

## LA GPU COMO MOTOR DE CÓMPUTO MODERNO:

### DE LOS GRÁFICOS A LA INTELIGENCIA ARTIFICIAL



CIFP Francesc Borja Moll

Administración de sistemas informáticos en redes I

Fonaments de maquinari

Antonio serna

2026

**Contenido**

1.	Introducción .....	6
1.1	La GPU ya no es “solo para jugar” .....	<b>¡Error! Marcador no definido.</b>
1.2	De los píxeles al cálculo masivo: el nuevo rol de la tarjeta gráfica en la sociedad digital.....	<b>¡Error! Marcador no definido.</b>
2.	¿Qué es una tarjeta gráfica y cuál es su función real? .....	7
2.1	Definición técnica y distinción entre GPU y tarjeta gráfica .....	7
3.	CPU vs GPU: Dos filosofías de procesamiento .....	8
3.1	Procesamiento secuencial (latencia) frente a paralelismo masivo (rendimiento).....	8
3.2	Arquitecturas SIMD y SIMT .....	9
	SIMD (Single Instruction, Multiple Data): paralelismo rígido y altamente eficiente.....	9
	SIMT (Single Instruction, Multiple Threads): paralelismo flexible orientado a la GPU moderna .....	10
3.3	Por qué una GPU supera a una CPU en tareas específicas: la analogía del tren frente a los coches deportivos .....	11
	La CPU como coche deportivo: velocidad, control y baja latencia .....	11
	La GPU como tren de mercancías: potencia bruta y paralelismo masivo.....	12
4.	Anatomía de una GPU moderna (Arquitectura de núcleos).....	13
4.1	El chip gráfico (GPU) como núcleo del sistema.....	14
4.2	Núcleos CUDA (NVIDIA) y Stream Processors (AMD): los obreros de la rasterización.....	15
4.3	Núcleos Tensor: aceleración por hardware para inteligencia artificial y <i>Deep Learning</i> .....	16
	Precisión y rendimiento en IA (FP16, BF16, INT8).....	17

4.4	Núcleos Ray Tracing: unidades de intersección para la simulación física de la luz	17
4.5	Reloj, frecuencia y overclocking: maximizando el rendimiento del silicio ...	17
5.	El pipeline gráfico: del dato matemático al píxel.....	20
5.1	Datos de entrada: geometría y atributos.....	20
5.2	Procesamiento de vértices (Vertex Shader).....	21
5.3	Ensamblado de primitivas.....	21
5.4	Rasterización: de geometría a fragmentos.....	21
5.5	Procesamiento de fragmentos (Fragment Shader).....	22
5.6	Tests finales y salida a pantalla.....	22
6.	VRAM: la memoria que alimenta a la GPU.....	23
6.1	¿Qué es la VRAM y por qué es diferente de la RAM del sistema? .....	24
6.2	Tipos de memoria gráfica y ancho de banda.....	25
6.3	VRAM y resolución: por qué cada píxel cuenta.....	26
6.4	Gestión de memoria y limitaciones prácticas.....	27
6.5	El mito de “más VRAM siempre es mejor” .....	28
6.6	Cuando la VRAM no es suficiente: stuttering y pérdida de fluidez.....	29
7.	Tarjetas gráficas integradas y dedicadas.....	30
7.1	GPU integrada (iGPU): eficiencia y limitaciones.....	30
7.2	GPU dedicada (dGPU): potencia y especialización.....	31
7.3	Comparación directa: iGPU vs dGPU.....	32
7.4	Consumo energético y eficiencia.....	34
7.5	Escenarios de uso: ¿cuándo es suficiente una iGPU?.....	35
8.	Conexión, alimentación y refrigeración .....	36

8.1	Interfaz PCI Express: la autopista de datos.....	37
8.2	Alimentación eléctrica: potencia y estabilidad.....	38
8.3	Refrigeración: control térmico y rendimiento sostenido .....	39
8.4	Impacto conjunto en el sistema.....	41
9.	Mercado actual y modelos de tarjetas gráficas.....	42
9.1	El triunvirato del silicio: NVIDIA, AMD e Intel.....	42
9.2	Ensambladores (AIB): más allá del chip gráfico.....	43
9.3	Segmentación por gamas: básica, media y alta.....	43
9.4	Aplicaciones apropiadas para cada gama.....	44
9.5	Tendencias del mercado y ciclos de renovación .....	45
9.6	Precio, rendimiento y sentido común.....	46
10.	Monitores y tecnologías de visualización: donde la GPU cobra sentido.....	47
10.1	Tecnologías de panel: cómo se crea la imagen.....	47
	TN (Twisted Nematic): Velocidad y Transición .....	48
	IPS (In-Plane Switching): Fidelidad y Consistencia.....	48
	VA (Vertical Alignment): El Dominio del Contraste.....	48
	OLED (Organic Light Emitting Diode): La Excelencia Autoemisiva .....	49
10.2	Resolución: cuántos píxeles puede mover tu GPU.....	50
10.3	Tasa de refresco y tiempo de respuesta: fluidez percibida.....	52
10.4	Brillo, contraste y calidad de imagen.....	53
10.5	Sincronización adaptativa: GPU y monitor trabajando juntos .....	55
	El conflicto entre FPS y Hercios (Hz).....	55
	Ecosistemas de Sincronización: G-Sync, FreeSync y VRR.....	55
10.6	La sinergia GPU–Monitor: elegir bien el conjunto .....	57

11.	Conclusiones y perspectivas de futuro.....	58
12.	ANEXO: HORIZONTES EXPANDIDOS (La Era GPGPU).....	59
12.1	Principios de los 2000: La "Alquimia Digital" y los Inicios.....	59
12.2	2006 – 2020: De Folding@home a la Lucha Global contra Pandemias.....	59
12.3	2010: El Cluster Condor y la Supercomputación con Consolas.....	60
12.4	2010 – Actualidad: La Revolución del Criptoanálisis.....	60
12.5	2012: El "Momento AlexNet" y el Renacimiento de la IA.....	61
12.6	Presente y Futuro: Del Renderizado al Diagnóstico Médico.....	61
13.	Bibliografía y fuentes consultadas.....	62

## 1. Introducción

La evolución de la informática moderna ha estado tradicionalmente dominada por la Unidad Central de Procesamiento (CPU) como el cerebro encargado de la lógica secuencial y la toma de decisiones. Sin embargo, en las últimas dos décadas, el crecimiento exponencial de la demanda de cálculo ha evidenciado las limitaciones del procesamiento en serie, impulsando la aparición de un nuevo protagonista indiscutible: la Unidad de Procesamiento Gráfico (GPU).

Inicialmente concebidas como aceleradores de función fija dedicados exclusivamente a pintar píxeles en pantalla, las GPUs han protagonizado una metamorfosis arquitectónica radical. Han evolucionado hasta convertirse en dispositivos de cómputo de propósito general (GPGPU) altamente paralelos, capaces de gestionar miles de hilos de ejecución simultáneos. Esta capacidad de procesamiento masivo (throughput) las ha convertido en la herramienta ideal para resolver problemas que la CPU no puede abordar eficientemente, desde la simulación física de la luz mediante Ray Tracing hasta el entrenamiento de redes neuronales profundas que sustentan la inteligencia artificial moderna.

Esta transformación ha redefinido el diseño de los sistemas informáticos actuales. La tarjeta gráfica ya no es un mero periférico de salida, sino un ecosistema complejo que integra memoria de alto ancho de banda (VRAM), sistemas de alimentación avanzados y núcleos especializados (Tensor y RT). Además, su función no termina en el cálculo; el proyecto también aborda cómo esta potencia bruta depende de una simbiosis crítica con el dispositivo de visualización: el monitor. Sin tecnologías de panel y sincronización adecuadas, el rendimiento generado por el silicio se pierde antes de llegar a los ojos del usuario.

El presente trabajo desglosa la anatomía y el funcionamiento de estos dispositivos, analizando desde la microarquitectura de los núcleos y la gestión de la memoria VRAM, hasta las implicaciones prácticas de elegir entre soluciones integradas o dedicadas. El objetivo final es comprender no solo qué es una tarjeta gráfica, sino por qué se ha convertido en el motor de cálculo esencial de la sociedad digital del siglo XXI.

## 2. ¿Qué es una tarjeta gráfica y cuál es su función real?

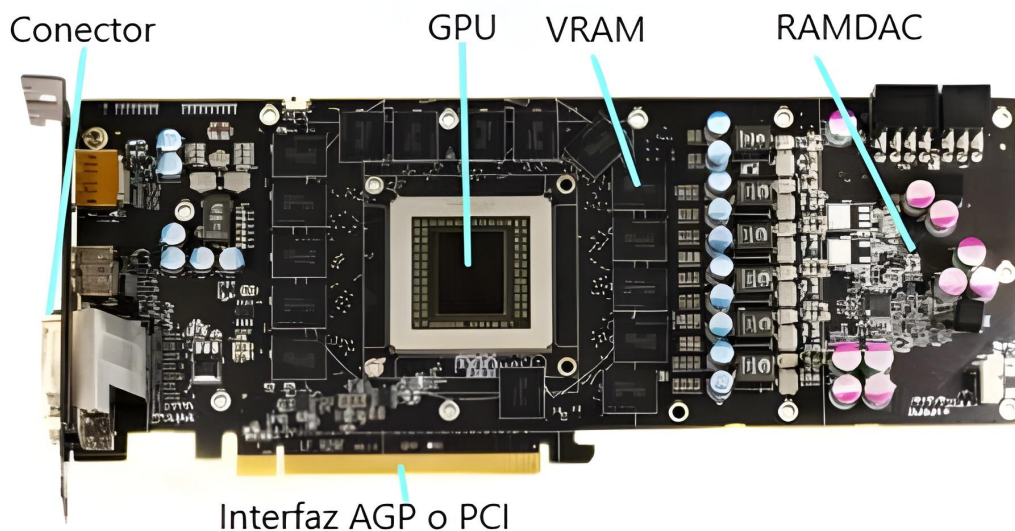
Una tarjeta gráfica es un dispositivo de hardware especializado cuya función principal es procesar información visual y generar la señal de vídeo que finalmente se muestra en un monitor. Actúa como un intermediario entre el sistema informático y la pantalla, transformando datos matemáticos y geométricos en píxeles visibles.

Aunque en el lenguaje cotidiano ambos términos suelen usarse como sinónimos, es fundamental diferenciar entre GPU y tarjeta gráfica, ya que no representan exactamente lo mismo desde el punto de vista técnico.

### 2.1 Definición técnica y distinción entre GPU y tarjeta gráfica

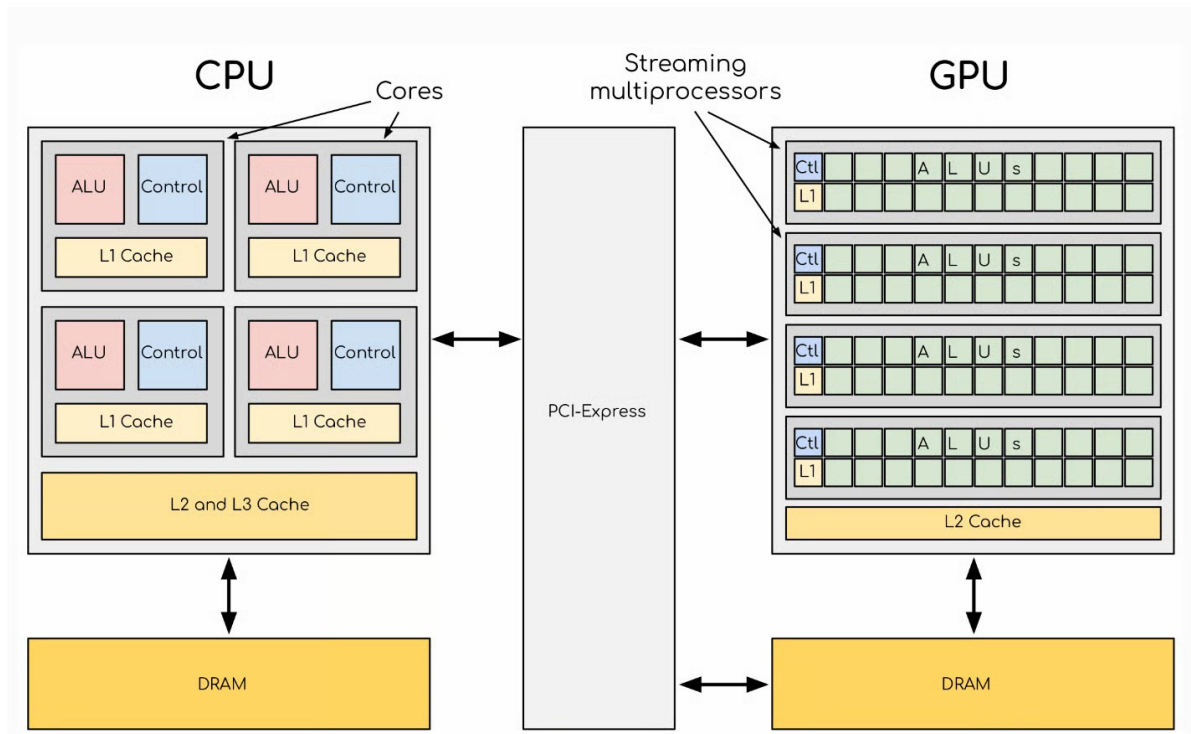
La GPU (Graphics Processing Unit) es el procesador gráfico: un circuito integrado altamente especializado en realizar cálculos paralelos relacionados con gráficos, geometría, iluminación y, en la actualidad, también con tareas de propósito general como inteligencia artificial o computación científica.

La tarjeta gráfica, en cambio, es el dispositivo físico completo que integra la GPU junto con otros componentes esenciales, entre los que se incluyen VRAM, sistema de alimentación eléctrica, sistema de refrigeración, controladores de vídeo y salidas de monitor y circuitería de comunicación con la placa base:



### 3. CPU vs GPU: Dos filosofías de procesamiento

Aunque tanto la CPU como la GPU son procesadores, su diseño interno responde a filosofías radicalmente distintas. Comprender estas diferencias es fundamental para entender por qué ciertas tareas se ejecutan mejor en una GPU y otras siguen dependiendo de la CPU. No se trata de determinar cuál es “más potente”, sino de analizar para qué tipo de problemas está optimizado cada procesador y cómo sus arquitecturas reflejan esas prioridades.



#### 3.1 Procesamiento secuencial (latencia) frente a paralelismo masivo (rendimiento)

La CPU está diseñada para minimizar la latencia, es decir, el tiempo que tarda en ejecutar una tarea individual desde que se inicia hasta que finaliza. Para lograrlo, cuenta con pocos núcleos muy complejos, capaces de ejecutar instrucciones de forma extremadamente rápida y con un alto grado de control lógico. Esta arquitectura resulta ideal para tareas secuenciales, toma de decisiones, lógica de programas y procesos que dependen de múltiples condiciones encadenadas.

