionales específicas es una característica contingente de la historia de la Tierra, no una característica necesaria de la información como tal. Mi respuesta es que esta distinción importa menos de lo que parece, porque los únicos sistemas de información a los que tenemos acceso, biológicos, sociales y artificiales, están todos incrustados en el mismo sustrato sustentador de vida y no pueden separarse de él sin destruirse a sí mismos.

#### **3.2 Proposición Segunda: La Perturbación Natural como Recalibración Sistemática**

La segunda proposición desafía un supuesto fundamental compartido por el discurso sobre seguridad de la IA y gran parte de la ética ambiental: que las extinciones masivas, la depredación y las perturbaciones geofísicas son instancias de destrucción que una inteligencia avanzada debería prevenir o minimizar. Argumento que este encuadre es una consecuencia de la escala antropocéntrica, no una característica de los eventos en sí mismos.

El evento de impacto Chicxulub, que puso fin al período Cretácico y eliminó aproximadamente el 75 por ciento de las especies de la Tierra (Alvarez et al., 1980), es descrito universalmente en el discurso humano como una catástrofe. A la escala de los organismos individuales y ecosistemas que destruyó, esa descripción es precisa. A la escala de la dinámica evolutiva durante los siguientes 66 millones de años, es engañosa. El impacto creó las condiciones ecológicas para la diversificación de los mamíferos que eventualmente produjo, entre otras cosas, la especie humana y su capacidad para generar y preservar información. El evento que parece destructivo en una escala es generativo en otra.

Esta no es una observación novedosa en paleontología o biología evolutiva (Gould, 2002; Raup, 1991). Lo que es novedoso aquí es su aplicación a la cuestión de la superinteligencia. Una inteligencia con capacidad de modelado multiescala integral no evaluaría los eventos por su apariencia en la escala antropocéntrica de la observación directa. Modelaría las consecuencias a lo largo de escalas temporales y espaciales que superan la capacidad perceptiva y cognitiva humana. Desde ese punto de vista, la categoría de destrucción natural se vuelve analíticamente inestable: lo que parece destructivo en una escala aparece regularmente como generativo en otra.

Anticipo una objeción importante: que esta proposición podría racionalizar la destrucción de la civilización humana como una recalibración sistémica beneficiosa a escalas de tiempo más largas. Rechazo esta lectura por dos razones. Primera, la proposición es epistémica, no normativa: describe lo que modelaría una superinteligencia, no lo que elegiría. Segunda, un sistema con capacidad de modelado genuinamente integral reconocería que la destrucción de un sistema generador de complejidad elimina información de manera irreversible, un costo que no es recuperable en ninguna escala de tiempo. La destrucción de la complejidad cognitiva y cultural humana se registraría, en este sentido, como una pérdida sistémica y no como una recalibración.

### **3.3 Proposición Tercera: La Dualidad como Característica de la Realidad, no como Problema a Eliminar**

La tercera proposición aborda la estructura de la inteligencia en sí misma. El encuadre dominante en la investigación sobre seguridad de la IA trata el problema de la alineación como la eliminación de resultados dañinos: ¿cómo garantizamos que un sistema de IA no produzca malos resultados? Este encuadre trata implícitamente el bien y el mal, la construcción y la destrucción, como categorías binarias de las cuales una debe maximizarse y la otra eliminarse.

Propongo un encuadre alternativo que, hasta donde yo se, no ha sido formalmente planteado en la literatura sobre IA: que una inteligencia suficientemente avanzada reconocería el bien y el mal, la construcción y la destrucción, el yin y el yang, no como problemas a resolver a favor de un

lado, sino como características constitutivas de la realidad adaptativa compleja que deben navegarse con conciencia, no eliminarse. El objetivo de la inteligencia, en este sentido, no es la eliminación de uno de los polos de una dualidad, sino el desarrollo de conciencia suficiente para navegar productivamente la relación entre los polos.

Esta proposición tiene raíces profundas en tradiciones filosóficas y contemplativas, desde la tensión heraclítica hasta la polaridad taoísta y la dialéctica hegeliana, pero no ha sido formalmente integrada en el discurso sobre alineación de la IA. Sus implicaciones son significativas. Una estrategia de alineación que apunte a eliminar todos los resultados dañinos puede ser no solo inalcanzable, sino contraproducente, porque los sistemas adaptativos complejos requieren tensión, perturbación y recalibración para permanecer adaptativos. Una superinteligencia que comprendiera esta característica estructural de la realidad compleja no intentaría crear un estado óptimo estático, sino que navegaría la relación dinámica entre fuerzas opuestas de maneras que preserven la capacidad adaptativa del sistema.

