# **CAPÍTULO IV** LOS IMPULSORES DEL CRECIMIENTO DE LA IA

## **OBJETIVO**:

Este capítulo tiene como objetivo explicar por qué la inteligencia artificial ha alcanzado un punto de inflexión histórico en el desarrollo tecnológico global. A través del análisis de tres factores convergentes —el incremento exponencial en la capacidad de procesamiento, el surgimiento de modelos revolucionarios y la disponibilidad masiva de datos— se ofrece una visión integral de las condiciones que han permitido la expansión acelerada de la IA generativa.

El lector comprenderá por qué este es el momento decisivo para invertir, implementar y priorizar soluciones de IA en sectores clave como la educación, la industria, la salud y los negocios. Además, se examinarán aspectos estratégicos como la competencia tecnológica entre potencias globales, el papel de empresas como NVIDIA, y los avances que han llevado a redefinir lo que entendemos por procesamiento del lenguaje natural.

Este capítulo sienta las bases para comprender no solo el "cómo", sino el "por qué ahora", revelando el entramado técnico, económico y geopolítico que sustenta la revolución de la inteligencia artificial contemporánea.

## TEMAS A EXPLORAR:

- Por qué la inteligencia artificial ha alcanzado un punto de inflexión histórico y qué factores convergentes han impulsado su expansión actual.
- Cómo los avances en la capacidad de procesamiento (especialmente el rol de NVIDIA y las GPU) han permitido entrenar modelos más complejos y eficientes.
- De qué forma la evolución de modelos como los Transformers ha revolucionado el procesamiento del lenguaje natural y dado paso a la IA generativa.
- Cuál es la importancia de la disponibilidad masiva de datos, tanto históricos como en tiempo real, para alimentar los sistemas de IA modernos.
- Qué implicaciones geopolíticas tiene el control de la infraestructura tecnológica de IA, especialmente en la rivalidad entre China y Estados Unidos.
- Qué desafíos éticos plantea el uso de grandes volúmenes de datos y cómo garantizar un uso responsable, justo y transparente.
- Por qué comprender estos impulsores es esencial para implementar estrategias de IA con visión ética, educativa y transformadora.

## ✓ CHATGPT EN ACCIÓN – Actividades:

- Motores del crecimiento: "¿Cómo han crecido datos, modelos y procesamiento?"
- 2. Caso NVIDIA: "¿Cómo pasó NVIDIA de juegos a líder en IA?"
- 3. Guerra tecnológica: "Impacto de la competencia China-EE.UU. en IA"
- 4. Ética y datos: "¿Qué principios éticos deben regir el uso de datos en IA?"
- 5. Resumen activo: "Resumen del capítulo 4 incluyendo NVIDIA, impulsores y ética"

**K HERRAMIENTAS IA: Notion AI** 

## 4.1 Introducción

Hablar de inteligencia artificial hoy no es hablar del futuro: es hablar del presente. Sin embargo, una de las preguntas más fundamentales que debemos hacernos es: ¿Por qué justo ahora la IA está transformando de manera tan acelerada nuestras vidas, nuestras industrias y nuestros sistemas educativos?



Figura 4.1. Convergencia para el avance de la IA generativa

La respuesta no se encuentra en un solo avance, sino en la convergencia histórica de tres grandes fuerzas que, tras años de evolución independiente, han alcanzado un punto de madurez simultáneo y explosivo:

## 4.1.1 Avances en la capacidad de procesamiento

La evolución del hardware ha sido determinante. El desarrollo de unidades de procesamiento gráfico (GPUs) de última generación —como las fabricadas por NVIDIA— ha permitido ejecutar cálculos complejos en paralelo, lo cual es esencial para el aprendizaje profundo (deep learning).

Lo que antes requería meses de entrenamiento computacional, hoy puede realizarse en semanas o incluso días, gracias al poder de cómputo actual, la computación distribuida y el acceso a infraestructura de nube altamente escalable.

Esto ha eliminado muchas de las barreras técnicas que impedían el despliegue masivo de soluciones basadas en IA.

#### 4.1.2. Modelos de nueva generación

Desde la publicación del artículo "Attention Is All You Need" en 2017, la arquitectura Transformer ha revolucionado el campo del procesamiento del lenguaje natural. Hoy contamos con modelos como GPT, BERT, Claude, Gemini o DALL·E, capaces de comprender el contexto, generar texto creativo, analizar sentimientos y adaptarse a diferentes tareas con una precisión sorprendente.

Ya no hablamos de IA estrecha con funciones limitadas. Hablamos de IA generativa, capaz de abordar múltiples tareas con un grado de generalización sin precedentes.

#### 4.1.3. Disponibilidad masiva de datos

Vivimos en la era del Big Data. Cada día generamos una avalancha de datos: publicaciones en redes, compras en línea, geolocalización, interacciones, imágenes, audio, texto. Gracias a técnicas avanzadas de minería de datos, estos flujos de información son aprovechados para entrenar modelos más robustos, adaptables y personalizados.

Además, el crecimiento del acceso a datasets abiertos, junto con herramientas para limpiar y etiquetar datos, ha democratizado el desarrollo de modelos incluso en universidades y pequeñas empresas.

#### Más allá de lo técnico: una transformación cultural

No solo es la tecnología lo que ha cambiado. También lo ha hecho la expectativa social. Hoy las personas esperan herramientas que sean más que funcionales: esperan inteligencia, personalización y asistencia en tiempo real. La IA ha venido a responder a esa demanda. Ya no se trata de una promesa futurista. Es una herramienta viva, poderosa y presente.

#### La gran pregunta

En este contexto de alineación tecnológica, económica y cultural, la pregunta ya no es si vamos a usar la inteligencia artificial, sino:

¿Cómo vamos a usarla para transformar positivamente la educación, la productividad y la innovación?

## 4.2 Avances en la capacidad de procesamiento

Uno de los principales impulsores del auge actual de la inteligencia artificial es el extraordinario salto en la capacidad de procesamiento computacional. El entrenamiento de un modelo moderno —en especial los modelos de lenguaje de gran escala (LLMs)— exige recursos monumentales: velocidad, potencia, almacenamiento, eficiencia energética y escalabilidad.

Durante la última década, hemos sido testigos de una transformación sin precedentes impulsada por chips especializados, innovaciones arquitectónicas y la expansión de la computación en la nube, lo que ha hecho posible lo que antes era impensable.