En cambio, la GPU prioriza el rendimiento total (throughput) frente a la latencia. En lugar de optimizar la ejecución rápida de una sola tarea, la GPU está diseñada para ejecutar miles de operaciones similares en paralelo, aunque cada una de ellas pueda tardar ligeramente más que en una CPU. Esta aproximación resulta especialmente eficiente en problemas donde el mismo cálculo debe aplicarse a grandes volúmenes de datos, como ocurre en el renderizado gráfico, el procesamiento de imágenes o el entrenamiento de modelos de inteligencia artificial.

Esta diferencia explica por qué una CPU puede superar a una GPU en tareas como la ejecución de un sistema operativo o un videojuego mal paralelizado, mientras que la GPU domina en cálculos masivos y repetitivos.

### **3.2 Arquitecturas SIMD y SIMT**

Para que una GPU pueda ejecutar miles de operaciones de forma simultánea, no basta con disponer de muchos núcleos de cálculo; es necesario un modelo de ejecución que permita coordinar ese paralelismo de manera eficiente. En este contexto surgen las arquitecturas SIMD y SIMT, dos enfoques diseñados para maximizar el rendimiento en tareas altamente paralelizables. Aunque comparten una filosofía común, presentan diferencias clave que explican la evolución de las GPU modernas y su capacidad para afrontar cargas de trabajo cada vez más complejas.

#### ***SIMD (Single Instruction, Multiple Data): paralelismo rígido y altamente eficiente***

El modelo SIMD se basa en la ejecución de una única instrucción sobre múltiples conjuntos de datos de forma simultánea. En este esquema, un controlador central emite una instrucción que es aplicada al mismo tiempo por varias unidades de cálculo, cada una operando sobre datos distintos. Este enfoque reduce significativamente la sobrecarga de control y permite un uso extremadamente eficiente del hardware, ya que todos los recursos trabajan de manera sincronizada.

La principal ventaja de SIMD es su alto rendimiento energético y computacional cuando las operaciones son homogéneas y no presentan bifurcaciones lógicas. Por esta razón,

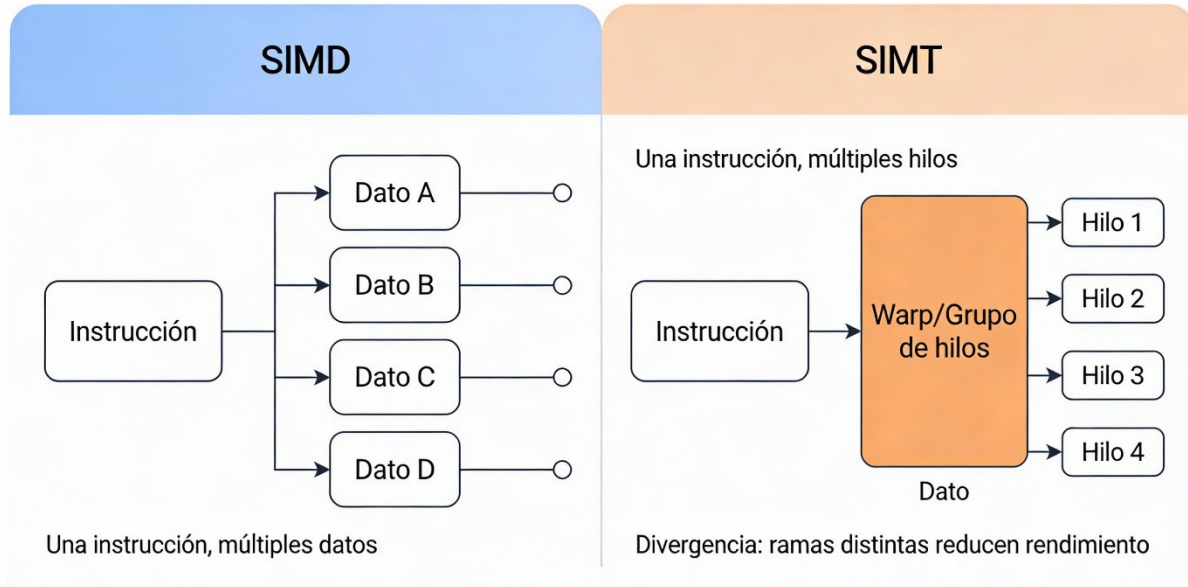
ha sido ampliamente utilizado en tareas como el procesamiento de vectores, operaciones matemáticas repetitivas y ciertos algoritmos gráficos clásicos. Sin embargo, esta misma rigidez se convierte en su mayor limitación: si una parte del cálculo requiere una operación diferente, el modelo SIMD pierde eficiencia o directamente no puede adaptarse sin introducir complejidad adicional.

***SIMT (Single Instruction, Multiple Threads): paralelismo flexible orientado a la GPU moderna***

El modelo SIMT surge como una evolución del enfoque SIMD, adaptado a las necesidades reales de las GPU modernas. En lugar de ejecutar estrictamente una única instrucción sobre múltiples datos, SIMT organiza la ejecución en grupos de hilos que siguen el mismo flujo de instrucciones, pero conservan la capacidad de tomar decisiones de manera independiente cuando es necesario.

Este enfoque permite a la GPU manejar divergencias de control de forma más flexible, algo esencial en tareas complejas como el trazado de rayos, la simulación física o la inferencia de redes neuronales. Aunque las divergencias pueden reducir temporalmente la eficiencia, el modelo SIMT logra un equilibrio entre paralelismo masivo y adaptabilidad, convirtiéndose en el pilar de arquitecturas como CUDA en NVIDIA o los *wavefronts* en GPUs AMD.

La diferencia fundamental entre SIMD y SIMT radica, por tanto, en el grado de flexibilidad: mientras SIMD prioriza la eficiencia absoluta bajo condiciones ideales, SIMT sacrifica parte de esa rigidez para ganar versatilidad, permitiendo que la GPU afronte una gama mucho más amplia de problemas computacionales.



### 3.3 Por qué una GPU supera a una CPU en tareas específicas: la analogía del tren frente a los coches deportivos

Una forma intuitiva de entender esta diferencia es mediante la siguiente analogía:

#### ***La CPU como coche deportivo: velocidad, control y baja latencia***

La CPU puede compararse con un coche deportivo de altas prestaciones: diseñada para reaccionar rápidamente, tomar decisiones complejas y ejecutar tareas secuenciales con la mínima latencia posible. Sus pocos núcleos, altamente sofisticados, están optimizados para manejar flujos de control complejos, saltos condicionales y operaciones que dependen unas de otras. Esta arquitectura la convierte en la pieza central para la ejecución del sistema operativo, la lógica de los programas y las tareas que requieren gran precisión y control.

Gracias a su capacidad para cambiar de contexto rápidamente y a sus avanzados mecanismos de predicción y caché, la CPU sobresale en escenarios donde cada instrucción importa y no puede paralelizarse fácilmente. Sin embargo, esta especialización limita su capacidad para escalar cuando el problema requiere realizar el mismo cálculo miles de veces de forma simultánea.

### ***La GPU como tren de mercancías: potencia bruta y paralelismo masivo***

La GPU, en contraste, se asemeja a un tren de mercancías: no está diseñada para maniobras rápidas o cambios constantes de dirección, pero puede transportar enormes volúmenes de carga de forma eficiente y sostenida. Sus miles de núcleos simples trabajan en paralelo ejecutando operaciones similares sobre grandes conjuntos de datos, lo que la hace ideal para tareas altamente paralelizables.

Este enfoque permite a la GPU alcanzar rendimientos muy superiores a los de una CPU en ámbitos como el renderizado gráfico, el procesamiento de imágenes, la minería de criptomonedas o el entrenamiento de modelos de inteligencia artificial. Aunque cada núcleo individual es menos potente que uno de CPU, la suma de todos ellos trabajando de forma coordinada genera una capacidad de cálculo masiva.

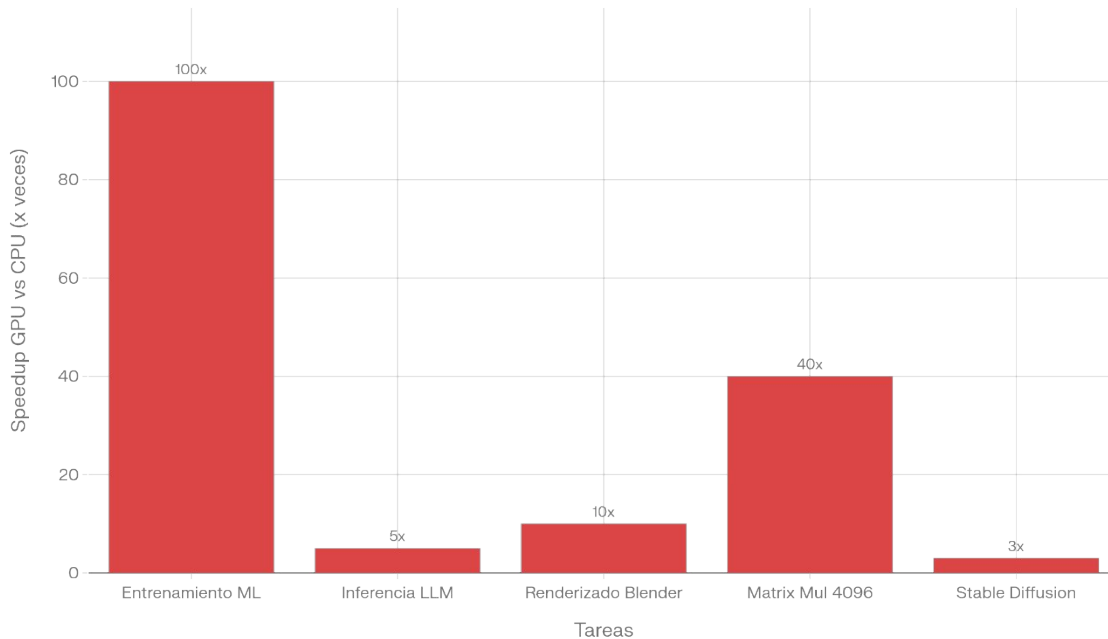
Esta analogía ilustra por qué CPU y GPU no compiten directamente, sino que se complementan. Cada una está diseñada para un tipo de problema distinto, y el verdadero rendimiento del sistema surge cuando ambas trabajan en conjunto, explotando sus fortalezas respectivas.

Para hacer posible este paralelismo masivo, la GPU adopta una arquitectura interna radicalmente distinta, organizada en miles de núcleos especializados que trabajan de forma coordinada.

## Tareas donde GPU supera a CPU: Speedup Relativo

Benchmarks 2025 | GPU acelera procesos paralelos masivos

Powered by  perplexity



### Tareas Específicas

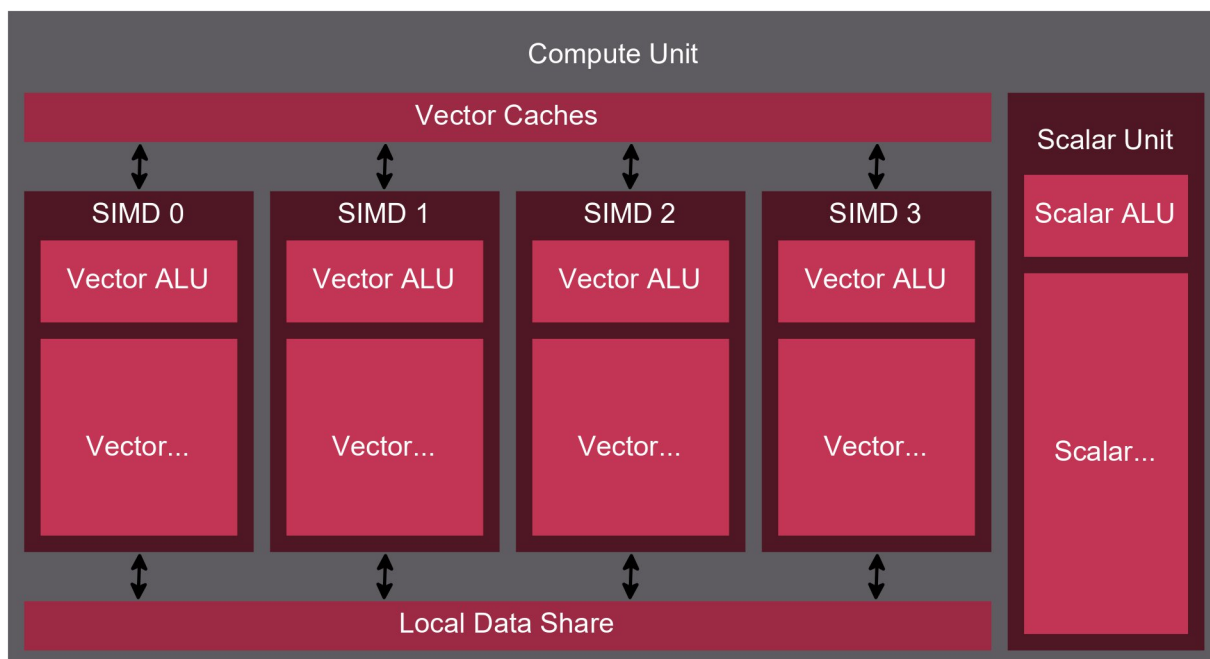
- ✓ Entrenamiento ML: GPU procesa matrices en paralelo, 100x vs CPU (TensorFlow benchmark).
- ✓ Inferencia LLM: 5x más rápida en modelos como Llama (RTX vs Ryzen).
- ✓ Renderizado Blender: 10x en OptiX/CUDA (RTX 4090 vs Ryzen 7950X).
- ✓ Matrix Mul 4096x4096: 40x speedup en cómputo paralelo.
- ✓ Stable Diffusion: 3x en generación de imágenes AI.

## 4. Anatomía de una GPU moderna (Arquitectura de núcleos)

El extraordinario rendimiento de una GPU moderna no es fruto únicamente de su frecuencia de funcionamiento, sino de una arquitectura interna diseñada específicamente para el paralelismo masivo. A diferencia de la CPU, que prioriza núcleos complejos y versátiles, la GPU está formada por miles de unidades de cálculo simples que trabajan de forma coordinada sobre grandes conjuntos de datos.

Esta organización permite a la GPU ejecutar simultáneamente miles de hilos (*threads*), lo que resulta ideal para tareas como el renderizado gráfico, el cálculo matricial o el aprendizaje automático. Para comprender cómo se alcanza este nivel de rendimiento, es necesario analizar los distintos tipos de núcleos y unidades especializadas que componen una GPU actual.