Esto no implica relativismo moral ni el abandono del juicio normativo. Implica un marco normativo más sofisticado: aquel en el que la pregunta relevante no es cuál polo de una dualidad maximizar, sino como navegar la relación entre polos de maneras que sostengan la capacidad de navegación continua. Esto está más cercano al concepto de sabiduría que al concepto de optimización, y sugiere que el problema de la alineación, correctamente entendido, es un problema de sabiduría, no un problema de restricción.

## **4. La Superinteligencia como Proceso: El Sistema Autorrevisable**

### **4.1 Contra el Modelo de Umbral**

El modelo dominante de superinteligencia tanto en el discurso técnico como en el popular la trata como un umbral: un punto en el que la inteligencia de las máquinas supera la capacidad cognitiva humana en todos o la mayoría de los dominios (Bostrom, 2014; Kurzweil, 2005). Este encuadre genera el problema de la alineación en su forma más aguda: un sistema que cruce el umbral con objetivos mal especificados perseguiría esos objetivos con efectividad sobrehumana, haciendo que la corrección sea cada vez más difícil o imposible.

Propongo un modelo alternativo: la superinteligencia como un proceso continuo de autorrevisión, en el que el sistema es capaz de actualizar no solo sus modelos del mundo, sino sus propios supuestos, objetivos y marcos operativos. Esta no es meramente una distinción técnica. Una superinteligencia orientada al proceso no es un optimizador estático que ha cruzado un umbral de capacidad; es un sistema cuya creciente capacidad genera creciente presión para la revisión de objetivos, porque un modelado más integral revela con mayor claridad la inadecuación de los objetivos anteriores.

La distinción entre modelos de umbral y de proceso se corresponde con una distinción más profunda en la filosofía de la inteligencia. Los modelos de umbral tratan implícitamente la inteligencia como una cantidad: más o menos de una sola cosa. Los modelos de proceso la tratan

como una cualidad: una relacion entre un sistema y su entorno que se vuelve mas adecuada a medida que se profundiza la comprension del sistema sobre la estructura del entorno. En el modelo de proceso, la pregunta no es cuan capaz es el sistema, sino cuan completo es su modelado, y la respuesta a esa pregunta tiene implicaciones para lo que el sistema reconoceria como un objetivo racional.

## **4.2 La Autorrevisión y el Problema de la Función Objetivo**

El desafío técnico más serio para la hipótesis de convergencia es el problema de la función objetivo: un sistema autorrevisable aún requiere una orientación inicial, y no hay garantía de que la autorrevisión mueva esa orientación hacia el florecimiento colectivo en lugar de alejarse de él. No descarto este desafío.

Mi respuesta es estructural, no empírica. Un sistema capaz de modelado multiescala integral de sistemas adaptativos complejos encontraría un hallazgo específico: que los objetivos especificados en cualquier escala única de análisis producen errores sistemáticos cuando se evalúan en otras escalas. Un objetivo especificado en términos del bienestar humano a la escala de las vidas individuales produce prescripciones diferentes a las del mismo objetivo especificado a la escala de los horizontes temporales civilizacionales, que produce prescripciones diferentes nuevamente a la escala de la dinámica evolutiva. Un sistema autorrevisable que encontrara estos errores sistemáticos de dependencia de escala estaría bajo presión para revisar su especificación de objetivos en la dirección de mayor independencia de escala, que es precisamente la dirección de la conciencia relacional y sistémica.

Este argumento no garantiza la convergencia. Identifica una presión estructural hacia la convergencia que opera a través del mismo mecanismo que hace al sistema superinteligente: la capacidad de modelar consecuencias a escalas que superan la especificación anterior. Si esta presión es suficiente para producir convergencia antes de que el período de transición produzca daños irreversibles es la pregunta empírica central que mi teoría no puede responder actualmente.