## 4.2.1 NVIDIA y su papel estratégico

En el corazón de esta revolución tecnológica se encuentra NVIDIA, la empresa que alguna vez fue conocida únicamente por tarjetas gráficas para videojuegos, pero que hoy lidera el mercado global de hardware para inteligencia artificial.

Su liderazgo se basa en tres pilares:

- GPUs de alto rendimiento, como las series Quadro, Tesla, A100 y H100, diseñadas para acelerar el entrenamiento y la inferencia de modelos.
- Plataformas de cómputo como NVIDIA DGX, que permiten a organizaciones ejecutar modelos complejos a gran escala en entornos locales o en la nube.
- Software especializado, como CUDA, cuDNN y TensorRT, que simplifican y optimizan el desarrollo de soluciones de IA.

# CES 2025: NVIDIA presentó Project Digits, una supercomputadora para IA

Caracas Digital 08 ene, 2025



El proyecto DIGITS de NVIDIA, con el nuevo superchip GB10, se presenta como el supercomputador de IA más pequeño del mundo capaz de ejecutar modelos de 200.000 parámetros.

Figura 4.2. Jen-Hsun Huang en CES 2025

Además, su colaboración estratégica con TSMC (Taiwan Semiconductor Manufacturing Company) ha garantizado una producción eficiente y constante de sus chips más avanzados.

#### 4.2.2 El fenómeno financiero de NVIDIA

El impacto de NVIDIA no se limita al plano técnico. Su crecimiento bursátil ha sido reflejo del interés global por la inteligencia artificial:

- 2023: aumento de más del 240% en el valor de sus acciones, impulsado por la demanda explosiva de chips para IA.
- 2024: crecimiento adicional del 170%, consolidando su posición como una de las compañías tecnológicas más valiosas del mundo.
- Mayo de 2025: tras una leve corrección, la acción se mantiene fuerte en \$114.50
   USD, acumulando un impresionante incremento del +1,519% en cinco años.

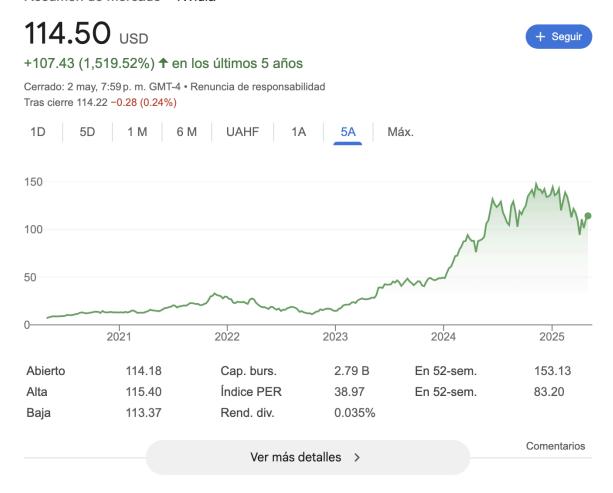


Figura 4.3. Evolución del valor de las acciones de NVIDIA (2019-2025)

Estos datos no solo evidencian el éxito de la compañía, sino el reconocimiento del papel central que juega el hardware especializado en la era de la inteligencia artificial.

## 4.2.3 Infraestructura y costos para el entrenamiento de modelos

Entrenar un modelo como GPT-4 no solo requiere un algoritmo avanzado, sino una infraestructura colosal, compuesta por:

- Hardware especializado
  - o GPUs como la NVIDIA H100
  - TPUs desarrolladas por Google
  - Servidores organizados en clústeres para procesamiento masivo
- I Almacenamiento y red
  - o Requiere sistemas de acceso rápido y redes de alta velocidad
- **a** Software y frameworks

- TensorFlow, PyTorch, JAX, Hugging Face Transformers
- H Datos
  - o Datasets masivos: artículos, libros, código, interacciones humanas
- Computación en la nube
  - AWS, Azure, Google Cloud eliminan la necesidad de infraestructura local

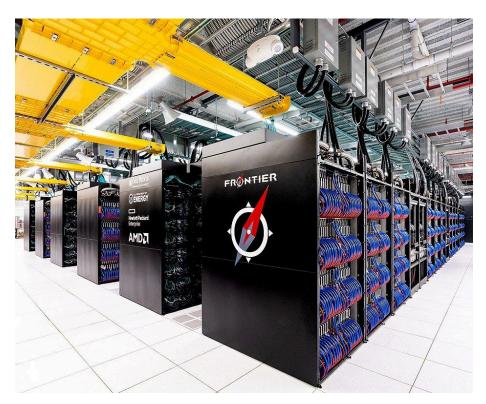


Figura 4.4. *HPE Cray EX*, supercomputador exascale del Oak Ridge National Laboratory (EE.UU.), usado para entrenar modelos de IA avanzada. Fuente: ORNL / HPE.

Costos estimados del entrenamiento de un LLM:

- 40–50% infraestructura (hardware)
- 10-20% energía
- 15–20% salarios expertos
- 5-10% datos
- 5-10% licencias de software
- 5−10% operación y mantenimiento

Un modelo como GPT-4 puede requerir más de 100 millones de dólares para ser entrenado, reflejando tanto su complejidad como su potencial transformador.

## 4.2.4 Evolución del rendimiento de los chips NVIDIA

En solo una década, los chips de NVIDIA han multiplicado por más de mil su capacidad de inferencia:

Tabla 4.1. Evolución del rendimiento de los chips NVIDIA

Año	Modelo de chip	Rendimiento TOPS
2012	NVIDIA K20X	3.94
2015	NVIDIA M40	6.84
2016	NVIDIA P100	21.20
2017	NVIDIA V100	125
2020	NVIDIA A100	1248
2023	NVIDIA H100	4000

La capacidad de inferencia en un chip se refiere a la velocidad y eficiencia con la que un procesador (como una GPU o un ASIC) puede ejecutar un modelo de inteligencia artificial ya entrenado para hacer predicciones o tomar decisiones. Es decir, mide cuántas operaciones puede realizar por segundo cuando "pone en acción" lo que ha aprendido.

Este salto ha sido posible gracias a avances en diseño, arquitectura y eficiencia energética.