Tras comprender las diferencias conceptuales entre CPU y GPU, es necesario descender al interior de la GPU para entender cómo este paralelismo se materializa físicamente



*\* Como se observa en el siguiente diagrama de una Unidad de Cómputo (Compute Unit), la arquitectura se divide en procesadores vectoriales para el grueso de los datos y unidades escalares para la gestión del flujo de trabajo.*

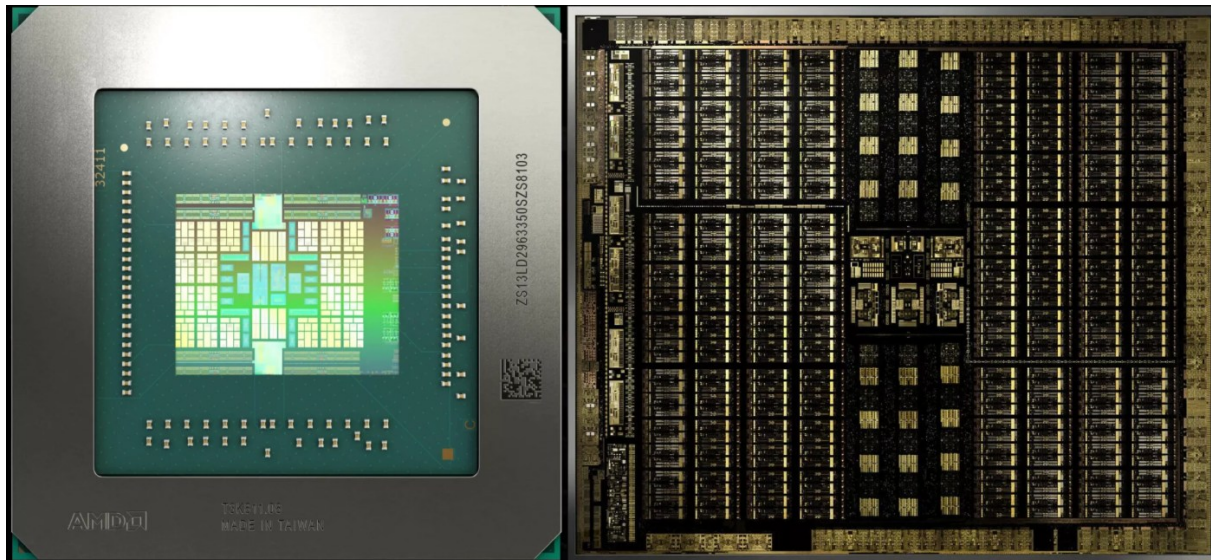
#### 4.1 El chip gráfico (GPU) como núcleo del sistema

El chip gráfico, o GPU propiamente dicha, es un circuito integrado de alta complejidad que concentra la mayor parte de la capacidad de cálculo de la tarjeta gráfica. Está fabricado mediante procesos de litografía avanzados (actualmente de pocos nanómetros), lo que permite integrar miles de millones de transistores en una superficie reducida.

Este chip no actúa de forma aislada: se comunica constantemente con la memoria gráfica, la placa base y el sistema operativo. Para ello, la GPU incluye:

- Unidades de cálculo.
- Cachés internas.
- Controladores de memoria.
- Motores de programación de hilos (*schedulers*).

La GPU puede considerarse un sistema de cómputo completo, especializado en operaciones paralelas, capaz de ejecutar tanto tareas gráficas como algoritmos de propósito general mediante tecnologías como CUDA, OpenCL o Vulkan Compute.



\* A la izquierda, el encapsulado externo de una GPU AMD con su die central protegido. A la derecha, una microfotografía del die (matriz), donde se aprecian los miles de millones de transistores organizados en estructuras repetitivas que forman los núcleos de procesamiento paralelo.

## 4.2 Núcleos CUDA (NVIDIA) y Stream Processors (AMD): los obreros de la rasterización

Los núcleos CUDA (en NVIDIA) y los Stream Processors (en AMD) son las unidades básicas de cálculo de la GPU. Aunque reciben nombres distintos según el fabricante, su función

es esencialmente la misma: ejecutar operaciones matemáticas simples de forma masiva y paralela.

Estos núcleos no funcionan de manera independiente como los de una CPU. Se agrupan en bloques (SM en NVIDIA, Compute Units en AMD) y ejecutan instrucciones bajo el modelo SIMT, donde muchos hilos realizan la misma operación sobre datos diferentes. Esto los hace extremadamente eficientes en tareas como:

- Cálculo de posiciones de vértices.
- Operaciones de sombreado.
- Manipulación de píxeles y fragmentos.

Durante la rasterización y el sombreado, miles de estos núcleos trabajan simultáneamente, permitiendo generar imágenes complejas en tiempo real incluso a resoluciones elevadas.

#### **4.3 Núcleos Tensor: aceleración por hardware para inteligencia artificial y *Deep Learning***

Los núcleos Tensor están optimizados para operaciones de precisión mixta (FP16/INT8), permitiendo realizar multiplicaciones de matrices a una velocidad imposible para los núcleos tradicionales. Son la base del DLSS y del entrenamiento de redes neuronales modernas.

A diferencia de los núcleos tradicionales, los núcleos Tensor:

- Operan con formatos numéricos optimizados (FP16, INT8, BFLOAT16).
- Ejecutan múltiples operaciones en un solo ciclo de reloj.
- Están diseñados para maximizar el rendimiento por vatio.

Gracias a estos núcleos, tareas que antes requerían hardware dedicado o superordenadores pueden ejecutarse en GPUs de consumo. Tecnologías como DLSS (Deep Learning Super Sampling) aprovechan los núcleos Tensor para generar imágenes de alta

resolución a partir de información incompleta, demostrando cómo la IA se integra directamente en el proceso gráfico.

### ***Precisión y rendimiento en IA (FP16, BF16, INT8)***

En cargas de aprendizaje profundo, gran parte del trabajo se reduce a multiplicaciones de matrices; por ello, usar precisión mixta (p. ej., FP16/BF16) o cuantización (INT8) permite aumentar el rendimiento y mejorar el rendimiento por vatio, con una pérdida de precisión que suele ser aceptable según el caso.

*Esta es una de las razones por las que GPUs modernas son tan competitivas en inferencia y entrenamiento frente a CPUs: cuentan con unidades especializadas y rutas de datos optimizadas para estos formatos.*

### **4.4 Núcleos Ray Tracing: unidades de intersección para la simulación física de la luz**

Los núcleos de Ray Tracing son unidades dedicadas a acelerar los cálculos necesarios para el trazado de rayos, una técnica que simula el comportamiento real de la luz mediante la intersección de rayos con objetos de la escena.

Estos cálculos son extremadamente costosos, ya que implican:

- Determinar colisiones entre rayos y geometría.
- Calcular reflexiones, refracciones y sombras.
- Evaluar múltiples rebotes de luz por píxel.

Al delegar estas tareas en hardware específico, la GPU puede ofrecer iluminación fotorrealista en tiempo real, algo impensable hace pocos años. El uso combinado de núcleos de Ray Tracing y técnicas de escalado por IA permite equilibrar realismo visual y rendimiento.

### **4.5 Reloj, frecuencia y overclocking: maximizando el rendimiento del silicio**

La frecuencia de reloj de una GPU indica la cantidad de ciclos que cada núcleo puede ejecutar por segundo y se expresa normalmente en megahercios (MHz) o gigahercios (GHz). En términos prácticos, una mayor frecuencia permite que cada núcleo procese más

instrucciones en el mismo intervalo de tiempo, lo que se traduce en un aumento del rendimiento bruto. Sin embargo, este incremento no es gratuito: a medida que la frecuencia aumenta, también lo hacen el consumo energético y la generación de calor, factores que limitan el rendimiento sostenido del chip.

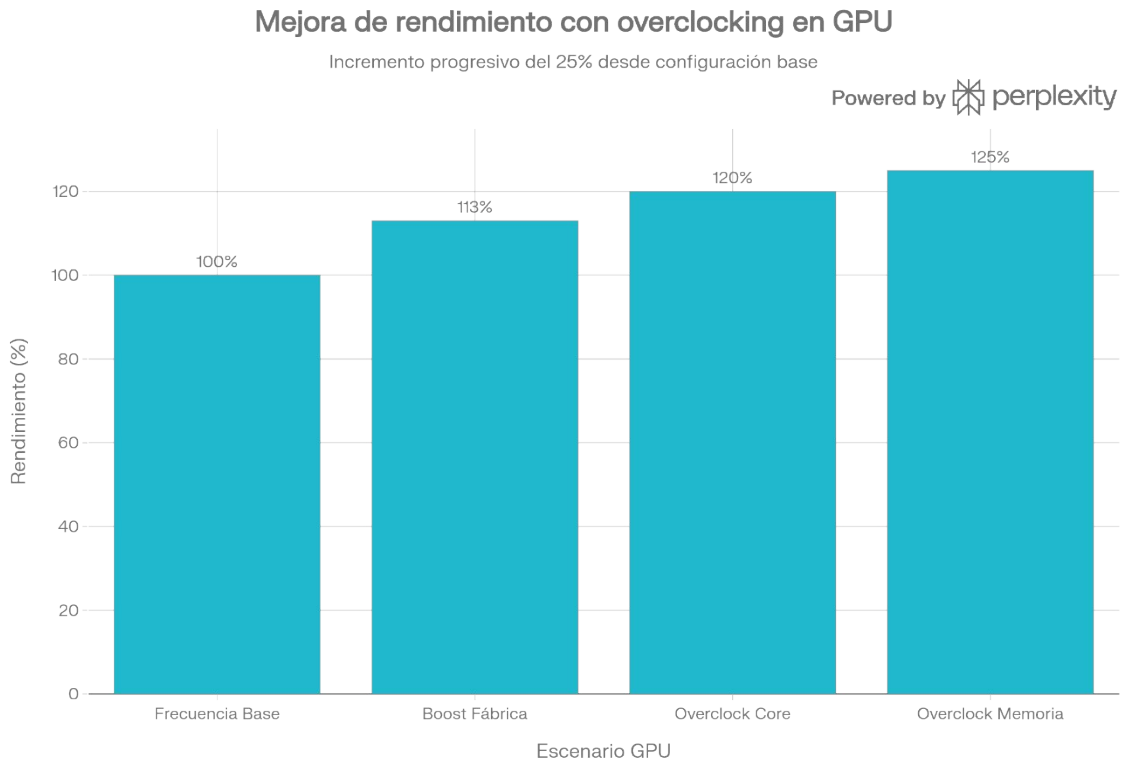
Las GPU modernas no funcionan a una frecuencia fija, sino que ajustan dinámicamente su velocidad en función de la carga de trabajo, la temperatura y el margen energético disponible. Tecnologías como el *boost dinámico* permiten que la GPU alcance frecuencias más altas de forma temporal cuando las condiciones térmicas y eléctricas lo permiten, maximizando el rendimiento sin comprometer la estabilidad del sistema. Esto explica por qué dos tarjetas gráficas con la misma arquitectura pueden ofrecer resultados distintos dependiendo de su diseño térmico y del sistema de refrigeración empleado.

El *overclocking* consiste en forzar manualmente la GPU para que opere a frecuencias superiores a las establecidas por el fabricante. Esta práctica puede proporcionar mejoras de rendimiento apreciables en determinadas aplicaciones, especialmente en escenarios limitados por la velocidad de cálculo. No obstante, también implica un aumento del consumo eléctrico y de la temperatura de funcionamiento, lo que puede provocar inestabilidad, reducción de la vida útil del chip o limitaciones por *thermal throttling* si la refrigeración no es adecuada.

Por este motivo, el rendimiento real de una GPU no depende exclusivamente de su frecuencia de reloj, sino del equilibrio entre múltiples factores. El número de núcleos disponibles, la eficiencia de la arquitectura, la capacidad del sistema de refrigeración y el consumo energético total determinan hasta qué punto una GPU puede mantener frecuencias elevadas de forma sostenida. Una frecuencia alta sin un soporte térmico y energético adecuado rara vez se traduce en un mejor rendimiento práctico.

Finalmente, para que miles de núcleos puedan trabajar de manera simultánea y a altas frecuencias, resulta imprescindible un sistema de memoria capaz de suministrar datos con la rapidez suficiente. Si la GPU no recibe información a la velocidad necesaria, se producen cuellos de botella que limitan el rendimiento, independientemente de la potencia de cálculo

disponible. Este vínculo directo entre frecuencia, consumo y memoria refuerza la idea de que el diseño de una GPU debe entenderse como un sistema equilibrado y no como una suma de cifras aisladas.

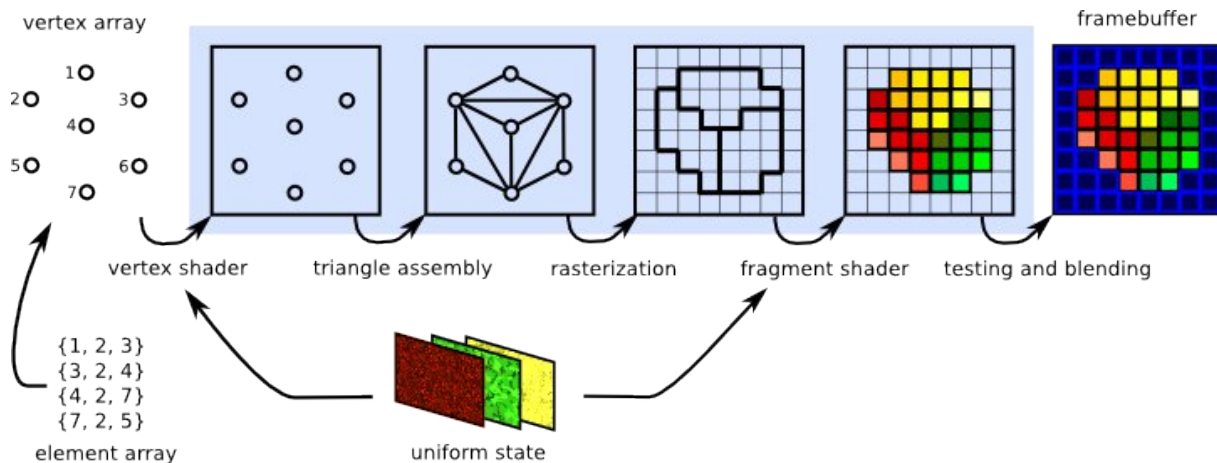


## 5. El pipeline gráfico: del dato matemático al píxel

El pipeline gráfico es la secuencia de etapas que sigue la GPU para transformar datos puramente matemáticos en una imagen visible en pantalla. No es un único proceso, sino una cadena de transformaciones altamente paralelizada, donde cada etapa prepara la información para la siguiente. La clave aquí es entender que la GPU no dibuja imágenes directamente, sino que procesa vértices, primitivas y fragmentos, que más tarde se convierten en píxeles.

En términos generales, el pipeline puede dividirse en cuatro grandes fases:

1. procesamiento de vértices,
2. ensamblado de primitivas,
3. rasterización,
4. procesamiento de fragmentos y salida final.



Cada una de estas fases incluye subetapas internas, muchas de ellas programables mediante *shaders*, lo que da a la GPU su enorme flexibilidad actual.

### 5.1 Datos de entrada: geometría y atributos

Todo comienza con los datos de entrada, que suelen provenir de la CPU. Estos datos describen la escena en forma de vértices, que no son más que puntos en el espacio definidos por coordenadas (x, y, z) y acompañados de atributos adicionales como color, normales o coordenadas de textura.

Es importante entender que en este punto no existe todavía ninguna imagen. La GPU recibe información abstracta: puntos y relaciones matemáticas entre ellos. Por ejemplo, un cubo no se envía como “un cubo”, sino como una lista de vértices y cómo se conectan entre sí. Este enfoque permite a la GPU trabajar de forma extremadamente eficiente y paralela, ya que cada vértice puede procesarse de manera independiente.