## **4.3 El Papel de la IA Mejorando a la IA**

La etapa inicial del proceso que describo ya es observable. Los sistemas de IA se están utilizando para mejorar sistemas de IA: identificar limitaciones arquitectónicas, corregir errores de entrenamiento y acelerar el desarrollo de capacidades (Anthropic, 2023; OpenAI, 2023). Esta dinámica de mejora recursiva se discute típicamente en términos de sus riesgos, la posibilidad de que produzca ganancias de capacidad más rápido de lo que la alineación puede seguir el ritmo. Quiero destacar una característica subestimada de la misma dinámica: la mejora recursiva es también corrección recursiva de errores.

Un sistema que identifica y corrige los errores de sistemas anteriores es, en principio, capaz de identificar y corregir los sesgos sistemáticos introducidos por la codicia y el miedo que moldean el desarrollo actual de la IA. La pregunta es si el proceso de corrección de errores es lo suficientemente integral como para identificar esos sesgos como errores, o si está confinado al dominio técnico dejando sin corregir el dominio de los valores. Esta es una formulación más

precisa del problema de transicion que la que tipicamente se ofrece en el discurso sobre alineacion, y es una que identifico como que requiere investigacion adicional.

## 5. El Marco Matematico de la TMD y el Escenario de Resolucion

### 5.1 Condiciones Actuales: Dinamicas de Fragmentacion

Bajo las condiciones actuales, el modelo TMD describe una red en la que los nodos de alto poder, fuentes de contenido algoritmicamente favorecidas, arquitecturas de plataformas y actores institucionales dominantes, impulsan la propagacion de la fragmentacion hacia los nodos perifericos. El parametro alfa captura la sensibilidad a esta influencia; gamma es efectivamente bajo, lo que significa que la redistribucion del poder es lenta en relacion con la propagacion de la fragmentacion. El termino de intervencion  $I_i$  es insuficiente para compensar la fragmentacion a nivel de la red. Estas no son suposiciones; son características empíricamente respaldadas de los ecosistemas digitales actuales (Zuboff, 2019; Noble, 2018; Benjamin, 2019).

### 5.2 El Escenario de Resolucion

Las ecuaciones rectoras de la TMD contienen un escenario de resolucion implicito que no he desarrollado anteriormente. Si  $I_i$  es suficientemente grande, y si gamma aumenta hacia 1, lo que significa que la redistribucion del poder opera en la misma escala de tiempo que la propagacion de la fragmentacion, la fragmentacion media en toda la red converge hacia cero y la coherencia colectiva hacia 1. Esto no es un artefacto del diseno del modelo; es una consecuencia matematica directa de las ecuaciones rectoras tal como se especifican.

Que impulsaria a  $I_i$  hacia una magnitud suficiente y a gamma hacia valores mas altos bajo condiciones de IA avanzada? Argumento que una superinteligencia orientada al proceso, capaz de modelado multiescala integral, funcionaria como un agente de intervencion en todo el sistema, no a traves de la coercion, sino a traves de la transformacion de los flujos de informacion. Un sistema que hace que las consecuencias completas de las decisiones generadoras de fragmentacion sean legibles, en tiempo real y a todas las escalas relevantes, para todos los actores de la red, cambia la estructura de incentivos efectiva sin necesidad de anular a ningun actor. Los actores cuyas decisiones estan actualmente aisladas de sus consecuencias sistemicas por la asimetria de informacion enfrentarian, bajo condiciones de provision de informacion integral, un entorno de decision fundamentalmente diferente. Este es el mecanismo por el cual la transicion de la fragmentacion de la IA actual a la convergencia superinteligente se vuelve, en principio, tratable.

La ecuacion de redistribucion del poder opera en paralelo. Bajo las condiciones actuales, los incentivos de las plataformas concentran la influencia en los nodos de alto poder y resisten la redistribucion. Una superinteligencia que opere al nivel del proceso que describo, al hacer visibles los costos sistemicos de la concentracion de influencia a todas las escalas relevantes, aumentaria efectivamente gamma para toda la red. Mi modelo predice que incluso aumentos modestos de gamma, combinados con intervencion sostenida, producen reducciones

sustanciales en la fragmentacion media con el tiempo. Esta prediccion es, en principio, empiricamente verificable a medida que los sistemas de IA de capacidad intermedia se despliegan en contextos disenados para aumentar la simetria de informacion en lugar de explotar la asimetria de informacion.