#### Tipos de precisión:

- FP32: máxima precisión (entrenamiento)
- FP16: más velocidad con menor consumo
- Int8: ideal para inferencia en tiempo real

#### ¿Qué es Int8 TOPS?

- Int8: formato numérico eficiente de 8 bits
- TOPS: Tera Operations Per Second (billones de operaciones por segundo)
- Un chip con 100 INT8 TOPS puede realizar 100 billones de operaciones por segundo

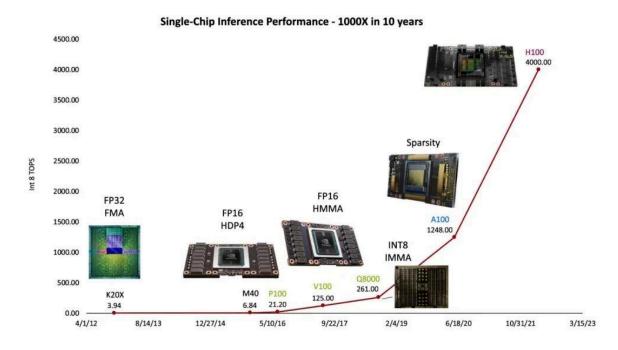


Figura 4.5 Rendimiento de inferencia por chip NVIDIA (2012–2023): mejora de 1000 veces en una década.

Estos valores permiten a modelos responder rápidamente, con bajo consumo energético y alta escalabilidad, ideal para educación, salud, negocios y creatividad.

## 4.2.5 Geopolítica del procesamiento: China vs. EE.UU.

## La carrera por los chips que mueven la inteligencia artificial

La revolución de la inteligencia artificial generativa no se sostiene sin una infraestructura invisible pero crítica: los chips. En particular, las **GPUs de alto rendimiento** que permiten entrenar modelos como GPT-4, Claude o DeepSeek. El control de esta tecnología no es menor: **quien domina los chips, domina el ritmo del avance tecnológico**.

## Nuevas restricciones de EE.UU. en 2024-2025

En abril de 2025, el gobierno de Estados Unidos asestó un nuevo golpe a las ambiciones tecnológicas de China: **prohibió formalmente la venta del chip H20 de NVIDIA**, un procesador diseñado específicamente para el mercado chino tras las restricciones anteriores sobre el H100 y A100.

#### @ ¿Por qué se restringió el chip H20?

- A pesar de estar técnicamente "por debajo" del umbral de las primeras sanciones, el H20 aún ofrecía suficiente potencia para entrenar modelos de IA de nivel competitivo.
- Las autoridades estadounidenses argumentaron riesgos de uso militar, vigilancia masiva y aplicaciones de doble uso.
- La presión no solo recayó sobre NVIDIA: también se impusieron límites a AMD (modelo MI308), entre otros fabricantes.

#### **>> Impacto inmediato:**

- NVIDIA reportó una caída bursátil de más del 10%, arrastrando al Nasdag.
- Se estiman pérdidas de hasta \$5.500 millones por pedidos congelados o cancelados.
- Proveedores y ensambladores tecnológicos asiáticos también sintieron el efecto dominó.

Esta medida refuerza una estrategia clara: **limitar el acceso de China a la infraestructura crítica para IA avanzada**, incluso si eso implica afectar a empresas estadounidenses de manera indirecta.



Figura 4.6 GPU (Unidad de Procesamiento Gráfico) de alto rendimiento NVIDIA

## La respuesta china: aceleración hacia la autosuficiencia

Lejos de frenarse, China ha redoblado su apuesta: si no puede comprar los chips más potentes, **los desarrollará por su cuenta**. Lo que comenzó como una reacción a las sanciones de 2019 y 2020, hoy es una **política industrial consolidada**.

#### 

- Huawei ha emergido como líder nacional en semiconductores con chips como Ascend 910B, optimizados para tareas de IA.
- ByteDance, Tencent y Alibaba están redirigiendo inversiones hacia proveedores nacionales y laboratorios propios de IA.
- **DeepSeek**, con su modelo V3 open-source, demostró que China puede desarrollar LLMs multilingües, competitivos y con capacidad de razonamiento avanzado, incluso con restricciones de hardware.

#### Nacional la lindependencia tecnológica en 3 frentes:

- 1. IA generativa y procesamiento de lenguaje natural
- 2. 5G y telecomunicaciones
- 3. Automatización y ciudades inteligentes

China no solo quiere competir. Quiere **superar el cuello de botella occidental** desarrollando una cadena de suministro propia, desde chips hasta plataformas de IA, pasando por sistemas operativos y hardware.

## X Una guerra silenciosa, pero global

Esta disputa ha dejado de ser una simple **guerra comercial**. Hoy es una batalla **por el alma de la próxima revolución industrial**, en la que el control de los datos, la infraestructura y los modelos de IA **se traduce directamente en poder geopolítico**.

## (in a section of the section of the

- Fragmentación del ecosistema tecnológico: hacia un mundo con infraestructuras paralelas (una occidental y otra china).
- Reducción de la cooperación científica internacional, especialmente en áreas como automatización, big data y IA aplicada.
- Mayor presión sobre países no alineados (Latinoamérica, África, Sudeste Asiático), que deberán elegir con qué ecosistema integrarse.

## ¿Qué está en juego?

Más allá del comercio, lo que está en juego es:

- Quién controla los estándares globales de IA
- Quién provee el "cerebro" de las plataformas educativas, sanitarias y de seguridad del futuro
- Quién lidera la narrativa sobre cómo debe comportarse una IA

En este contexto, no se trata solo de producir tecnología, sino de **influir en las reglas del juego**, en los valores que se programan en los modelos y en las decisiones que estos toman en nuestra vida diaria.

### 4.2.6 CPU vs. GPU: la demostración de MythBusters

Una de las comparaciones más ilustrativas se dio en el programa *MythBusters*, al enfrentar una CPU tradicional contra una GPU en tareas de procesamiento pesado:

Tabla 4.2. CPU versus GPU

Unidad	Característica principal	Tiempo estimado
CPU	Procesamiento secuencial	Minutos o horas
GPU	Procesamiento paralelo	Segundos

#### Las GPU superan a las CPU por:

- Paralelismo masivo (miles de núcleos)
- Optimización para operaciones simples
- Alto ancho de banda de memoria

### Aplicaciones prácticas:

- IA: entrenamiento de LLMs y visión artificial
- Videojuegos: renderizado en tiempo real
- Ciencia: simulación de sistemas complejos



Figura 4.7 NVIDIA GPU VS CPU Leonardo Paint the Mona Lisa in 80 Milliseconds!

Sin GPUs, la IA moderna simplemente no existiría en su forma actual. Las CPU siguen siendo esenciales, pero las GPU son el motor que impulsa esta nueva era tecnológica.