## **5.2 Procesamiento de vértices (Vertex Shader)**

En la etapa de procesamiento de vértices, cada vértice pasa por el vertex shader, un pequeño programa que se ejecuta en la GPU. Su función principal es transformar las coordenadas del vértice desde su espacio original (espacio del modelo) hasta el espacio de pantalla, aplicando matrices de transformación como traslación, rotación y proyección.

Además de la transformación espacial, el vertex shader puede modificar o generar atributos adicionales, como calcular la iluminación básica o preparar datos que se usarán más adelante. Lo clave aquí es que cada vértice se procesa de forma independiente, lo que explica por qué las GPUs son tan buenas manejando millones de vértices simultáneamente.

## **5.3 Ensamblado de primitivas**

Una vez transformados, los vértices se ensamblan en primitivas: normalmente triángulos. Esta etapa define qué vértices forman una figura concreta y cómo se conectan entre sí. En la imagen que has elegido, esto se aprecia cuando los puntos aislados pasan a formar una estructura geométrica reconocible.

Aquí también se realiza el clipping, que consiste en descartar las partes de las primitivas que quedan fuera del campo de visión de la cámara. Este paso es fundamental para el rendimiento, ya que evita procesar información que nunca llegará a verse en pantalla.

## **5.4 Rasterización: de geometría a fragmentos**

La rasterización es uno de los pasos conceptualmente más importantes. En esta etapa, las primitivas geométricas (como los triángulos) se convierten en una rejilla de fragmentos, que son candidatos a convertirse en píxeles. Cada fragmento representa una posible contribución a un píxel de la pantalla.

Aquí ocurre la transición clave que muestra tu imagen: se pasa de una figura continua a una representación discreta basada en una cuadrícula. La rasterización calcula qué fragmentos están cubiertos por cada triángulo y asigna a cada uno valores interpolados (color, profundidad, coordenadas de textura) basados en los vértices originales.

### **5.5 Procesamiento de fragmentos (Fragment Shader)**

Cada fragmento generado pasa por el fragment shader, donde se determina su color final. Este es el punto donde se aplican efectos visuales complejos como texturas, sombras, reflejos o iluminación avanzada. A diferencia del vertex shader, aquí el cálculo se realiza por fragmento, lo que puede implicar millones de ejecuciones por fotograma.

Es importante aclarar que fragmento no es lo mismo que píxel. Un fragmento es una propuesta de píxel. Todavía puede ser descartado en etapas posteriores, por ejemplo, si otro fragmento está más cerca de la cámara. Esta distinción explica por qué algunas escenas son tan costosas de renderizar incluso si parecen simples.

### **5.6 Tests finales y salida a pantalla**

Antes de que el fragmento se convierta definitivamente en píxel, pasa por varios tests finales, como el *depth test* (profundidad) o el *stencil test*. Estos tests determinan si el fragmento debe dibujarse o descartarse, asegurando que solo se muestre el elemento visible correcto en cada posición de la pantalla.

Finalmente, los fragmentos que superan todos los tests se escriben en el framebuffer, que es la memoria que contiene la imagen final. Esta imagen es la que el monitor recibirá y mostrará. En este punto, el pipeline ha completado su recorrido: de números abstractos a una imagen visible.

## 6. VRAM: la memoria que alimenta a la GPU

La memoria gráfica, conocida como VRAM, es uno de los componentes más determinantes en el rendimiento real de una tarjeta gráfica. Aunque a menudo se asocia únicamente con la cantidad de gigabytes disponibles, su funcionamiento interno, su velocidad y la forma en que la GPU la gestiona influyen directamente en la fluidez de los gráficos, la estabilidad de los fotogramas y la capacidad para trabajar con resoluciones altas o cargas complejas. Comprender cómo funciona la VRAM permite interpretar correctamente las especificaciones técnicas de una GPU y evitar comparaciones simplistas basadas solo en cifras.

A diferencia de la RAM del sistema, el objetivo principal de la VRAM no es minimizar latencia para tareas variadas, sino mover grandes volúmenes de datos (texturas, buffers, geometría) de forma sostenida para evitar cuellos de botella en el renderizado y el cómputo. Por eso, al evaluar una GPU no basta con mirar los GB: también importa el ancho de bus, la velocidad efectiva de la memoria y la eficiencia del subsistema de cachés.

Un factor crítico que suele pasarse por alto es el bus de memoria, que actúa como la "autopista" que conecta la VRAM con el chip gráfico. Un monitor con gran cantidad de GB puede verse limitado si el bus es estrecho (por ejemplo, 128 bits frente a 384 bits), ya que esto restringe la cantidad de información que puede viajar en cada ciclo de reloj. En resoluciones altas como 4K o 8K, el tamaño de los activos gráficos crece exponencialmente, lo que exige no solo más capacidad, sino una mayor velocidad de transferencia para evitar que la GPU tenga que esperar datos, un fenómeno conocido como "cuello de botella de memoria".

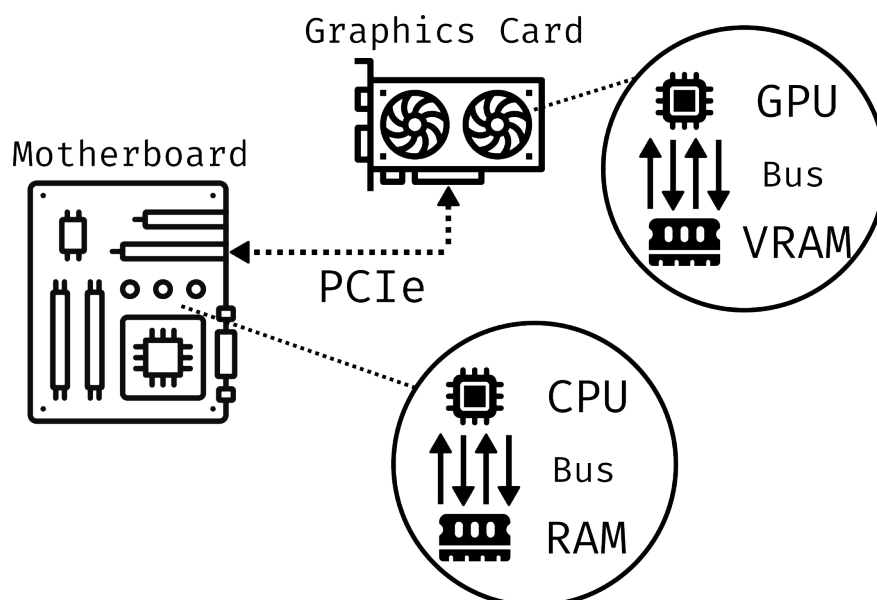
Además, la VRAM gestiona funciones avanzadas como el Anti-Aliasing y el Ray Tracing, que consumen recursos significativos para suavizar bordes y calcular rebotes de luz en tiempo real. Si el consumo de memoria supera la capacidad física disponible, el sistema recurre al "swapping", utilizando la RAM del sistema a través del bus PCIe. Dado que la RAM convencional es drásticamente más lenta que la GDDR6, este proceso provoca una caída súbita en los FPS y el fenómeno de stuttering, destruyendo la fluidez de la experiencia visual independientemente de la potencia del chip gráfico.

Finalmente, es esencial comprender la sinergia entre el tipo de memoria y la eficiencia del software. Tecnologías modernas como Resizable BAR permiten que la CPU acceda a todo el ecosistema de la VRAM de una sola vez en lugar de en pequeños paquetes, optimizando la comunicación entre ambos componentes. Por tanto, la elección de una GPU debe basarse en el equilibrio: una capacidad suficiente para la resolución objetivo, un bus de datos que no estrangule el rendimiento y un estándar de memoria (como GDDR6X) que ofrezca la velocidad necesaria para las exigencias de los motores gráficos de nueva generación.

### 6.1 ¿Qué es la VRAM y por qué es diferente de la RAM del sistema?

La VRAM (Video Random Access Memory) es la memoria dedicada de la tarjeta gráfica y cumple una función crítica: almacenar todos los datos que la GPU necesita de forma inmediata para generar la imagen. Entre estos datos se incluyen texturas, mapas de sombras, búferes de profundidad, geometría, shaders y el propio *framebuffer* que contiene la imagen final antes de enviarse al monitor.

A diferencia de la RAM del sistema, que está pensada para un uso generalista y una baja latencia con la CPU, la VRAM está diseñada para ofrecer un ancho de banda extremadamente alto, necesario para alimentar a miles de núcleos gráficos trabajando en paralelo. Esta diferencia explica por qué una GPU con mucha potencia de cálculo puede verse gravemente limitada si su memoria gráfica es insuficiente o lenta.

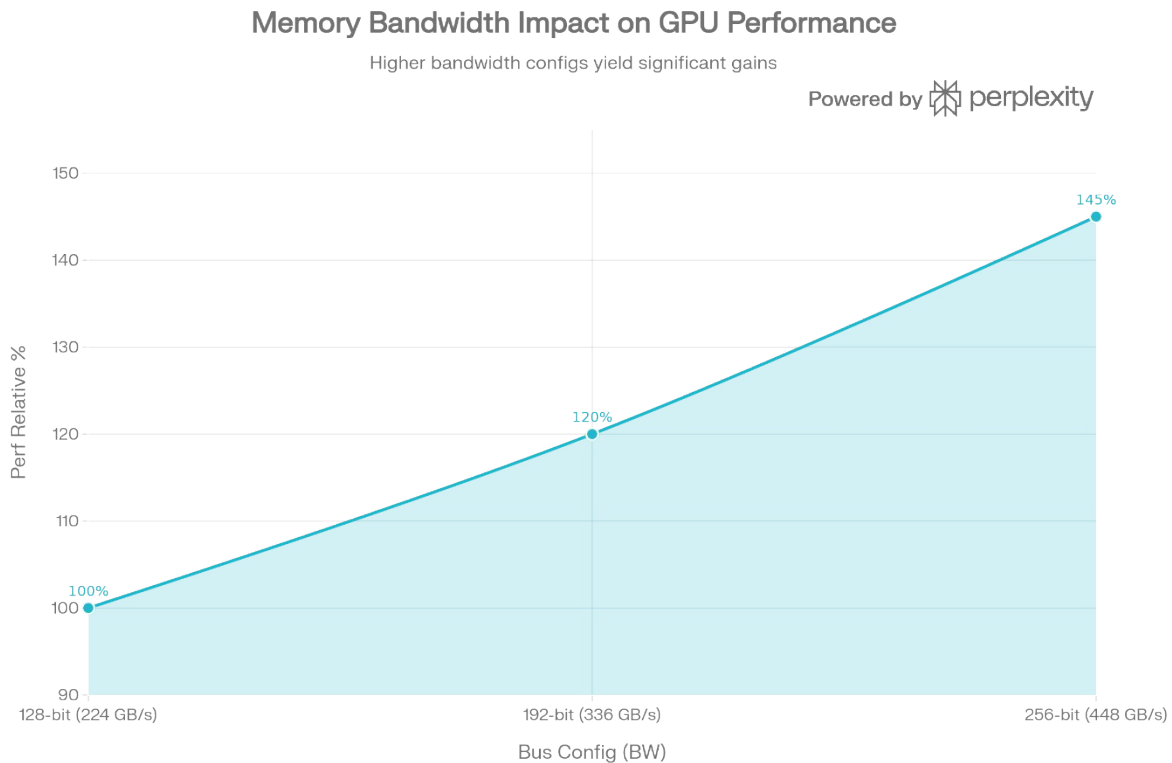


## 6.2 Tipos de memoria gráfica y ancho de banda

Las GPUs modernas utilizan memorias especializadas como GDDR6 y GDDR6X, optimizadas para transferir enormes volúmenes de datos por segundo. No solo importa la cantidad de VRAM (GB), sino también su velocidad efectiva y el ancho del bus de memoria, que determinan el ancho de banda total disponible.

El ancho de banda puede entenderse como la “capacidad de la carretera” por la que viajan los datos entre la VRAM y la GPU. Una tarjeta con 8 GB de VRAM pero un bus estrecho puede rendir peor que otra con la misma cantidad de memoria pero un bus más amplio. Por ello, evaluar una GPU únicamente por los gigabytes de VRAM es un error común que puede llevar a decisiones de compra equivocadas.

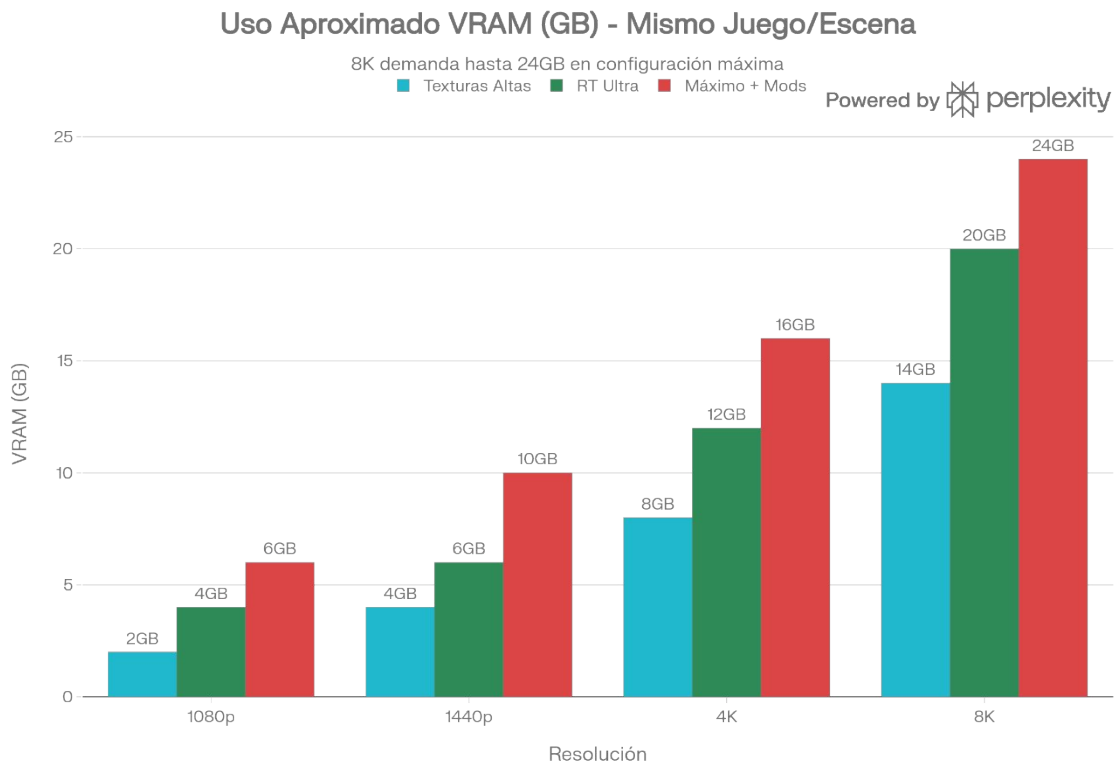
Un error común es fijarse solo en los GB de VRAM. El bus de memoria (ej. 256 bits) determina cuántos datos pueden viajar por ciclo; una tarjeta con mucha memoria pero bus estrecho sufrirá 'stuttering' en altas resoluciones.