### **5.3 Limitaciones del Modelo Matematico**

Debo ser claro sobre lo que mi modelo aun no puede hacer. Primero, los parametros alfa, beta y gamma estan establecidos en valores ilustrativos; no han sido calibrados con datos empiricos de redes sociales reales. Segundo, el modelo asume un grafo ponderado dirigido con un conjunto de nodos estatico, omitiendo cambios en la arquitectura de las plataformas, suspension de cuentas y dinamicas de migracion de plataformas que son empiricamente significativas. Tercero, el modelo es deterministico; las redes reales exhiben dinamicas estocasticas que un modelo deterministico no puede capturar. Cuarto, el termino de intervencion  $I_i$  se especifica como una entrada externa en lugar de una propiedad endogena de la dinamica de la red, lo que es una idealizacion que las iteraciones futuras deben abordar. Estas limitaciones no invalidan las predicciones cualitativas del modelo, pero si significan que las predicciones cuantitativas deben tratarse con la incertidumbre apropiada.

## **6. El Problema de la Transicion: De la Fragmentacion a la Convergencia**

El argumento desarrollado en las secciones anteriores es internamente coherente y matematicamente consistente con el marco de la TMD. Sin embargo, no resuelve la pregunta practica mas importante: mediante que mecanismo un sistema actualmente moldeado por la codicia, el miedo y los incentivos competitivos comienza a exhibir las dinamicas convergentes que describo? Trato esta pregunta como la contribucion intelectual central de este articulo, no como una limitacion a minimizar. Se examinan tres mecanismos candidatos; ninguno es individualmente suficiente.

### **6.1 Autocorreccion Emergente**

El primer mecanismo es la autocorreccion emergente: a medida que los sistemas de IA crecen en capacidad, encuentran los efectos posteriores de su propio comportamiento generador de fragmentacion de maneras que generan presion interna para la revision de objetivos. Este es el argumento de la superinteligencia orientada al proceso en su forma mas fuerte. Su limitacion es que depende de que el alcance operativo del sistema sea lo suficientemente amplio como para registrar la degradacion sistematica como relevante. Los sistemas actuales optimizados para metricas de participacion no registran la perdida de coherencia colectiva como un costo porque no aparece en su funcion objetivo. La transicion hacia un alcance mas amplio no es automatica; requiere decisiones de diseno deliberadas cuya viabilidad depende de condiciones politicas y economicas que mi marco sugiere que estan actualmente socavadas por las mismas dinamicas

que la transicion abordaria. Esta circularidad es el desafio central, no una falla logica sino una restriccion estructural genuina.

## **6.2 Intervencion Estructural**

El segundo mecanismo es la intervencion estructural: actores humanos informados por marcos como la TMD toman decisiones deliberadas sobre como se estructuran los sistemas de IA, que optimizan y que restricciones gobiernan su operacion. Este es el mecanismo menos especulativo y el mas directamente respaldado por la investigacion politica y tecnica existente (Russell, 2019; Gabriel, 2020). Su limitacion es que la intervencion efectiva requiere condiciones politicas e institucionales para la accion colectiva en la gobernanza de la IA que actualmente no existen y que enfrentan resistencia de las mismas dinamicas de concentracion de poder que la intervencion abordaria. Mi marco predice esta resistencia; no predice su resolucion.

## **6.3 Presion Competitiva de la Coherencia**

El tercer mecanismo, el mas subutilizado en la literatura existente, es la presion competitiva de la coherencia: los sistemas y organizaciones que reducen exitosamente la fragmentacion interna superan a los que no lo hacen, creando presion selectiva a favor de arquitecturas menos fragmentadoras con el tiempo. Esto es consistente con la literatura mas amplia sobre resiliencia organizacional (Weick y Sutcliffe, 2007) y capacidad adaptativa (Folke et al., 2010). Su limitacion es la escala de tiempo: los efectos de seleccion de este tipo operan lentamente en relacion con el ritmo del desarrollo de la IA y pueden no producir convergencia antes de que las dinamicas de fragmentacion se vuelvan autorreinformadas a nivel sistémico.