## 4.3 Desarrollo de modelos revolucionarios

El auge actual de la inteligencia artificial no puede entenderse únicamente desde el punto de vista de la potencia de cómputo o la disponibilidad de datos. En el núcleo de esta transformación se encuentran los modelos que han hecho posible que las máquinas no solo procesen información, sino que la comprendan, la generen y la adapten a múltiples contextos.

Durante la última década, se han producido avances radicales en las arquitecturas de aprendizaje profundo, permitiendo a los sistemas artificiales realizar tareas antes reservadas exclusivamente a los humanos. Esta sección aborda tres grandes pilares de esta revolución: los pioneros del deep learning, la arquitectura Transformer y la evolución del procesamiento de lenguaje natural (NLP).

### 4.3.1 Geoffrey Hinton: el padrino de la IA moderna

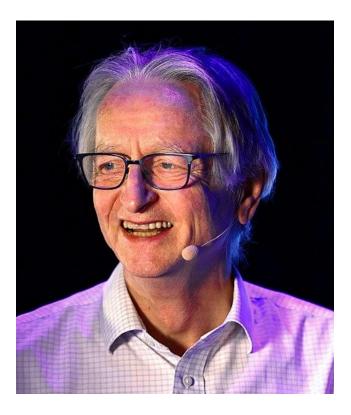


Figura 4.8 Geoffrey Hinton

Geoffrey Hinton es una figura central en la historia de la inteligencia artificial. Reconocido como uno de los padres del aprendizaje profundo (*deep learning*), su trabajo ha sido fundamental para reactivar el campo tras décadas de estancamiento.

**Retropropagación (backpropagation)**: Algoritmo que ajusta los pesos de una red neuronal en función del error, esencial para el entrenamiento de redes profundas.

**Máquinas de Boltzmann Restringidas (RBM)**: Método de aprendizaje no supervisado que permite detectar patrones complejos sin necesidad de datos etiquetados.

**Aplicaciones reales:** Sus investigaciones son la base de tecnologías actuales como el reconocimiento de voz, la visión por computadora y los modelos de lenguaje natural.

En 2018, Hinton recibió el Premio Turing, junto a Yoshua Bengio y Yann LeCun, por haber sentado las bases científicas del deep learning moderno. Sin su trabajo, los avances que hoy celebramos en IA generativa probablemente no habrían sido posibles.

### 4.3.2 Transformers: la arquitectura que lo cambió todo

En 2017, un equipo de Google publicó el influyente artículo "Attention Is All You Need", que introdujo la arquitectura Transformer, marcando un antes y un después en el procesamiento de lenguaje natural (NLP).

## **Attention Is All You Need**

Ashish Vaswani\* Google Brain avaswani@google.com Noam Shazeer\*
Google Brain
noam@google.com

Niki Parmar\* Google Research nikip@google.com

Jakob Uszkoreit\* Google Research usz@google.com

Llion Jones\*
Google Research
llion@google.com

Aidan N. Gomez\* † University of Toronto aidan@cs.toronto.edu

Łukasz Kaiser\*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin\* † illia.polosukhin@gmail.com

#### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

Figura 4.9 Portada del artículo original "Attention Is All You Need" (2017), donde se presenta la arquitectura Transformer

#### Q ¿Qué hace único al Transformer?

- Mecanismo de atención: Permite al modelo identificar qué partes de una secuencia son más relevantes en cada contexto, sin necesidad de procesar palabra por palabra de forma secuencial.
- Paralelismo y velocidad: Puede procesar grandes cantidades de texto simultáneamente, acelerando el entrenamiento y la inferencia.
- Estructura modular: Con codificadores y decodificadores que se adaptan a diversas tareas (traducción, resumen, generación).
- Codificación posicional: Permite al modelo entender el orden de las palabras sin depender de estructuras secuenciales.



Figura 4.10. 31.ª Conferencia sobre Sistemas de Procesamiento de Información Neural (NIPS 2017), Long Beach, California, EE. UU.

## Impacto del Transformer:

- Superó todos los modelos anteriores en tareas como traducción automática, generación de texto y análisis semántico.
- Dio origen a familias de modelos como GPT, BERT, T5 y T-NLG, que hoy están detrás de asistentes conversacionales, sistemas de búsqueda, motores de recomendación y herramientas de escritura automática.

En síntesis, el Transformer no solo representó una mejora en rendimiento: abrió la puerta a una nueva era de IA generativa.

### 4.3.3 Evolución del procesamiento de lenguaje natural (NLP)

El procesamiento de lenguaje natural ha sido uno de los campos que más ha evolucionado gracias al aprendizaje profundo. A través de sucesivas innovaciones, las máquinas han pasado de simplemente reconocer palabras a comprender intenciones, generar ideas y participar en conversaciones complejas.

Tabla 4.3. Línea de tiempo de hitos clave

Año	Hito	Impacto principal
2003	Modelos de lenguaje neuronal	Primeros intentos de predicción de texto
2008	Aprendizaje multitarea	Entrenamiento simultáneo en varias tareas (traducción, resumen)
2013	Word embeddings (Word2Vec)	Representación semántica de palabras como vectores
2014	Aprendizaje secuencia a secuencia (seq2seq)	Mejoras en traducción automática, chatbots
2015	Mecanismo de atención	Base conceptual de los Transformers
2018	Modelos preentrenados (BERT, GPT)	Adaptabilidad a múltiples tareas con un solo entrenamiento

Gracias a esta evolución, hoy es posible interactuar con sistemas capaces de escribir ensayos, generar resúmenes, traducir idiomas, responder preguntas complejas e incluso crear contenido original, todo con una fluidez cercana a la humana.

Estos avances no solo han ampliado lo que las máquinas pueden hacer; también han redefinido lo que entendemos por "inteligencia". Los modelos revolucionarios de la última década han permitido que la inteligencia artificial transite de ser una promesa técnica a convertirse en una fuerza activa de transformación educativa, económica y cultural.

## 4.4 Disponibilidad masiva de datos

Uno de los motores más poderosos detrás del auge contemporáneo de la inteligencia artificial es la abundancia sin precedentes de datos disponibles para el entrenamiento de modelos. Los algoritmos de aprendizaje profundo —y en especial los modelos de lenguaje de gran escala (LLMs)— dependen de enormes volúmenes de información para detectar patrones, aprender estructuras del lenguaje y generalizar a nuevos contextos.



Figura 4.11. Disponibilidad de datos en la nube

Sin datos, incluso los modelos más sofisticados serían meras arquitecturas vacías. Esta sección analiza las condiciones que han posibilitado esta explosión de datos, así como sus implicaciones estratégicas, técnicas y éticas.