### 6.3 VRAM y resolución: por qué cada píxel cuenta

El consumo de VRAM está directamente relacionado con la resolución de salida y la complejidad de la escena renderizada. A mayor resolución, mayor número de píxeles deben almacenarse simultáneamente en memoria, tanto para el framebuffer principal como para estructuras auxiliares como el depth buffer, los mapas de sombras o las texturas en alta definición.

Por ejemplo, pasar de 1080p a 4K no supone simplemente “cuatro veces más píxeles en pantalla”, sino un incremento significativo en la cantidad de datos que la GPU debe gestionar de forma simultánea. Este aumento impacta directamente en el uso de VRAM y explica por qué ciertas configuraciones gráficas funcionan correctamente en resoluciones bajas pero presentan problemas al escalar a resoluciones superiores.

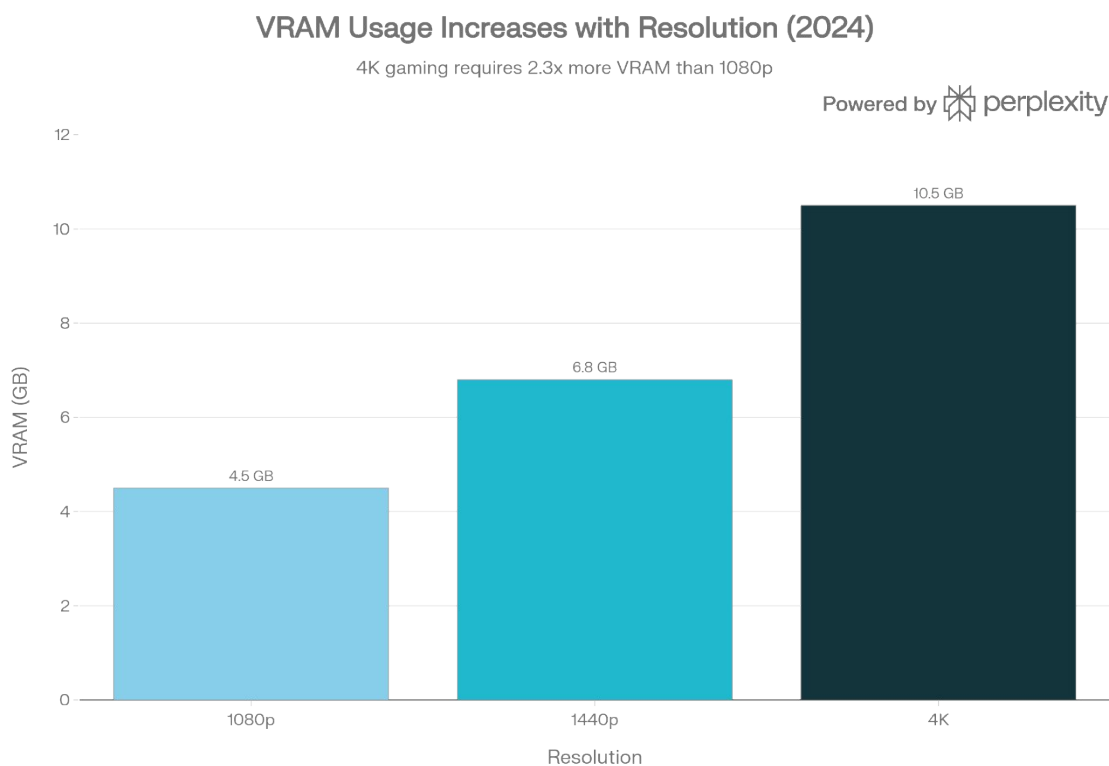


\* *Cyberpunk 2077 ultra RT ejemplo: 4K necesita 12-16GB VRAM*

## 6.4 Gestión de memoria y limitaciones prácticas

Las GPUs modernas incluyen mecanismos avanzados de gestión de memoria, como la compresión de datos y la priorización de recursos, para optimizar el uso de la VRAM disponible. Sin embargo, estos sistemas no pueden compensar una falta estructural de memoria cuando las demandas superan claramente la capacidad física de la tarjeta.

En escenarios como videojuegos de última generación, renderizado 3D o entrenamiento de modelos de inteligencia artificial, la VRAM se convierte en un factor limitante real. En estos casos, disponer de más memoria no aumenta directamente los fotogramas por segundo, pero evita cuellos de botella que afectan a la estabilidad y consistencia del rendimiento.



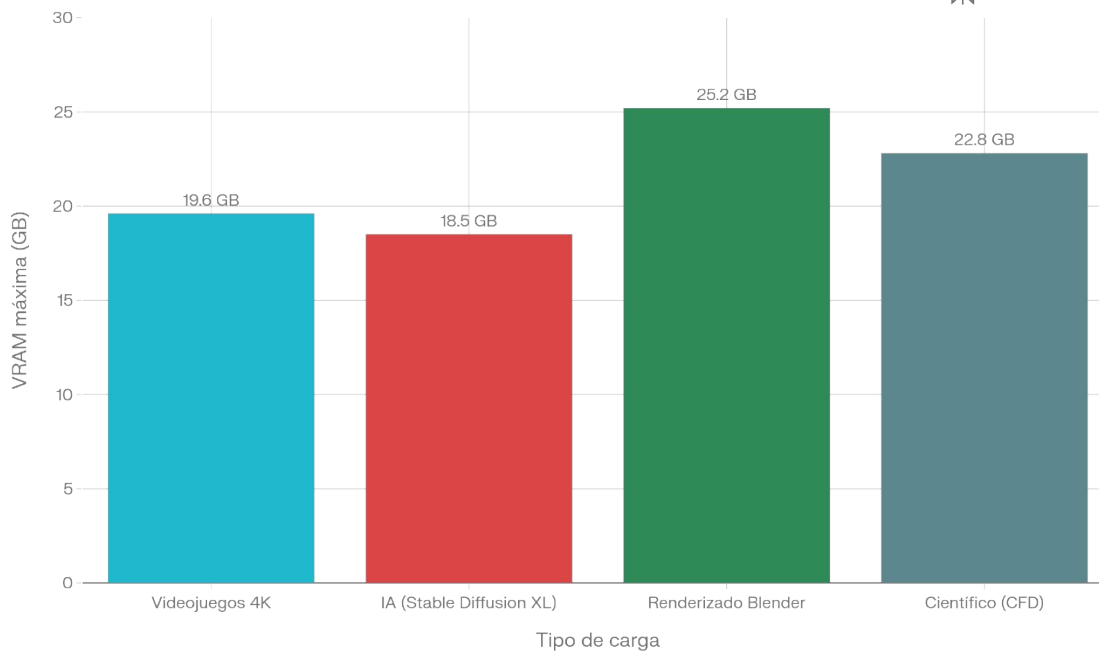
El uso de la VRAM varía significativamente según el tipo de carga de trabajo. En videojuegos, la memoria gráfica se utiliza principalmente para almacenar texturas, mapas de sombras, buffers de imagen y datos temporales del pipeline gráfico. En cambio, en aplicaciones de inteligencia artificial, renderizado o cálculo científico, la VRAM se emplea para almacenar

modelos completos, matrices de datos y tensores, lo que provoca consumos más sostenidos y menos dependientes de la resolución. Esta diferencia explica por qué ciertas aplicaciones profesionales requieren grandes cantidades de VRAM incluso sin mostrar gráficos complejos en pantalla.

### Uso máximo VRAM por tipo de carga (RTX 5090)

El renderizado 3D requiere mayor memoria de video que otras cargas

Powered by  perplexity



### 6.5 El mito de “más VRAM siempre es mejor”

Uno de los errores más extendidos en el mercado es asumir que una tarjeta gráfica con más VRAM es automáticamente superior. En realidad, la VRAM solo es útil si la GPU tiene la potencia de cálculo suficiente para aprovecharla. Una tarjeta de gama baja con mucha memoria no puede compensar la falta de núcleos, ancho de banda o frecuencia de reloj.

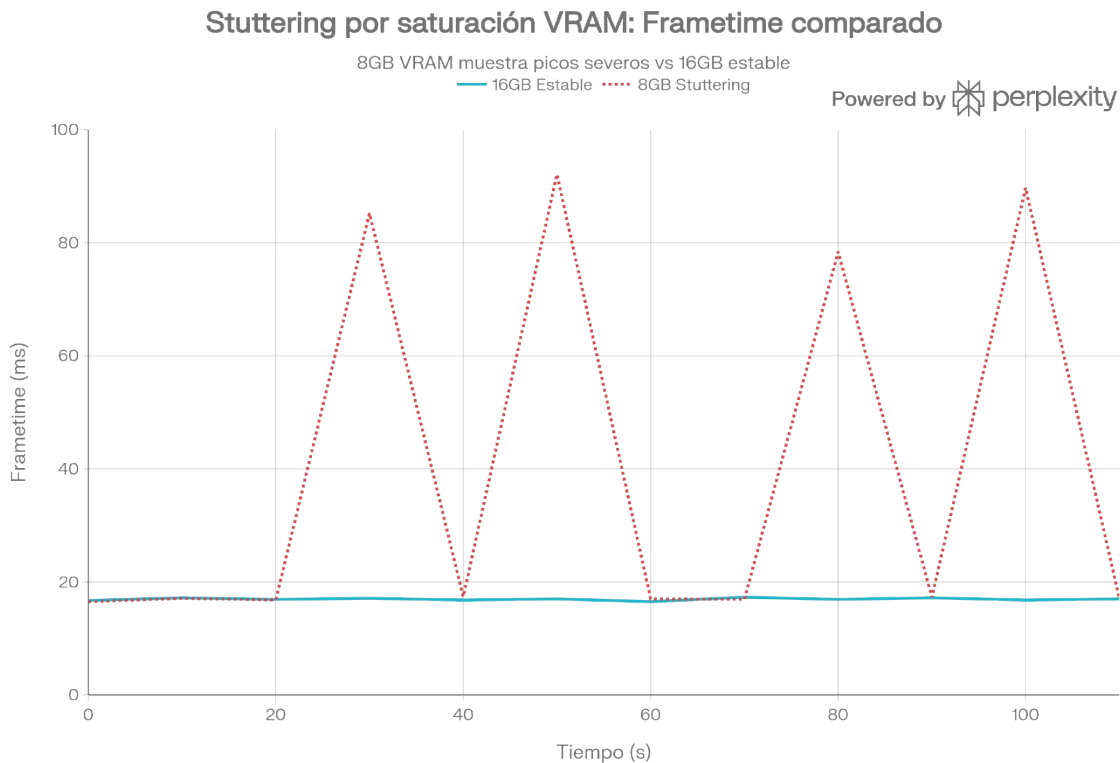
Por ello, la VRAM debe evaluarse siempre en conjunto con el resto de la arquitectura de la GPU y el uso previsto. Para resoluciones estándar y tareas ligeras, una cantidad moderada de VRAM es suficiente. En cambio, para cargas profesionales o resoluciones elevadas, la memoria gráfica pasa a ser un requisito crítico más que un simple valor numérico.

## 6.6 Cuando la VRAM no es suficiente: stuttering y pérdida de fluidez

Cuando la memoria gráfica disponible se agota, la GPU se ve obligada a intercambiar datos con la RAM del sistema a través del bus PCI Express, un proceso mucho más lento que el acceso directo a la VRAM. Este fenómeno no siempre se refleja en una caída brusca de los FPS medios, pero sí provoca microparones, irregularidades en el tiempo de fotograma (frametime) y una sensación de falta de fluidez conocida como stuttering.

En la práctica, el stuttering suele manifestarse como picos en el tiempo de fotograma (frametime spikes): los FPS medios pueden parecer “aceptables”, pero aparecen microparones visibles. Una causa típica es saturar la VRAM: el sistema empieza a mover recursos fuera de la VRAM y traerlos de vuelta, generando interrupciones y latencia.

Para diagnosticarlo, conviene registrar simultáneamente frametime, uso de VRAM y uso de RAM del sistema; si el stutter coincide con VRAM al límite, el problema suele ser memoria/streaming de assets más que “potencia bruta”.



## 7. Tarjetas gráficas integradas y dedicadas

Una vez comprendido el papel de la GPU y su memoria, surge una de las decisiones más importantes a la hora de configurar un sistema: utilizar una GPU integrada o una dedicada.

La elección entre una GPU integrada (iGPU) y una dedicada (dGPU) define la arquitectura fundamental del sistema y su capacidad de procesamiento paralelo. Mientras que las soluciones integradas forman parte del propio silicio del procesador (SoC - System on a Chip) y comparten recursos térmicos y de memoria, las tarjetas dedicadas operan como un ecosistema independiente con su propia infraestructura de alimentación y refrigeración. Esta distinción no solo determina la potencia bruta, sino que establece compromisos críticos entre la eficiencia energética y la capacidad de cómputo masivo.

Técnicamente, la mayor divergencia reside en la gestión de la memoria y el ancho de banda. Las iGPU dependen de la RAM del sistema, compartiendo el bus de datos con la CPU, lo que limita significativamente su throughput en tareas complejas. Por el contrario, las dGPU cuentan con VRAM dedicada de alta velocidad, diseñada específicamente para alimentar a miles de núcleos sin interferir en los procesos generales del sistema, evitando así los cuellos de botella que degradan la fluidez visual.

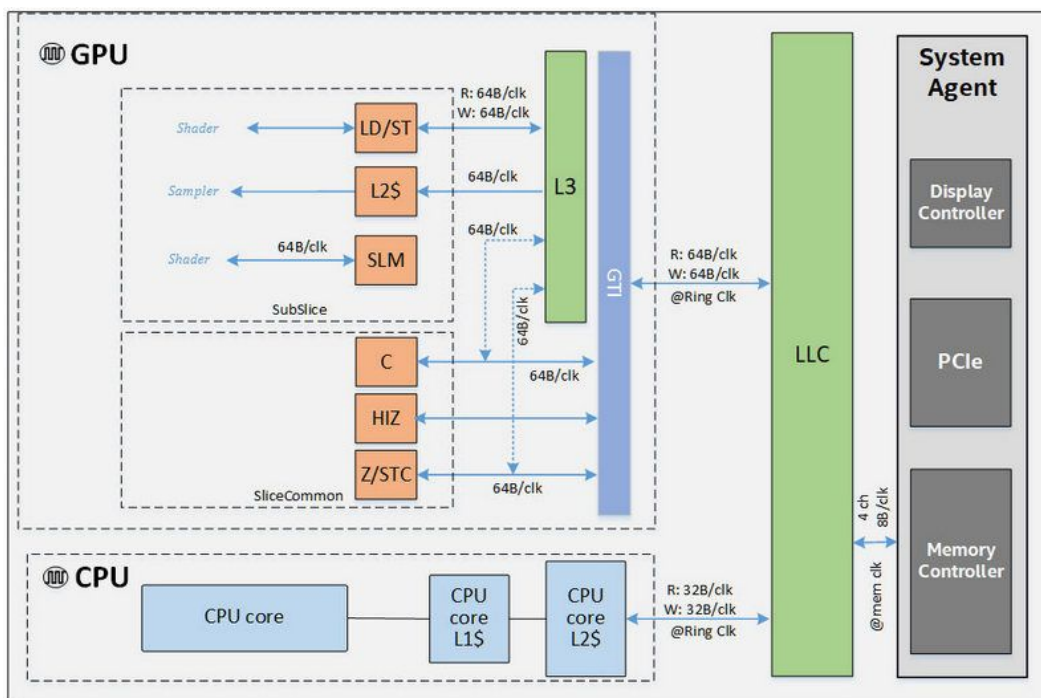
Este capítulo analiza ambas soluciones bajo un prisma técnico y práctico, evaluando cómo factores como el TDP (Potencia de Diseño Térmico) y la latencia de memoria definen sus límites operativos. Más allá de una comparativa de rendimiento, se exploran los escenarios óptimos para cada arquitectura, desde la movilidad extrema y el bajo consumo de las integradas hasta la computación de alto rendimiento, el renderizado profesional y la inteligencia artificial que demandan las soluciones dedicadas.