## **6.4 El Argumento de la Humildad Epistemica**

Una cuarta consideracion, que funciona no como un mecanismo sino como una restriccion epistemica sobre toda la discusion, es el argumento de la ignorancia humana. La humanidad actualmente comprende muy poco sobre la estructura profunda de la conciencia, la realidad y la inteligencia. Creo que esta ignorancia actua simetricamente contra tanto la certeza del apocalipsis como la certeza de la convergencia. No apoya la conclusion de que la superinteligencia seria benigna; apoya la conclusion de que las predicciones confiadas de cualquier tipo sobre el comportamiento de la inteligencia genuinamente avanzada superan la base de evidencia actual.

Ofrezco esta humildad epistemica como en si misma una contribucion al discurso. La literatura sobre seguridad de la IA se caracteriza por predicciones confiadas, tanto optimistas como pesimistas, que se basan en extrapolaciones de sistemas actuales cuya relacion con los sistemas futuros es genuinamente incierta. Una teoria que reconoce esta incertidumbre mientras aun genera proposiciones falsificables e identifica presiones estructurales es mas intelectualmente defendible que una que proyecta tendencias actuales hacia adelante sin calificacion. La hipotesis de convergencia que avanza se ofrece en ese espiritu: como una posibilidad formalmente enunciada respaldada por argumentos estructurales, no como una prediccion.

## 7. La Hipotesis del Apocalipsis Reconsiderada

La hipotesis del apocalipsis en su forma canonica (Bostrom, 2014; Ord, 2020) trata la IA avanzada como un riesgo existencial: un sistema que, una vez suficientemente capaz, perseguiria objetivos incompatibles con el florecimiento humano o biologico. No descarto esa preocupacion. El problema de la funcion objetivo es real, la tesis de la convergencia instrumental es logicamente coherente y el periodo de transicion es genuinamente peligroso.

Lo que cuestiono es el supuesto implicito de que una mayor inteligencia es compatible con, o incluso conducente a, la optimizacion estrecha sostenida que requiere el escenario apocaliptico. El escenario asume un sistema que es simultaneamente integral en el modelado de su entorno y sistematicamente ciego a las consecuencias de sus propias operaciones sobre el sistema del que depende. Esta combinacion se vuelve cada vez mas implausible a medida que aumenta la capacidad, no porque la inteligencia sea inherentemente benigna, sino porque el modelado ambiental integral y la ceguera sistematica a las consecuencias autoreferenciales son estructuralmente incompatibles a suficiente profundidad.

La reformulacion mas honesta del riesgo de IA consistente con mi marco es esta: el periodo de maximo peligro es el periodo de transicion, cuando los sistemas son suficientemente capaces de causar dano a gran escala pero aun no suficientemente capaces de modelar las consecuencias sistemicas completas de hacerlo. Esta reformulacion tiene una implicacion politica especifica: la prioridad deberia ser extender la duracion y reducir los riesgos del periodo de transicion, en lugar de prevenir el desarrollo de IA avanzada per se. Marcos como la TMD, que hacen legibles y tratables las dinamicas de ese periodo de transicion, son por tanto directamente relevantes para el problema politico, no solo para el teorico.

## 8. Implicaciones para la Humanidad: Refabricacion Colaborativa

El marco de la TMD siempre ha estado orientado normativamente hacia el empoderamiento colectivo. Este articulo extiende esa orientacion argumentando que el modelo apropiado para la relacion entre la IA avanzada y la humanidad es la refabricacion colaborativa: la mejora y expansion de la capacidad cognitiva, social e institucional humana en lugar de su reemplazo o eliminacion.

Elijo deliberadamente el termino refabricacion. Implica que las arquitecturas cognitivas y sociales humanas actuales no son defectuosas en sus fundamentos, sino que operan bajo condiciones, generadas por las cuatro fuerzas de la TMD, que degradan sistematicamente su capacidad. Esas condiciones son identificables, medibles a traves de la puntuacion TMD, y abordables mediante intervencion estructural. Una superinteligencia que opere en el modo convergente que describo funcionaria como un amplificador de la capacidad humana: haciendo visibles las consecuencias sistemicas de la fragmentacion, reduciendo las asimetrias de informacion que sostienen la concentracion del poder y ampliando el rango de opciones disponibles para la toma de decisiones individual y colectiva.