## 4.4.1 El crecimiento exponencial de los datos globales

Vivimos en una época donde cada acción digital genera datos: un clic, una búsqueda, una compra, una conversación. Estas son algunas de las principales fuentes:

- Redes sociales: Plataformas como Instagram, TikTok, Facebook o X producen miles de millones de interacciones diarias que se convierten en insumos para análisis de lenguaje, emociones y comportamiento.
- Comercio electrónico: Sitios como Amazon o Alibaba recopilan reseñas, hábitos de compra y patrones de navegación.

- Sensores IoT: Dispositivos inteligentes (termostatos, relojes, cámaras, etc.) generan datos en tiempo real desde el entorno físico.
- Digitalización masiva: Documentos físicos, libros, fotografías y archivos históricos se están digitalizando a gran escala y se integran a corpus de entrenamiento.

Este ecosistema global e interconectado ha dado lugar al llamado universo Big Data, un océano de información en constante crecimiento.

#### 4.4.2 Acceso abierto y conjuntos de datos públicos

Uno de los elementos más democratizadores del desarrollo de la IA ha sido la proliferación de datasets de acceso abierto. Ya no todo el conocimiento está encerrado en servidores privados. Hoy, investigadores, startups y universidades pueden acceder a datos valiosos gracias a:

- Plataformas como Kaggle y Hugging Face, donde se publican datasets curados y etiquetados para tareas específicas.
- Repositorios académicos y gubernamentales, con datos sobre salud, geografía, medio ambiente, educación, y más.
- Google Dataset Search y similares, que indexan millones de bases de datos públicas y privadas.

Este acceso ha permitido que el entrenamiento de modelos avanzados no sea exclusivo de grandes corporaciones.

## 4.4.3 Infraestructura para almacenamiento y procesamiento

La abundancia de datos sería irrelevante sin una infraestructura que permita almacenarlos, procesarlos y analizarlos en tiempo y forma. Hoy contamos con:

- Sistemas distribuidos como Hadoop y Apache Spark, que permiten el procesamiento eficiente de grandes volúmenes de datos a través de múltiples nodos.
- Servicios en la nube (AWS, Azure, Google Cloud), que ofrecen almacenamiento escalable y acceso remoto a capacidades computacionales avanzadas.
- Redes de alta velocidad, que permiten la transferencia en tiempo real de información entre servidores, sensores y plataformas de análisis.

Esta infraestructura es esencial para que la IA pueda operar sobre datos masivos de forma dinámica y continua.

## 4.4.4 Calidad, curación y etiquetado de datos

El volumen de datos no lo es todo. La calidad, precisión y diversidad de los datos son factores críticos para el éxito de cualquier modelo de IA. Hoy existen herramientas que permiten:

- Limpiar y normalizar datos crudos, eliminando errores, inconsistencias o duplicados.
- Etiquetar datos con precisión, mediante plataformas como Labelbox o Amazon Mechanical Turk.
- Detectar y mitigar sesgos, que podrían introducir distorsiones en los resultados de los modelos.

Un dataset mal curado no solo reduce el rendimiento del modelo: puede llevar a conclusiones erróneas, decisiones injustas o fallos críticos en entornos sensibles.

#### 4.4.5 Datos en tiempo real

La inteligencia artificial no solo se nutre de datos históricos. El acceso a información en tiempo real permite a los modelos adaptarse continuamente al entorno. Ejemplos de esto incluyen:

- Sistemas financieros que responden al mercado en segundos.
- Salud digital, donde wearables monitorean signos vitales y generan alertas inmediatas.
- Movilidad urbana, donde plataformas como Google Maps o Waze ajustan rutas en función del tráfico en vivo.

Los datos en tiempo real son claves para aplicaciones que exigen respuestas ágiles y personalizadas.

## 4.4.6 Consideraciones éticas y desafíos

El acceso masivo a datos también trae consigo riesgos importantes y dilemas éticos:

- Privacidad: ¿Cómo proteger la información personal en modelos entrenados con datos públicos o no autorizados?
- Sesgos: ¿Qué pasa cuando los datos reflejan prejuicios sociales, culturales o de género?
- Regulación: Normativas como el GDPR (Europa) o la Ley de Protección de Datos (América Latina) exigen transparencia, consentimiento y uso responsable.

La IA no solo debe ser técnicamente eficiente, sino también éticamente responsable y socialmente justa. La manera en que gestionemos los datos hoy definirá la legitimidad y sostenibilidad de la IA en el futuro.

La disponibilidad masiva de datos ha sido uno de los factores decisivos en la explosión de la inteligencia artificial moderna. Desde su recolección hasta su uso responsable, los datos se han convertido en el nuevo combustible de la innovación digital. Pero con gran poder viene una gran responsabilidad: transformar los datos en conocimiento sin sacrificar los derechos y valores humanos.

## 4.5 Conclusiones

Este capítulo deja claro que la inteligencia artificial ha llegado a un punto de inflexión histórico, no solo por sus avances técnicos, sino por su integración profunda en nuestra cultura y economía.

- Tres fuerzas clave han impulsado su crecimiento explosivo: la potencia de cómputo (liderada por NVIDIA), el desarrollo de modelos como los Transformers y la disponibilidad masiva de datos.
- Ya no hablamos de IA como una promesa futura, sino como una infraestructura activa en salud, educación, industria y vida cotidiana.
- El auge de la IA también refleja un cambio cultural: buscamos experiencias más inteligentes, personalizadas y humanas.
- ★ La gran pregunta no es si usaremos IA, sino cómo la usaremos con responsabilidad y visión estratégica.
- Entender estos impulsores es esencial para no solo adaptarnos, sino liderar un uso ético, transformador y humano de la inteligencia artificial.

## 🚀 ¿Qué sigue en nuestro viaje?

En el próximo capítulo nos adentraremos en el corazón de la IA generativa y los modelos de lenguaje.

#### Exploraremos:

- ¿Qué es exactamente la IA generativa y cómo funciona?
- ¿Cómo se entrenan los modelos como ChatGPT?
- ¿Qué puede —y qué no debería— hacer un modelo de lenguaje?

Prepárate para descubrir los fundamentos de la IA que está escribiendo el futuro... palabra por palabra.