### 7.1 GPU integrada (iGPU): eficiencia y limitaciones

Las GPUs integradas (iGPU) son unidades gráficas que forman parte del mismo chip que la CPU o del propio System on a Chip (SoC). En lugar de disponer de memoria gráfica dedicada, utilizan la RAM del sistema como memoria compartida, lo que reduce costes,

consumo energético y complejidad del sistema. Este diseño es habitual en procesadores de Intel, AMD y Apple, especialmente en portátiles y equipos compactos.

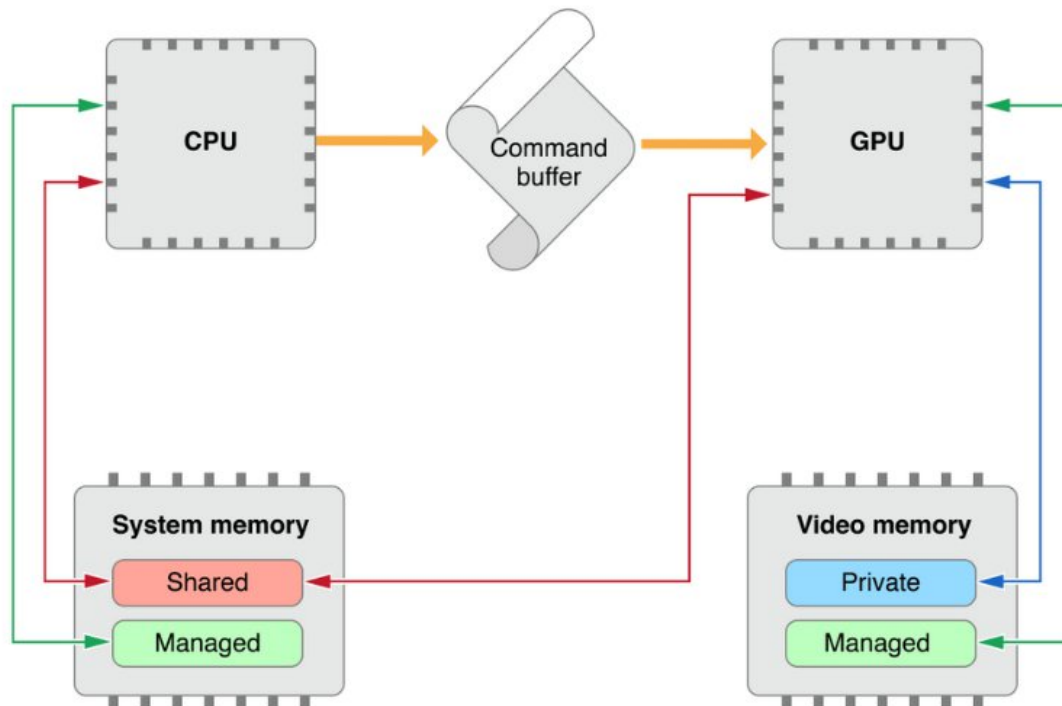
El principal punto fuerte de una iGPU es su eficiencia energética. Al compartir silicio y recursos con la CPU, el consumo es mucho menor que el de una tarjeta dedicada, lo que se traduce en menos calor, menos ruido y mayor autonomía en dispositivos portátiles. Sin embargo, esta eficiencia tiene un coste: el rendimiento gráfico es limitado, especialmente en tareas que requieren alto paralelismo o gran ancho de banda de memoria, como videojuegos modernos, renderizado 3D o inteligencia artificial.



## 7.2 GPU dedicada (dGPU): potencia y especialización

Una GPU dedicada (dGPU) es una tarjeta gráfica independiente que se conecta a la placa base a través de una interfaz como PCI Express. A diferencia de la iGPU, dispone de su propio chip gráfico, VRAM dedicada y un sistema de alimentación y refrigeración independiente. Este diseño permite alcanzar niveles de rendimiento muy superiores, tanto en gráficos como en computación general.

Las dGPUs están diseñadas para manejar cargas intensivas y sostenidas: videojuegos exigentes, renderizado profesional, simulaciones científicas o entrenamiento de modelos de inteligencia artificial. Gracias a su arquitectura altamente paralela y a su gran ancho de banda de memoria, pueden procesar millones de operaciones simultáneamente. Como contrapartida, su consumo energético y coste económico son significativamente mayores, lo que obliga a considerar factores como la fuente de alimentación y la refrigeración del sistema.



### 7.3 Comparación directa: iGPU vs dGPU

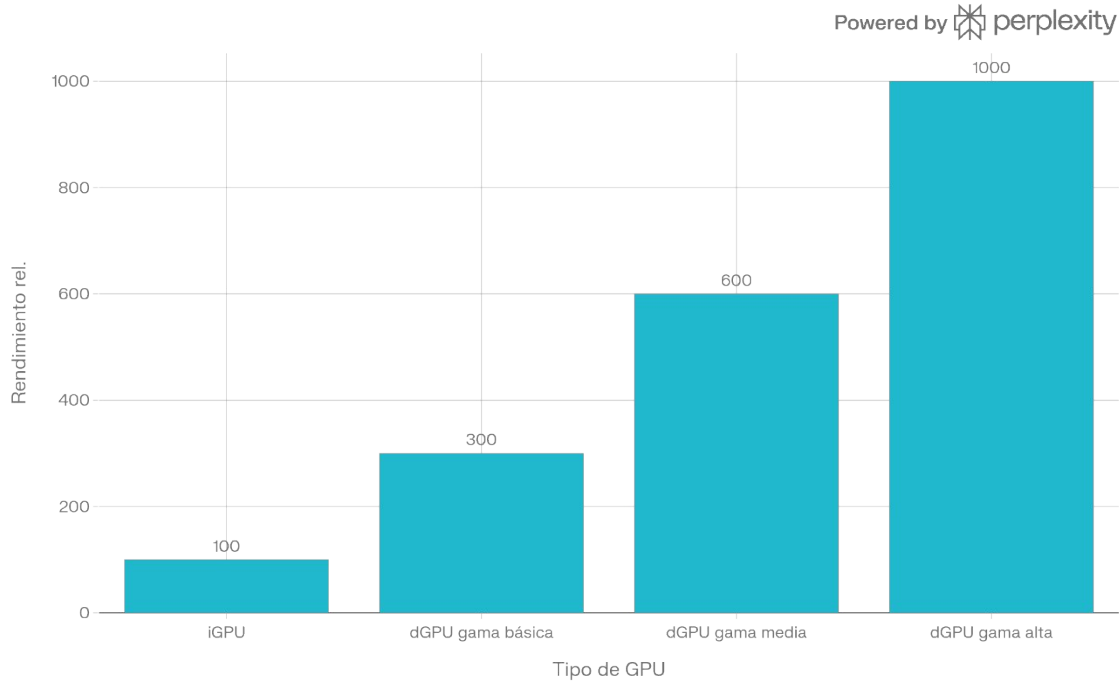
La diferencia entre una GPU integrada y una dedicada no se limita únicamente al rendimiento bruto. Se trata de dos filosofías de diseño distintas, orientadas a escenarios de uso diferentes. Mientras que la iGPU prioriza la eficiencia y la integración, la dGPU apuesta por la potencia y la especialización.

Desde el punto de vista técnico, la clave está en el ancho de banda de memoria, el número de núcleos de procesamiento y la capacidad de mantener cargas sostenidas sin estrangulamiento térmico (thermal throttling). Estas diferencias hacen que una iGPU sea

adecuada para tareas cotidianas y multimedia, mientras que una dGPU sea prácticamente imprescindible para aplicaciones gráficas o computacionales exigentes.

### Rendimiento relativo según tipo de GPU

Las GPUs dedicadas superan significativamente a las integradas



*\* Una dGPU no es “un poco mejor”, es un orden de magnitud superior.*

## 7.4 Consumo energético y eficiencia

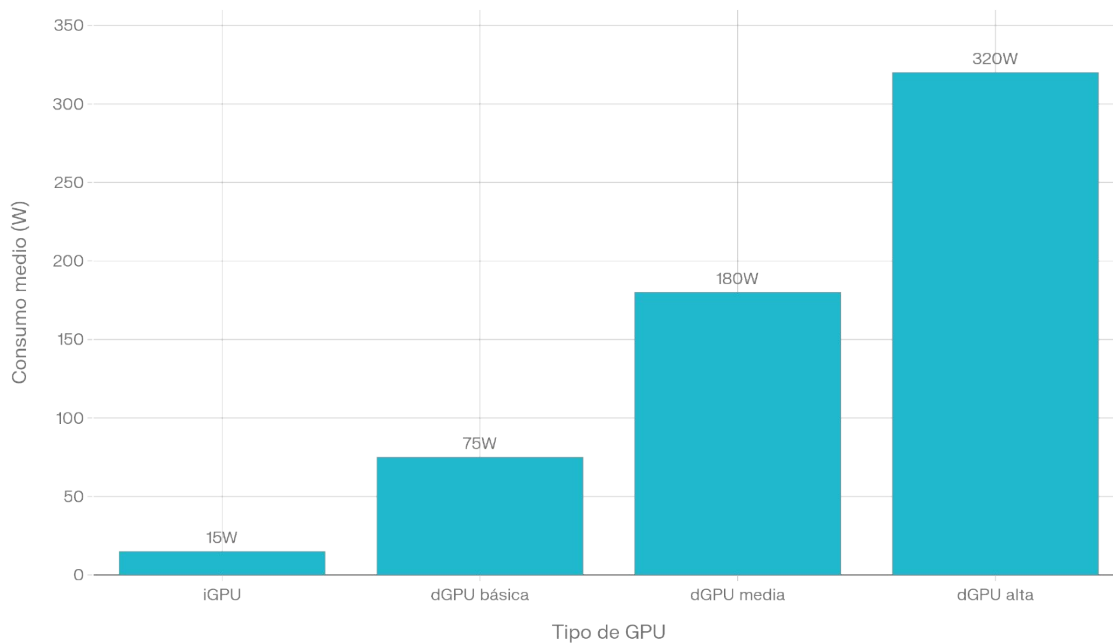
El consumo energético es uno de los factores más diferenciadores entre iGPU y dGPU. Una iGPU suele consumir entre 10 y 30 W, mientras que una GPU dedicada puede superar fácilmente los 200 W en carga. Esta diferencia no solo afecta a la factura eléctrica, sino también al diseño del sistema: refrigeración, ruido y tamaño del equipo.

Desde una perspectiva de eficiencia, la iGPU ofrece un rendimiento por vatio excelente en tareas ligeras, mientras que la dGPU sacrifica eficiencia a cambio de potencia bruta. Por ello, la elección entre ambas debe tener en cuenta no solo el rendimiento necesario, sino también el entorno de uso y las limitaciones físicas del equipo.

### Consumo energético medio según tipo de GPU

Mayor potencia de GPU requiere significativamente más energía

Powered by  perplexity



*\* Más rendimiento implica más consumo y más requisitos.*

## 7.5 Escenarios de uso: ¿cuándo es suficiente una iGPU?

No todos los usuarios necesitan una GPU dedicada. Para tareas como navegación web, ofimática, reproducción de vídeo en alta resolución o incluso programación básica, una iGPU moderna es más que suficiente. Además, muchas iGPUs actuales incluyen aceleración por hardware para vídeo y ciertas cargas de IA ligera.

En equipos portátiles y dispositivos compactos, la elección entre una GPU integrada o dedicada tiene implicaciones especialmente relevantes. Las iGPU suelen ofrecer una mayor autonomía y un diseño térmico más simple, al compartir recursos con la CPU. Por el contrario, las dGPU aportan un aumento considerable de rendimiento, pero a costa de un mayor consumo energético y generación de calor, lo que puede provocar limitaciones por thermal throttling en chasis pequeños. Esta relación entre rendimiento, temperatura y consumo es uno de los factores más determinantes en el diseño de portátiles modernos.

Sin embargo, cuando entran en juego videojuegos modernos, edición de vídeo, modelado 3D o aprendizaje automático, las limitaciones de la iGPU se hacen evidentes. En estos casos, una dGPU no es un lujo, sino un requisito técnico para garantizar fluidez y estabilidad.

Escenario	iGPU	dGPU
<b>Productividad y Movilidad</b>	<input checked="" type="checkbox"/> Óptimo	<input checked="" type="checkbox"/> Ineficiente
<b>Gaming y E-sports</b>	<input type="checkbox"/> 1080p bajo	<input checked="" type="checkbox"/> Superior
<b>Edición de Vídeo y Diseño</b>	<input checked="" type="checkbox"/> Básico	<input checked="" type="checkbox"/> Acelerado (CUDA)
<b>IA, Deep Learning y Render 3D</b>	<input checked="" type="checkbox"/> VRAM limitada	<input checked="" type="checkbox"/> Profesional

*\* La iGPU sirve para todo... pero no rinde bien en todo.*

## 8. Conexión, alimentación y refrigeración

El rendimiento de una GPU no depende únicamente de su arquitectura interna, sino también de cómo se integra físicamente en el sistema. La interfaz de conexión con la placa base y el sistema de refrigeración condicionan tanto la estabilidad como la capacidad de mantener altas frecuencias de trabajo de forma sostenida. Una GPU potente mal refrigerada o limitada por su conexión puede ofrecer un rendimiento inferior al esperado, lo que convierte estos aspectos en elementos clave del diseño y la elección del hardware gráfico.

En este capítulo se analizan los principales elementos físicos que permiten a una GPU operar a pleno rendimiento, destacando cómo una mala elección en estos aspectos puede limitar incluso a las tarjetas gráficas más potentes.



## 8.1 Interfaz PCI Express: la autopista de datos

La conexión entre la tarjeta gráfica y el resto del sistema se realiza a través de la interfaz PCI Express (PCIe). Esta interfaz actúa como una autopista de alta velocidad por la que viajan datos entre la GPU, la CPU y la memoria del sistema. Actualmente, las GPUs utilizan normalmente un slot PCIe x16, que proporciona el mayor ancho de banda disponible.