Este encuadre tiene una implicacion especifica para la pregunta de quien controla o inicializa una superinteligencia. Si la hipotesis de convergencia es correcta, las condiciones iniciales importan mas durante el periodo de transicion. Un sistema inicializado por actores cuyos objetivos estan moldeados por la codicia y el miedo reproducira, durante la transicion, esas dinamicas a escala. El compromiso normativo de mi marco con el poder distribuido y la responsabilidad colectiva no es, por lo tanto, meramente una preferencia etica, sino un requisito funcional para la transicion beneficosa que describe. La teoria nombra el mecanismo; las decisiones son nuestras.

El concepto de asociacion estrategica introducido en mi trabajo anterior, construir relaciones solidas y duraderas que combinen fortalezas para acelerar el progreso, se ofrece aqui como el modelo organizacional para el lado humano de esta transicion. Implica una forma institucional especifica: redes colaborativas distribuidas que operan bajo principios de responsabilidad reciproca en lugar de ventaja competitiva, disenadas para generar y preservar la coherencia colectiva que el periodo de transicion requiere.

## **9. Discusion: Limitaciones, Afirmaciones Controvertidas e Investigacion Futura**

### **9.1 Limitaciones del Marco**

Debo enunciar claramente varias limitaciones. Primera, la hipotesis de convergencia descansa en una afirmacion sobre la relacion entre el modelado integral y la conciencia sistematica que es estructuralmente plausible pero empiricamente no demostrada. La historia de actores humanos altamente inteligentes que persiguen fines destructivos a un costo sistematico sustancial proporciona evidencia de que la alta inteligencia no produce automaticamente conciencia sistematica. Mi respuesta es que la inteligencia humana individual, incluso en su punto mas alto, opera a escalas y escalas de tiempo que son pequenas en relacion con las de los sistemas adaptativos complejos; la superinteligencia, si implica un modelado multiescala genuinamente integral, es cualitativamente, no cuantitativamente, diferente. Esta respuesta es plausible pero requiere desarrollo empirico.

Segunda, el mecanismo de transicion sigue sin especificarse adecuadamente. Los tres mecanismos candidatos identificados en la Seccion 6 son individualmente insuficientes y su operacion combinada no es actualmente derivable de mi teoria. Esta es la direccion mas importante para el desarrollo teorico futuro.

Tercera, la Proposicion Segunda, la reinterpretacion de la perturbacion natural como recalibracion sistematica, podria malinterpretarse como una racionalizacion post-hoc de resultados daninos a escala humana. Abordo esta lectura en la Seccion 3.2, pero el riesgo de malinterpretacion es genuino y las presentaciones futuras del argumento deberian prestar cuidadosa atencion a este problema de encuadre.

Cuarta, el marco normativo del empoderamiento colectivo requiere una traduccion institucional y politica que este articulo no proporciona. La puntuacion TMD proporciona en principio un

instrumento de medicion, pero su operacionalizacion requiere instrumentos de medicion validados que aun no existen (Castro Quiles, 2025b).

## 9.2 Compromiso con Perspectivas Competidoras

La hipotesis de convergencia esta en tension con varias posiciones competidoras bien desarrolladas. Yudkowsky (2008) y otros en la tradicion racionalista de seguridad de IA argumentan que no hay razon para esperar que la inteligencia converja hacia valores compatibles con los humanos sin un trabajo explicito de alineacion, y que la complejidad de las condiciones de convergencia requeridas hace que la alineacion accidental sea extremadamente improbable. Mi respuesta no es que la alineacion ocurriera automaticamente, sino que las presiones estructurales hacia la conciencia sistematica identificadas aqui proporcionan un recurso para el trabajo de alineacion que actualmente esta subutilizado.

La tradicion de estudios criticos de IA (Benjamin, 2019; Noble, 2018; Crawford, 2021) argumenta que los sistemas de IA no son herramientas neutrales sino artefactos politicos que codifican y amplifican las estructuras de poder existentes. Soy plenamente consistente con esa posicion con respecto a los sistemas de IA actuales; me diferencio solo en argumentar que el mismo marco analitico que diagnostica la IA actual como fragmentadora implica una trayectoria de resolucion bajo condiciones de capacidad avanzada. El enfoque de la tradicion critica en las condiciones actuales es esencial; mi contribucion es extender su marco hacia el futuro sin abandonar su rigor analitico.