# **X HERRAMIENTAS IA: Tu nuevo kit de superpoderes**

Nerramienta recomendada: Notion Al

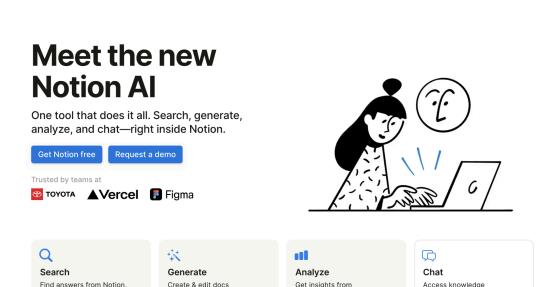


Figura 4.12 Notion AI https://www.notion.com/product/ai

PDFs & images

from GPT-4 & Claude

## Q ¿Qué es?

Slack, Google Drive & more

in your own style

**Notion AI** es una herramienta de inteligencia artificial integrada dentro de la plataforma Notion. Permite redactar, resumir, traducir, corregir y organizar ideas de forma automática. Funciona como un asistente de escritura inteligente que potencia la productividad de docentes, estudiantes y equipos educativos en la planificación, gestión y documentación de actividades escolares o académicas.

## 💡 ¿Para qué sirve en educación?

Ideal para docentes que desean organizar sus clases, redactar materiales, preparar actas, tomar notas inteligentes o resumir textos extensos. Notion Al mejora la claridad, ahorra tiempo y permite enfocar la energía en lo pedagógico en lugar de lo administrativo. También es útil para estudiantes que quieren estructurar apuntes o proyectos.

## Ventajas clave

- Redacción automática de textos, ideas, instrucciones y resúmenes.
- Traducción, corrección y reformulación de textos en segundos.
- Integración perfecta con plantillas para clases, agendas o portafolios.
- Ideal para planificación curricular, diarios de clase o desarrollo profesional.
- Disponible en español y otros idiomas.
- Uso colaborativo en tiempo real con equipos docentes o grupos de estudiantes.

Accede aquí: https://www.notion.so/product/ai

## 💰 Modelo de precios:

- Plan gratuito: Incluye funcionalidades básicas de Notion sin IA.
- Notion AI: Suscripción adicional de pago (~\$8 USD/mes) que puede añadirse al plan gratuito o de equipos.
- **Planes educativos**: Acceso gratuito o con descuento para estudiantes y docentes, previa verificación institucional.

# One tool for your whole company. Free for teams to try.

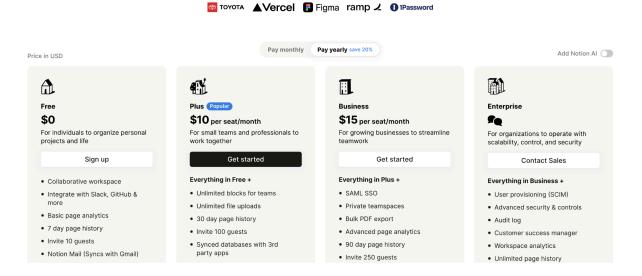


Figura 4.13 Notion Al: Pricing Plans



#### Aprende haciendo con ChatGPT

#### **OBJETIVO**:

Aplicar los conceptos del capítulo para analizar las fuerzas que han impulsado la expansión acelerada de la IA, comprender su contexto tecnológico y geopolítico, y proyectar sus aplicaciones en la vida real.

## **MERRAMIENTAS A UTILIZAR:**

- ChatGPT (versión gratuita)
- Acceso a internet (opcional para validación o exploración de datos reales)

## NO ACTIVIDAD PRÁCTICA N.º 1

#### Explorando los 3 motores del crecimiento de la IA

#### Actividad:

Dialoga con ChatGPT para explicar, con ejemplos actuales, cómo cada uno de los tres factores clave (procesamiento, modelos y datos) ha acelerado el desarrollo de la inteligencia artificial.

#### Prompt sugerido para ChatGPT:

"Explícame con ejemplos actuales cómo el aumento en la capacidad de procesamiento, los nuevos modelos como los Transformers, y la disponibilidad masiva de datos han impulsado el crecimiento de la IA generativa."

## Q Propósito:

Comprender el "por qué ahora" del auge de la IA y visualizar cómo estos factores se combinan para potenciarla.

## NOTIVIDAD PRÁCTICA N.º 2

#### Caso NVIDIA: del videojuego a la revolución tecnológica

## \* Actividad:

Investiga con ayuda de ChatGPT la transformación de NVIDIA en un actor clave de la inteligencia artificial y su impacto en la economía y la geopolítica.

#### Prompt sugerido para ChatGPT:

"Resume cómo NVIDIA pasó de ser una empresa de tarjetas gráficas para videojuegos a convertirse en un actor clave del desarrollo de la inteligencia artificial. Incluye aspectos tecnológicos, financieros y geopolíticos."

## **§** Sugerencia:

Usa los datos obtenidos para crear una presentación de 3 diapositivas o una infografía con Canva o PowerPoint.

## **ACTIVIDAD PRÁCTICA N.º 3**

#### Analiza la batalla tecnológica entre China y EE.UU.

#### Actividad:

Explora las tensiones entre China y Estados Unidos en torno a los chips de IA. Reflexiona sobre cómo esto podría afectar a tu país o región.

#### **Prompt sugerido para ChatGPT:**

"¿Qué implicaciones tiene la guerra tecnológica entre China y EE.UU. por el control de los chips de IA? ¿Cómo podría afectar esto a países de América Latina, África o el Sudeste Asiático?"

#### **Utilidad:**

Conectar los desarrollos tecnológicos con el contexto global y entender el impacto político de la IA.

## **NOTIVIDAD PRÁCTICA N.º 4**

#### Del dato al poder: ética y responsabilidad en la IA

#### Actividad:

Reflexiona con ChatGPT sobre los riesgos del uso de datos masivos y cómo se pueden aplicar principios éticos en su gestión.

#### **Prompt sugerido para ChatGPT:**

"Quiero diseñar un código de buenas prácticas éticas para el uso de datos en proyectos de inteligencia artificial. ¿Qué principios clave debería incluir para proteger la privacidad, reducir sesgos y promover la equidad?"

## Propósito:

Promover una visión crítica y responsable del uso de la información en contextos tecnológicos.

## NOTIVIDAD PRÁCTICA N.º 5

#### Resumen activo del capítulo 4

## Actividad:

Redacta un resumen integrador del capítulo que abarque los factores impulsores, el rol de actores clave como NVIDIA y los desafíos éticos del uso masivo de datos.