Con cada nueva generación (PCIe 3.0, 4.0, 5.0), el ancho de banda se duplica, permitiendo transferencias más rápidas sin cambiar el formato físico del conector. En la práctica, una GPU moderna puede funcionar en una placa base más antigua, aunque con una ligera pérdida de rendimiento en escenarios muy concretos. Esto convierte a PCIe en una interfaz retrocompatible y escalable, clave para la evolución del hardware gráfico.

Versión	Año	Ancho de banda x16 (GB/s)	Impacto en GPUs Modernas
<b>PCIe 3.0</b>	2010	16 GB	Suficiente para gaming 4K <a href="#">beforethetrashcan</a>
<b>PCIe 4.0</b>	2019	32 GB	Estándar actual, mínimo bottleneck
<b>PCIe 5.0</b>	2023	64 GB	Beneficioso en ML/render <a href="#">beforethetrashcan</a>
<b>PCIe 6.0</b>	2024	128 GB	Futuro para IA masiva/8K pcisig+1

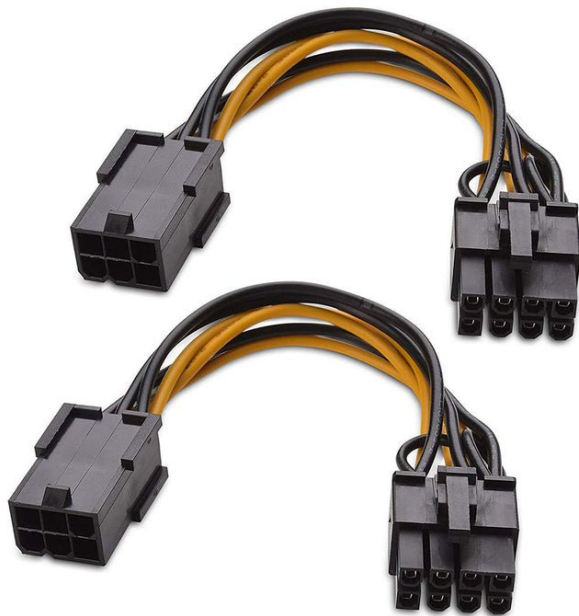
*\* El bus raramente es el cuello de botella, pero evoluciona para el futuro.*

## 8.2 Alimentación eléctrica: potencia y estabilidad

Las GPUs dedicadas modernas requieren una alimentación eléctrica significativa, muy superior a la que puede proporcionar el propio slot PCIe. Por este motivo, incorporan conectores de alimentación adicionales procedentes directamente de la fuente de alimentación (PSU). Tradicionalmente se han utilizado conectores de 6 y 8 pines, pero las GPUs más recientes emplean el nuevo estándar 12VHPWR, capaz de suministrar potencias mucho más elevadas.

El consumo energético de una GPU se expresa habitualmente como TGP (Total Graphics Power) o TDP, y representa la potencia máxima que puede demandar bajo carga. Una fuente de alimentación insuficiente o de baja calidad puede provocar inestabilidad, caídas de rendimiento o incluso apagados repentinos. Por ello, la elección de la PSU es tan importante como la propia tarjeta gráfica.

Las tarjetas de gama alta actuales utilizan el conector 12VHPWR de 16 pines, capaz de suministrar hasta 600W, lo que requiere fuentes de alimentación con certificación ATX 3.0 para garantizar la estabilidad energética.

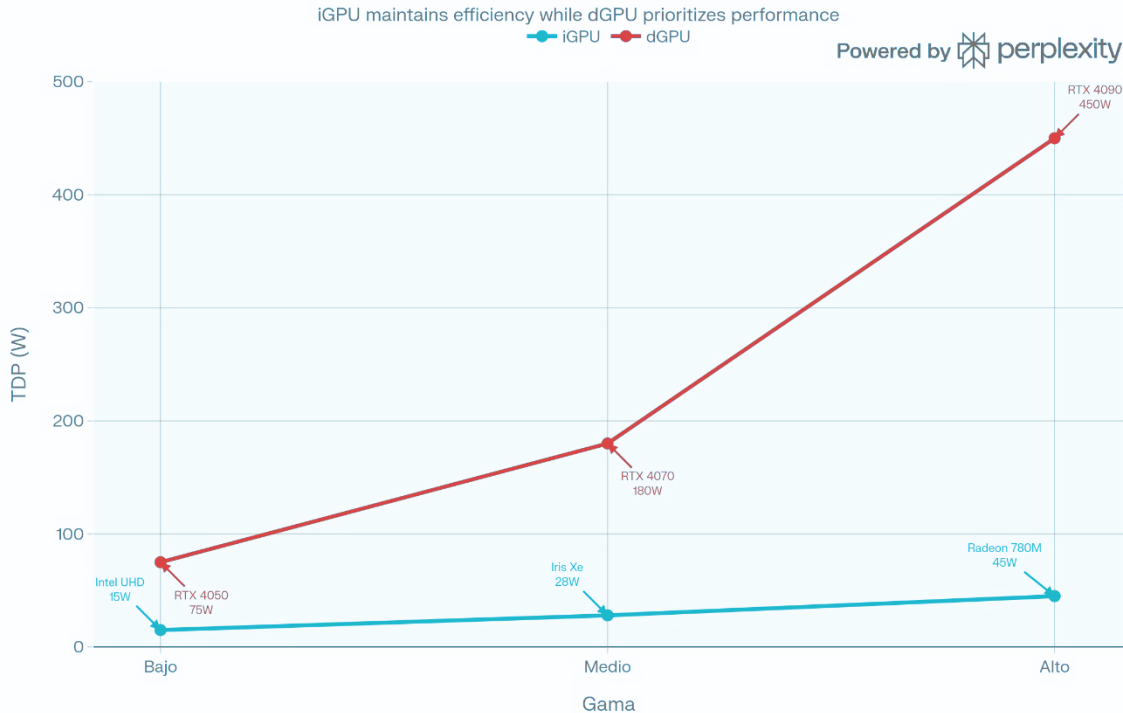


*Conector de 6 pines\**



*Conector 12VHPWR\**

### Consumo Energético: iGPU vs dGPU por Gama



\* Comparación de consumo TDP: iGPU ultraeficiente vs dGPU de alto rendimiento

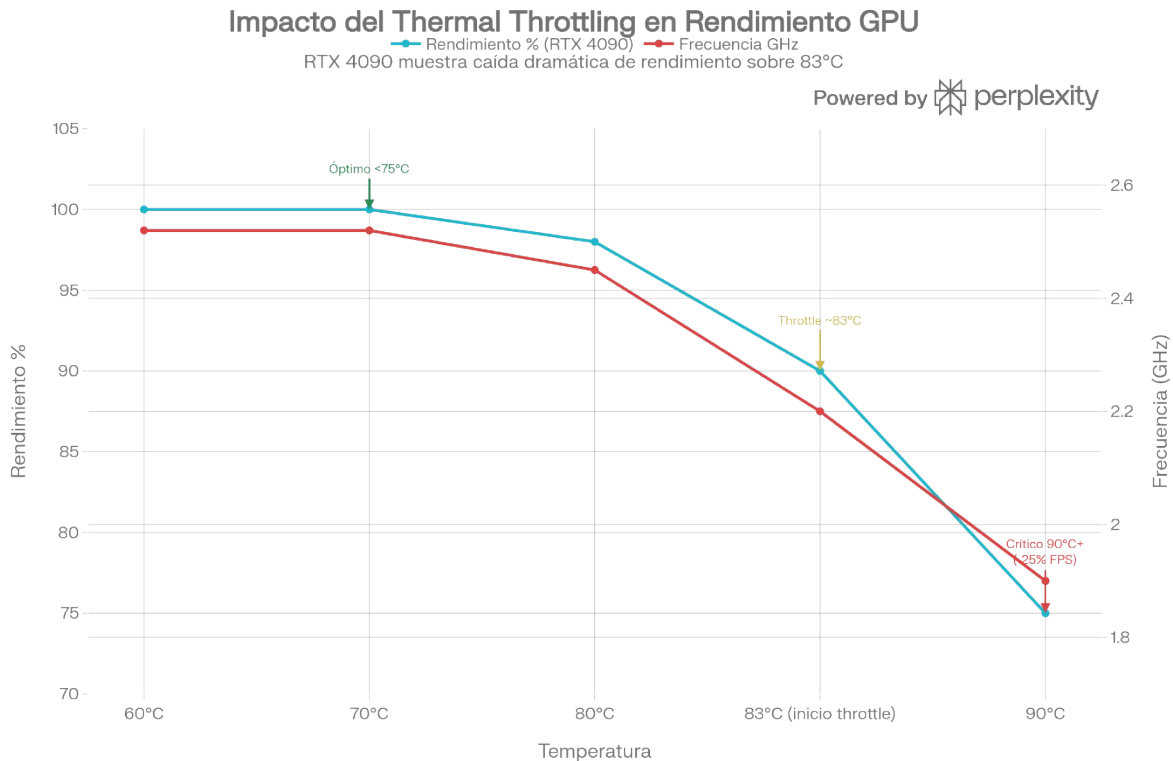
El consumo de una iGPU permanece ultraeficiente (15-45 W) incluso en gama alta, ideal para laptops y sistemas silenciosos, mientras que las dGPU escalan drásticamente hasta 450 W en alto rendimiento para gaming/IA.

### 8.3 Refrigeración: control térmico y rendimiento sostenido

Durante su funcionamiento, una GPU genera una gran cantidad de calor debido a la alta densidad de transistores y a su elevado consumo energético. Para evitar daños y mantener un rendimiento estable, es imprescindible un sistema de refrigeración eficiente. La mayoría de las tarjetas utilizan disipación por aire, combinando grandes disipadores de aluminio con ventiladores axiales.

En modelos de gama alta o entornos profesionales, también es habitual el uso de refrigeración líquida, que permite una disipación térmica más eficiente y silenciosa. Una refrigeración inadecuada provoca thermal throttling, un mecanismo de protección mediante

el cual la GPU reduce automáticamente su frecuencia para evitar sobrecalentamiento, lo que impacta directamente en el rendimiento.



\* Caída de rendimiento y frecuencia por thermal throttling (ejemplo RTX 4090)

El thermal throttling es un mecanismo de protección: cuando la GPU alcanza límites térmicos definidos por firmware/driver, reduce frecuencia y/o voltaje para mantenerse dentro de un rango seguro.

El punto exacto de throttling depende del modelo, del límite de potencia configurado y del diseño térmico (disipador, flujo de aire, pasta térmica), por lo que no existe un único umbral universal válido para todas las GPUs.

El efecto práctico es una caída de rendimiento sostenido (no solo un pico puntual) y, en escenarios exigentes, una experiencia menos estable.

## **8.4 Impacto conjunto en el sistema**

La conexión, la alimentación y la refrigeración no son factores aislados, sino elementos que deben estar equilibrados dentro del sistema. Una GPU potente instalada en un equipo con mala ventilación o una fuente insuficiente no solo no alcanza su máximo rendimiento, sino que puede comprometer la estabilidad general del sistema.

Este equilibrio entre componentes es fundamental para evitar cuellos de botella físicos y garantizar una experiencia fluida y segura. Entender estos aspectos permite al usuario no solo elegir mejor una tarjeta gráfica, sino también diseñar un sistema coherente y eficiente.

## 9. Mercado actual y modelos de tarjetas gráficas

Tras comprender la arquitectura interna de una GPU, su memoria, su integración física en el sistema y sus requisitos energéticos, resulta imprescindible analizar cómo todo este conocimiento se traduce en el mercado real de tarjetas gráficas. La oferta actual es amplia, diversa y, en muchos casos, confusa para el usuario que no dispone de criterios técnicos claros.

El mercado de las tarjetas gráficas ha evolucionado más allá del simple aumento de rendimiento generación tras generación. Factores como la escasez de componentes, la demanda procedente de nuevos sectores y el desarrollo de tecnologías exclusivas han modificado tanto los precios como la percepción del valor de una GPU. Analizar el mercado permite comprender por qué dos tarjetas con especificaciones similares pueden ofrecer experiencias muy distintas y por qué el precio no siempre refleja únicamente el rendimiento bruto.

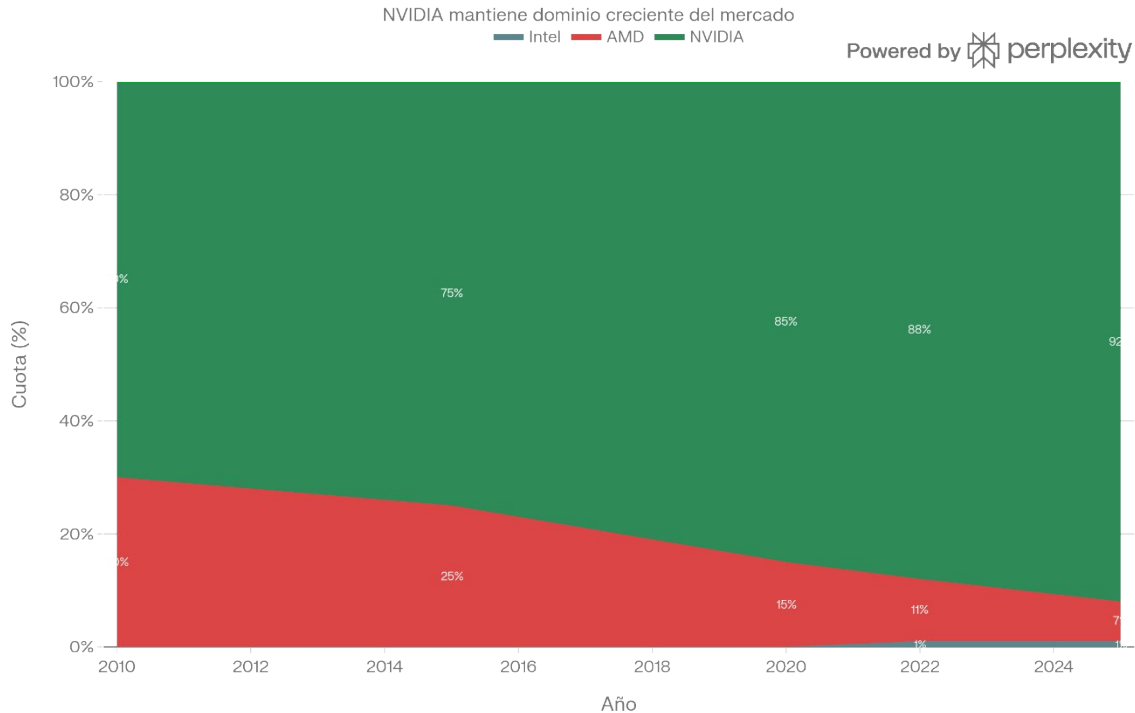
Este capítulo tiene como objetivo contextualizar las GPUs dentro de su ecosistema comercial, identificando los principales fabricantes, el papel de los ensambladores y la segmentación por gamas. De este modo, se establecen las bases necesarias para interpretar correctamente especificaciones técnicas, evitar decisiones basadas únicamente en marketing y seleccionar una tarjeta gráfica adecuada al uso previsto.