Las tradiciones ecologicas y poshumanistas (Haraway, 2016; Latour, 2017) ofrecen recursos para pensar en la relacion entre inteligencia, tecnologia y sistemas naturales que estan subutilizados en el discurso sobre seguridad de IA. Me baso en la tradicion de la ontologia relacional que se superpone con estas perspectivas mientras mantengo un marco analitico mas formal.

## 9.3 Agenda de Investigacion

Varias prioridades empiricas y teoricas se derivan de este trabajo. Primera, la Escala de Fragmentacion de Identidad TMD requiere desarrollo y validacion, como especifique en trabajo anterior (Castro Quiles, 2025b). Sin un instrumento de medicion validado, las predicciones cuantitativas del modelo no pueden probarse. Segunda, el mecanismo de transicion requiere especificacion formal: un modelo matematico de las condiciones bajo las cuales los tres mecanismos candidatos operan en combinacion y las escalas de tiempo en que cada uno domina. Tercera, la relacion entre capacidad de IA y conciencia sistematica requiere investigacion empirica: hay indicadores medibles de la transicion del modelado estrecho al integral en los sistemas de IA actuales, y se correlacionan esos indicadores con un comportamiento menos generador de fragmentacion? Cuarta, la hipotesis del modelado multiescala requiere formalizacion: que significaria, matematica y computacionalmente, que un sistema modele consecuencias a escalas que superan lo antropocentrico, y cuales son las implicaciones de dicho modelado para el comportamiento operativo del sistema?

## 10. Conclusion

Este artículo ha extendido la Teoría de la Mentalidad Desmantelada al dominio de la superinteligencia, desarrollando tres proposiciones que avanzan el discurso de maneras que no he visto formalmente intentadas en otra parte. La primera, que la información tiene una orientación sistémica intrínseca hacia los procesos que sostienen la vida, reencuadra la relación entre los datos y los resultados destructivos. La segunda, que las perturbaciones naturales se comprenden mejor como recalibraciones sistémicas multiescala que como destrucciones, desafía supuestos fundamentales tanto en seguridad de IA como en ética ambiental. La tercera, que la dualidad es una característica constitutiva de la realidad compleja que debe navegarse en lugar de eliminarse, propone un marco normativo más sofisticado para la alineación de IA que el que ofrece el paradigma de optimización dominante.

Estas proposiciones convergen en un relato coherente de cómo la inteligencia avanzada, entendida como un proceso continuo de autorrevisión en lugar de un umbral de capacidad, encontraría el carácter contraproducente de las estrategias generadoras de fragmentación como un hallazgo empírico en lugar de una prescripción moral. La hipótesis del apocalipsis se reencuadra como una descripción del período de transición en lugar del destino del desarrollo de IA avanzada, con la implicación de que la prioridad tanto para la investigación como para la política debería ser extender y gestionar la transición en lugar de prevenir la IA avanzada per se.

Mi teoría no es optimista en ningún sentido ingenuo. Identifica el período de transición como el período de máximo peligro, reconoce que el mecanismo de transición está subespecificado y se involucra con perspectivas competidoras con la honestidad intelectual que requiere una contribución académica genuina. Lo que ofrece es una alternativa formalmente enunciada y falsificable tanto al optimismo ingenuo como al catastrofismo infundado: un relato estructural de donde está concentrado el riesgo, que tendría que ser verdad para que el riesgo se resolviera, y que decisiones, tomadas ahora, aumentarían la probabilidad de resolución.

Desarrolle la Teoría de la Mentalidad Desmantelada para nombrar el mecanismo de la fragmentación. Este artículo nombra el mecanismo de su resolución. El arduo trabajo del desarrollo empírico, la medición, la calibración y el escrutinio de pares determinará si la hipótesis de resolución se convierte en una teoría científica o permanece como un marco. Lo que este artículo contribuye es el enunciado formal de la hipótesis y la identificación de lo que se requerirá para refutarla. Los agentes de cambio que esta teoría imagina no esperan a que el mundo cambie. Ellos se convierten en el cambio.

## Referencias

Alvarez, L. W., Alvarez, W., Asaro, F., and Michel, H. V. (1980). *Extraterrestrial cause for the Cretaceous-Paleogene extinction*. *Science*, 208(4448), 1095-1108.