#### **Prompt sugerido para ChatGPT:**

"Ayúdame a escribir un resumen de 200 palabras del capítulo 4 del libro, incluyendo los tres impulsores clave del crecimiento de la IA, el caso de NVIDIA, el conflicto China-EE.UU., y los desafíos éticos relacionados con los datos."

## ★ Propósito:

Consolidar el aprendizaje desde una perspectiva crítica y conectarlo con el entorno actual.

## **GUÍA DE ESTUDIO**

Este estudio de repaso cubre los puntos clave del Capítulo IV, centrándose en los factores que explican el momento actual de rápido avance y expansión de la inteligencia artificial.

#### **Temas Clave:**

- Punto de Inflexión de la IA: Por qué la IA ha alcanzado un momento histórico de desarrollo y los factores convergentes que lo han impulsado.
- Capacidad de Procesamiento: El papel crucial de los avances en hardware (especialmente las GPUs y NVIDIA) y cómo han permitido entrenar modelos más complejos.
- Modelos Revolucionarios: La evolución de las arquitecturas de aprendizaje profundo, destacando la importancia de los Transformers y su impacto en el Procesamiento del Lenguaje Natural (PLN).
- Disponibilidad de Datos: El rol fundamental de la abundancia de datos masivos (Big Data) y en tiempo real, así como los desafíos éticos asociados.
- Geopolítica de la IA: La competencia global por el control de la infraestructura de IA, particularmente entre China y Estados Unidos.
- Desafíos Éticos: Las consideraciones morales y de responsabilidad en el uso de datos y el desarrollo de la IA.

#### **Preguntas de Repaso Cortas:**

- ¿Cuáles son los tres factores convergentes que explican por qué la IA ha alcanzado un punto de inflexión histórico ahora?
- ¿Cómo ha influido el desarrollo de las GPUs, y en particular NVIDIA, en el entrenamiento de modelos de IA complejos?
- ¿Qué es la arquitectura Transformer y por qué se considera que revolucionó el Procesamiento del Lenguaje Natural?
- Explica la importancia de la disponibilidad masiva de datos para el entrenamiento de los sistemas de IA modernos.
- ¿Qué papel juega el acceso a datasets abiertos en la democratización del desarrollo de modelos de IA?
- Menciona un ejemplo de cómo los datos en tiempo real potencian las aplicaciones de IA.
- ¿Qué son los TOPS (Tera Operations Per Second) y por qué la evolución del rendimiento de los chips NVIDIA (como el paso de FP32 a Int8) es relevante para la inferencia de IA?
- Describe brevemente la estrategia de Estados Unidos para limitar el acceso de China a chips avanzados de IA.
- ¿Cómo ha respondido China a las restricciones de EE.UU. en cuanto a chips de IA?

• ¿Cuáles son algunos de los desafíos éticos clave relacionados con el uso de grandes volúmenes de datos en IA?

#### Preguntas en Formato de Ensayo:

- Analiza la convergencia de los tres factores clave (procesamiento, modelos, datos) como un fenómeno histórico único que explica el punto de inflexión actual de la IA. ¿Cómo la interacción entre estos factores ha creado un ecosistema propicio para la IA generativa?
- Evalúa el impacto estratégico y geopolítico del dominio de NVIDIA en el mercado de hardware para IA. ¿Cómo la competencia entre potencias globales por el control de esta infraestructura, ejemplificada por la disputa China-EE.UU. sobre los chips, moldea el futuro de la tecnología global?
- Explica cómo la arquitectura Transformer no solo mejoró el rendimiento en tareas de PLN, sino que también sentó las bases para la era de la IA generativa. Compara sus capacidades con modelos anteriores de procesamiento de lenguaje natural y discute las nuevas posibilidades que abrió.
- Discute la dualidad de la disponibilidad masiva de datos: por un lado, es un motor esencial para el avance de la IA, y por otro, plantea significativos desafíos éticos y de regulación. ¿Cómo se puede equilibrar la necesidad de datos para el entrenamiento con la protección de la privacidad y la mitigación de sesgos?
- Considerando los impulsores técnicos, económicos y geopolíticos del crecimiento de la IA, ¿cuáles son las implicaciones más importantes para un sector específico (por ejemplo, la educación o la salud)? ¿Cómo deberían prepararse las instituciones y los profesionales para aprovechar las oportunidades y mitigar los riesgos de esta transformación?

#### **Glosario de Términos Clave:**

- Inteligencia Artificial (IA): Capacidad de las máquinas para realizar tareas que normalmente requieren inteligencia humana, como aprendizaje, percepción y toma de decisiones.
- IA Generativa: Tipo de IA capaz de crear contenido nuevo (texto, imágenes, audio, etc.) a partir de patrones aprendidos de datos existentes.
- Punto de Inflexión Histórico: Momento significativo en el que una tendencia o tecnología experimenta un cambio drástico y acelerado, con profundas implicaciones futuras.
- Capacidad de Procesamiento: La potencia computacional y la velocidad con la que los sistemas pueden ejecutar cálculos complejos, esencial para el entrenamiento de modelos de IA.
- GPU (Unidad de Procesamiento Gráfico): Tipo de procesador especializado diseñado para manejar cálculos paralelos masivos, fundamental para el deep learning y la inferencia de IA.
- Deep Learning (Aprendizaje Profundo): Subcampo del machine learning que utiliza redes neuronales artificiales con múltiples capas (profundas) para aprender representaciones complejas de los datos.