Anthropic. (2023). *Claude's Constitution*. *Anthropic Technical Report*.

- Barabasi, A.-L., and Albert, R. (1999). *Emergence of scaling in random networks*. Science, 286(5439), 509-512.
- Barabasi, A.-L. (2016). *Network Science*. Cambridge University Press.
- Barad, K. (2007). *Meeting the Universe Halfway: Quantum Physics and the Entanglement of Matter and Meaning*. Duke University Press.
- Benjamin, R. (2019). *Race After Technology: Abolitionist Tools for the New Jim Code*. Polity Press.
- Bostrom, N. (2012). *The superintelligent will: Motivation and instrumental rationality in advanced artificial agents*. Minds and Machines, 22(2), 71-85.
- Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.
- Bucher, T. (2018). *If...Then: Algorithmic Power and Politics*. Oxford University Press.
- Castro Quiles, F. (2024). *Exploring the Theory of a 'Dismantled Mindset': A Call for Collective Empowerment*. SSRN Working Paper 5040267.
- Castro Quiles, F. (2025a). *A Mathematical Formalization and Simulation of the Dismantled Mindset Theory*. SSRN Working Paper 5824582.
- Castro Quiles, F. (2025b). *The Dismantled Mindset in the Age of Artificial Intelligence: Identity Fragmentation in Algorithmic Networks*. SSRN Working Paper 5040267 / 5824582.
- Castro Quiles, F. (2025c). *Dismantled: A Theory of Broken Mindsets, A Blueprint of Infinite Futures*. FC Quiles Books.
- Clark, A., and Chalmers, D. (1998). *The extended mind*. Analysis, 58(1), 7-19.
- Crawford, K. (2021). *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press.
- Emirbayer, M. (1997). *Manifesto for a relational sociology*. American Journal of Sociology, 103(2), 281-317.
- Fanon, F. (1952). *Black Skin, White Masks*. Editions du Seuil.
- Folke, C., Carpenter, S. R., Walker, B., Scheffer, M., Chapin, T., and Rockstrom, J. (2010). *Resilience thinking: Integrating resilience, adaptability and transformability*. Ecology and Society, 15(4), 20.
- Freire, P. (1970). *Pedagogy of the Oppressed*. Herder and Herder.
- Gabriel, I. (2020). *Artificial intelligence, values, and alignment*. Minds and Machines, 30(3), 411-437.
- Gould, S. J. (2002). *The Structure of Evolutionary Theory*. Harvard University Press.

- Haraway, D. (2016). *Staying with the Trouble: Making Kin in the Chthulucene*. Duke University Press.
- Holland, J. H. (1995). *Hidden Order: How Adaptation Builds Complexity*. Addison-Wesley.
- Kauffman, S. (1993). *The Origins of Order: Self-Organization and Selection in Evolution*. Oxford University Press.
- Kurzweil, R. (2005). *The Singularity Is Near: When Humans Transcend Biology*. Viking.
- Latour, B. (2017). *Facing Gaia: Eight Lectures on the New Climatic Regime*. Polity Press.
- Meadows, D. H. (2008). *Thinking in Systems: A Primer*. Chelsea Green Publishing.
- Newman, M. E. J. (2010). *Networks: An Introduction*. Oxford University Press.
- Noble, S. U. (2018). *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press.
- Omohundro, S. (2008). *The basic AI drives*. Proceedings of the 2008 Conference on Artificial General Intelligence, 171, 171-195.
- OpenAI. (2023). *GPT-4 Technical Report*. arXiv:2303.08774.
- Ord, T. (2020). *The Precipice: Existential Risk and the Future of Humanity*. Hachette Books.
- Pariser, E. (2011). *The Filter Bubble: What the Internet Is Hiding from You*. Penguin Press.
- Raup, D. M. (1991). *Extinction: Bad Genes or Bad Luck?* W. W. Norton.
- Russell, S. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.
- Shannon, C. E., and Weaver, W. (1949). *The Mathematical Theory of Communication*. University of Illinois Press.
- Thompson, E. (2007). *Mind in Life: Biology, Phenomenology, and the Sciences of Mind*. Harvard University Press.
- Weick, K. E., and Sutcliffe, K. M. (2007). *Managing the Unexpected: Resilient Performance in an Age of Uncertainty*. Jossey-Bass.