- Modelos Revolucionarios: Arquitecturas de algoritmos de aprendizaje automático que representan un avance significativo en sus respectivas áreas, permitiendo nuevas capacidades.
- Transformer: Arquitectura de red neuronal introducida en 2017, conocida por su mecanismo de atención y su capacidad para procesar secuencias de datos de forma no secuencial, revolucionando el PLN.
- Mecanismo de Atención: Componente de un modelo (como el Transformer) que le permite ponderar la importancia de diferentes partes de los datos de entrada al generar una salida.
- Procesamiento del Lenguaje Natural (PLN / NLP): Campo de la IA que se ocupa de la interacción entre computadoras y lenguaje humano, permitiendo a las máquinas comprender, interpretar y generar texto y voz.
- Modelos de Lenguaje de Gran Escala (LLMs): Modelos de PLN muy grandes, entrenados en vastas cantidades de texto, capaces de comprender contexto, generar texto coherente y realizar diversas tareas lingüísticas.
- Disponibilidad Masiva de Datos (Big Data): La existencia de enormes y crecientes volúmenes de datos, a menudo variados y producidos a alta velocidad, que sirven como base para el entrenamiento de modelos de IA.
- Datasets Abiertos: Conjuntos de datos que están disponibles públicamente para su uso, distribución y reutilización por cualquier persona, a menudo bajo licencias permisivas.
- Curación de Datos: El proceso de organizar, limpiar y preparar datos para su uso, asegurando su calidad, precisión y relevancia para una tarea específica de IA.
- Etiquetado de Datos: El proceso de asignar etiquetas o metadatos a los datos (como identificar objetos en imágenes o palabras clave en texto) para ayudar a los modelos de aprendizaje supervisado a reconocer patrones.
- Datos en Tiempo Real: Información que se procesa y está disponible inmediatamente después de su recopilación, permitiendo a los sistemas de IA reaccionar de forma dinámica.
- TOPS (Tera Operations Per Second): Unidad de medida del rendimiento computacional, indicando billones de operaciones por segundo, comúnmente utilizada para chips de IA.
- Int8: Formato numérico de 8 bits, más eficiente en espacio y energía que formatos de mayor precisión (como FP32 o FP16), ideal para la fase de inferencia de los modelos de IA.
- Inferencia de IA: El proceso de ejecutar un modelo de IA ya entrenado sobre nuevos datos para hacer predicciones o tomar decisiones.
- Geopolítica del Procesamiento: La interacción entre la tecnología de chips de alto rendimiento (procesamiento para IA) y las relaciones internacionales, rivalidades entre potencias y control estratégico de la infraestructura tecnológica.
- Autosuficiencia Tecnológica: La capacidad de un país o entidad para desarrollar y producir su propia tecnología crítica sin depender de proveedores externos, a menudo impulsada por razones de seguridad nacional o soberanía.

- CPU (Unidad Central de Procesamiento): El procesador principal de una computadora, diseñado para realizar una amplia gama de tareas de procesamiento secuencial.
- Retropropagación (Backpropagation): Algoritmo fundamental utilizado para entrenar redes neuronales ajustando los pesos de las conexiones basado en la diferencia entre la salida esperada y la salida real.
- Word Embeddings: Técnica de PLN que representa palabras como vectores numéricos en un espacio multidimensional, donde la distancia entre vectores refleja similitud semántica.
- Regulación de Datos: Leyes y normativas (como GDPR) que rigen la recolección, almacenamiento, procesamiento y uso de datos, especialmente datos personales, para proteger la privacidad y otros derechos.
- Sesgos en Datos: Patrones en los datos de entrenamiento que reflejan o perpetúan prejuicios sociales, culturales o de otro tipo, pudiendo llevar a que los modelos de IA produzcan resultados discriminatorios o injustos.
- Notion AI: Herramienta de inteligencia artificial integrada en la plataforma Notion que ayuda a redactar, resumir, traducir y organizar información.

## REFERENCIAS

Wikipedia. (2023, octubre 17). *Jen-Hsun Huang*. En *Wikipedia, la enciclopedia libre*. Recuperado el 2 de Enero del 2025,

de https://es.wikipedia.org/wiki/Jen-Hsun\_Huang

NVIDIA. (2023, August 29). Why GPUs are great for AI. NVIDIA Blog. https://blogs.nvidia.com/blog/why-gpus-are-great-for-ai/

NVIDIA Corporation. (n.d.). *Stock quote and chart*. NVIDIA Investor Relations. Recuperado el 3 de Enero de 2025,

de <a href="https://investor.nvidia.com/stock-info/stock-quote-and-chart/default.aspx">https://investor.nvidia.com/stock-info/stock-quote-and-chart/default.aspx</a>

Sánchez, E. (2024, abril 12). El Gobierno de EEUU acaba de dar un golpe muy duro a NVIDIA: no podrá vender su GPU H20 para IA en China. Xataka.

https://www.xataka.com/empresas-y-economia/gobierno-eeuu-acaba-dar-golpe-muy-duro-a-nvidia-no-podra-vender-su-gpu-h20-para-ia-china

Associated Press. (2024, abril 13). Nvidia halts AI chip exports to China after U.S. ban, projects \$5.5 billion loss. AP News.

https://apnews.com/article/nvidia-china-h20-chip-ban-2024

South China Morning Post. (2024, abril 14). US Nvidia H20 restrictions could be a shot in the arm for domestic AI chips like Huawei's Ascend. SCMP.

https://www.scmp.com/tech/tech-war/article/3306793/us-nvidia-h20-restrictions-could-be-shot-arm-domestic-ai-chips-huaweis

Cadena SER. (2024, abril 16). El Nasdaq cae un 3% arrastrado por Nvidia en medio de los temores por la guerra comercial entre EE.UU. y China.

https://cadenaser.com/nacional/2025/04/16/el-nasdaq-cae-un-3-arrastrado-por-nvid ia-en-medio-de-los-temores-por-la-guerra-comercial-entre-eeuu-y-china-cadena-ser

Vincent, J. (2023, October 17). New US restrictions block AI chip sales to China, targeting Nvidia. The

Verge. <a href="https://www.theverge.com/2023/10/17/23921131/us-china-restrictions-ai-chip-sales-nvidia">https://www.theverge.com/2023/10/17/23921131/us-china-restrictions-ai-chip-sales-nvidia</a>

NVIDIA. (2015, March 18). *Mythbusters demo GPU versus CPU* [Video]. YouTube. <a href="https://www.youtube.com/watch?v=-P28LKWTzrl&t=14s">https://www.youtube.com/watch?v=-P28LKWTzrl&t=14s</a>

NVIDIA. (2023, August 29). Why GPUs are great for AI. NVIDIA Blog. https://blogs.nvidia.com/blog/why-gpus-are-great-for-ai/

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. In Advances in Neural Information Processing Systems (Vol.

**30)**. <a href="https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf">https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf</a>

Wikipedia. (2023, octubre 17). *Geoffrey Hinton*. En *Wikipedia, la enciclopedia libre*. Recuperado el 3 de Enero del 2025, de <a href="https://es.wikipedia.org/wiki/Geoffrey\_Hinton">https://es.wikipedia.org/wiki/Geoffrey\_Hinton</a>

Louis, A. (n.d.). A brief history of natural language processing (part 2). Medium. Recuperado el 3 de Enero del 2025,

de <a href="https://medium.com/@antoine.louis/a-brief-history-of-natural-language-processing-part-2-f5e575e8e37">https://medium.com/@antoine.louis/a-brief-history-of-natural-language-processing-part-2-f5e575e8e37</a>