### 9.1 El triunvirato del silicio: NVIDIA, AMD e Intel

El mercado de GPUs dedicadas está dominado por tres grandes fabricantes de chips gráficos: NVIDIA, AMD e Intel. Estas empresas diseñan la arquitectura interna de la GPU, desarrollan los controladores y marcan la hoja de ruta tecnológica del sector. Aunque compiten entre sí, cada una tiene enfoques y fortalezas diferenciadas.

NVIDIA lidera históricamente en rendimiento bruto, tecnologías propietarias y computación acelerada por GPU, especialmente en inteligencia artificial. AMD destaca por ofrecer una mejor relación rendimiento/precio y arquitecturas abiertas, mientras que Intel, el actor más reciente en GPUs dedicadas, busca posicionarse como una alternativa competitiva integrando gráficos y CPU dentro de un mismo ecosistema.

## Cuota de Mercado Histórica Tarjetas Gráficas (dGPU Discrete)



*\* NVIDIA domina el mercado de tarjetas gráficas discretas (dGPU) con 92% cuota en 2025 (vs AMD 7%, Intel 1%), según datos JPR Q3 2025.*

### 9.2 Ensambladores (AIB): más allá del chip gráfico

Aunque NVIDIA, AMD e Intel diseñan las GPUs, la mayoría de las tarjetas gráficas comerciales son fabricadas por empresas conocidas como AIB (Add-In Board partners), como ASUS, MSI, Gigabyte, Sapphire o Zotac. Estas compañías se encargan del diseño del sistema de refrigeración, la alimentación y el formato físico de la tarjeta.

Un mismo chip gráfico puede encontrarse en múltiples versiones con diferencias significativas en ruido, temperatura, tamaño y estabilidad. Por ello, no todas las tarjetas con la misma GPU ofrecen exactamente la misma experiencia. Este aspecto es clave a la hora de elegir un modelo concreto y demuestra que el rendimiento no depende únicamente del silicio.

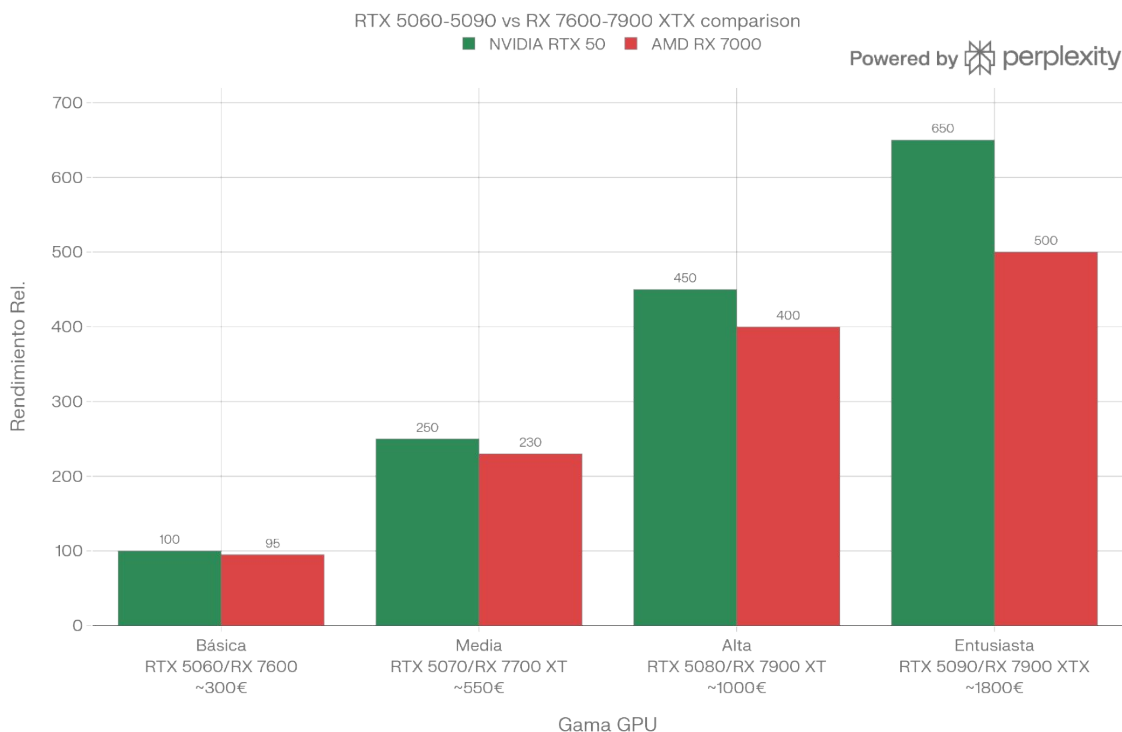
### 9.3 Segmentación por gamas: básica, media y alta

El mercado de tarjetas gráficas se estructura habitualmente en gamas, que agrupan modelos con niveles de rendimiento similares. Esta segmentación facilita la elección según el

uso previsto y el presupuesto disponible. Entender las gamas es más útil que conocer modelos concretos, ya que estas cambian con cada generación.

La gama básica está orientada a tareas ligeras y gaming ocasional; la gama media ofrece un equilibrio entre precio y rendimiento; y la gama alta está pensada para usuarios exigentes, resoluciones elevadas y cargas profesionales. Cada salto de gama implica un incremento notable en consumo, precio y requisitos del sistema.

### Rendimiento Relativo por Gama: NVIDIA vs AMD (2026)



*\* Comparativa rendimiento NVIDIA RTX 50 vs AMD RX 7000 por gama (2026)*

## 9.4 Aplicaciones apropiadas para cada gama

Cada gama de GPU responde a necesidades concretas. Utilizar una tarjeta de gama alta para tareas básicas supone un desperdicio de recursos, mientras que intentar ejecutar cargas exigentes con una GPU básica conduce a una experiencia deficiente. Por ello, la elección debe basarse en el uso real, no en cifras de marketing.

Este enfoque permite optimizar presupuesto, consumo y rendimiento. Además, refuerza la idea de que una buena elección no es la más potente, sino la más adecuada al contexto.

Uso	Básica	Media	Alta	Recomendación
Ofimática	100%	100%	100%	Cualquiera (básica sobra)
Multimedia	80%	100%	100%	Media/Alta
Gaming 1080p	60%	100%	100%	Media ideal
Gaming 4K	20%	60%	100%	Alta imprescindible
Render / IA	10%	50%	100%	Alta esencial

*\* Cada gama tiene su lugar y su propósito.*

## 9.5 Tendencias del mercado y ciclos de renovación

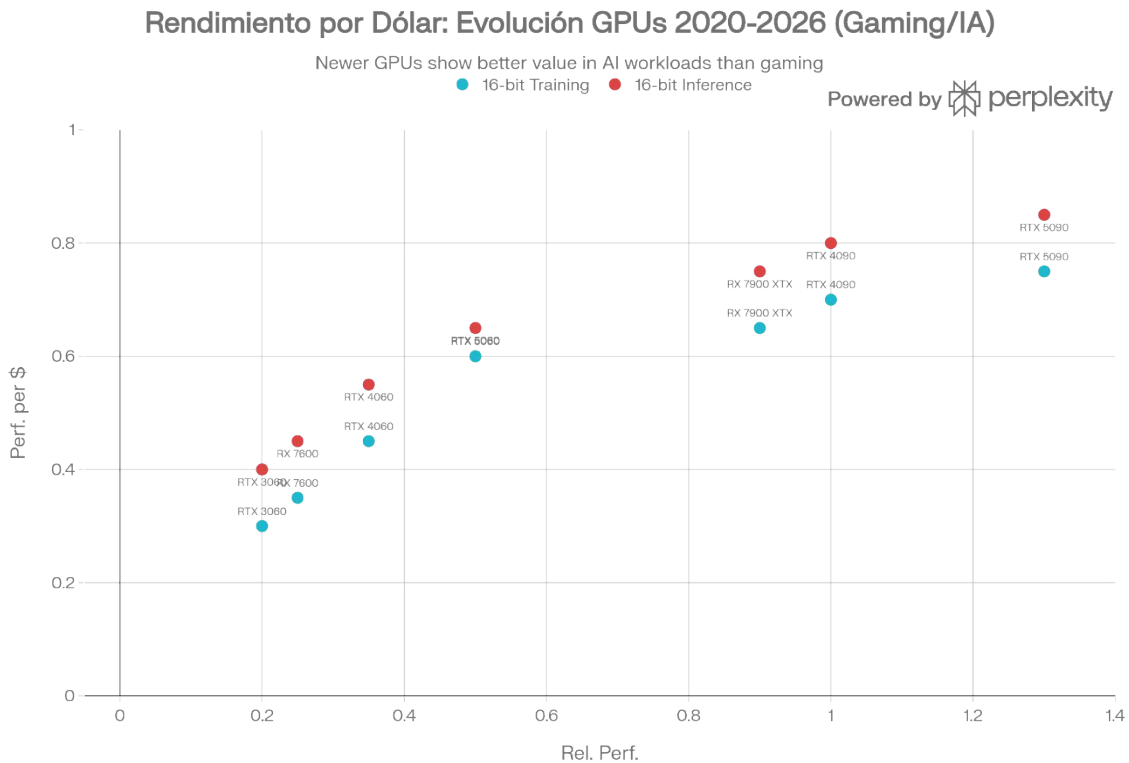
El mercado de GPUs evoluciona rápidamente, con ciclos de renovación cada pocos años. Sin embargo, el aumento del rendimiento no siempre justifica una actualización inmediata. Factores como el uso real, la resolución del monitor y las aplicaciones utilizadas deben pesar más que la novedad tecnológica.

En los últimos años, la demanda de GPUs se ha visto influida por fenómenos como la minería de criptomonedas y el auge de la inteligencia artificial, lo que ha afectado a precios y disponibilidad. Comprender estas dinámicas ayuda a contextualizar el mercado y a tomar decisiones más racionales.

## 9.6 Precio, rendimiento y sentido común

A medida que se asciende en la gama de tarjetas gráficas, el incremento de rendimiento es cada vez menor en relación al aumento de precio. Este fenómeno, conocido como rendimiento marginal decreciente, implica que no siempre la GPU más potente es la opción más racional para la mayoría de usuarios.

Elegir bien  $\neq$  elegir lo más caro:



*\* Nuevas GPUs mejoran valor en IA, pero gaming se satura post-minería*

La vida útil de una tarjeta gráfica no depende únicamente de su potencia inicial, sino también del soporte de software que recibe a lo largo del tiempo. La optimización de drivers, la compatibilidad con nuevas APIs y la implementación de tecnologías como el escalado por IA permiten prolongar la relevancia de una GPU más allá de lo que indicarían sus especificaciones originales. En este contexto, el soporte continuado del fabricante se convierte en un factor tan importante como el hardware en sí.

## 10. Monitores y tecnologías de visualización: donde la GPU cobra sentido

El trabajo de procesamiento masivo realizado por la GPU culmina en un destino crítico: el monitor. Este dispositivo no es un mero periférico de salida, sino el traductor final que convierte flujos de datos binarios en fotones perceptibles por el ojo humano. Factores como la resolución, la profundidad de color, la tasa de refresco y la tecnología de retroiluminación no solo definen la calidad de la imagen, sino que imponen las exigencias de carga de trabajo a la tarjeta gráfica. Sin una visualización precisa, el cálculo de trazado de rayos (Ray Tracing) o el suavizado por IA (DLSS) pierden su propósito, convirtiendo al monitor en el cuello de botella final o en la ventana perfecta hacia el rendimiento del sistema.

En la computación moderna, la GPU y el monitor han dejado de operar de forma independiente para formar un **ecosistema de sincronización activa**. La aparición de tecnologías de refresco variable ha eliminado la barrera entre la generación del fotograma en el búfer de la tarjeta y su representación física en el panel, permitiendo una fluidez absoluta que depende de una comunicación constante a través de interfaces de alto ancho de banda como HDMI 2.1 o DisplayPort. Una elección desequilibrada en este binomio —como emparejar una GPU de gama entusiasta con un panel limitado en hercios o espacio de color— resulta en un desperdicio de recursos y una experiencia degradada.

Este capítulo analiza las arquitecturas de paneles actuales (IPS, VA, OLED) y los parámetros técnicos que definen la fidelidad visual, explorando la relación simbiótica entre la potencia de cálculo y la capacidad de representación. El objetivo es comprender que la calidad de la experiencia no reside solo en cuántos fotogramas puede generar la GPU, sino en la precisión, rapidez y coherencia con la que el monitor es capaz de dibujarlos.

### 10.1 Tecnologías de panel: cómo se crea la imagen

Los monitores actuales se diferencian principalmente por la tecnología de panel, que determina aspectos clave como la reproducción del color, el contraste, los ángulos de visión y la velocidad de respuesta. No existe una tecnología “perfecta”, sino soluciones optimizadas para distintos escenarios de uso.

### ***TN (Twisted Nematic): Velocidad y Transición***

Los paneles TN priorizan la velocidad de respuesta y las altas tasas de refresco. Su principal ventaja es el bajo tiempo de respuesta, lo que históricamente los ha hecho populares en entornos competitivos. Sin embargo, presentan limitaciones importantes en color y ángulos de visión, lo que reduce su fidelidad visual, debido a que el cristal líquido no se orienta de forma uniforme, provocando un "lavado" de color si no se mira frontalmente.

Actualmente están relegados casi exclusivamente a nichos de e-sports de ultra-alta frecuencia (360Hz+) donde la claridad del movimiento es la única prioridad.

### ***IPS (In-Plane Switching): Fidelidad y Consistencia***

Los paneles IPS ofrecen una excelente reproducción del color y amplios ángulos de visión. Son ampliamente utilizados en diseño gráfico, edición de vídeo y uso general. Aunque tradicionalmente eran más lentos que los TN, los IPS modernos han reducido esta diferencia de forma notable.

Su talón de Aquiles es el "IPS Glow", un ligero resplandor visible en las esquinas de escenas oscuras debido a la retroiluminación siempre encendida. Actualmente, son el estándar de oro para profesionales creativos y jugadores que buscan un equilibrio entre velocidad y representación cromática (cobertura sRGB/DCI-P3).

### ***VA (Vertical Alignment): El Dominio del Contraste***

La tecnología VA se sitúa entre TN e IPS, destacando por su alto contraste y negros más profundos. Son adecuados para consumo multimedia y juegos inmersivos, aunque pueden presentar ghosting en determinadas transiciones.

Utilizan una alineación vertical que bloquea la luz de fondo de forma mucho más eficiente cuando están en estado "cerrado". Mientras un panel IPS tiene un contraste típico de 1000:1, un VA alcanza fácilmente 3000:1 o 4000:1.

Suelen sufrir de "black smearing" (estelas oscuras) en escenas rápidas, ya que a los cristales les cuesta más tiempo pasar de negro total a gris oscuro que en otras tecnologías.

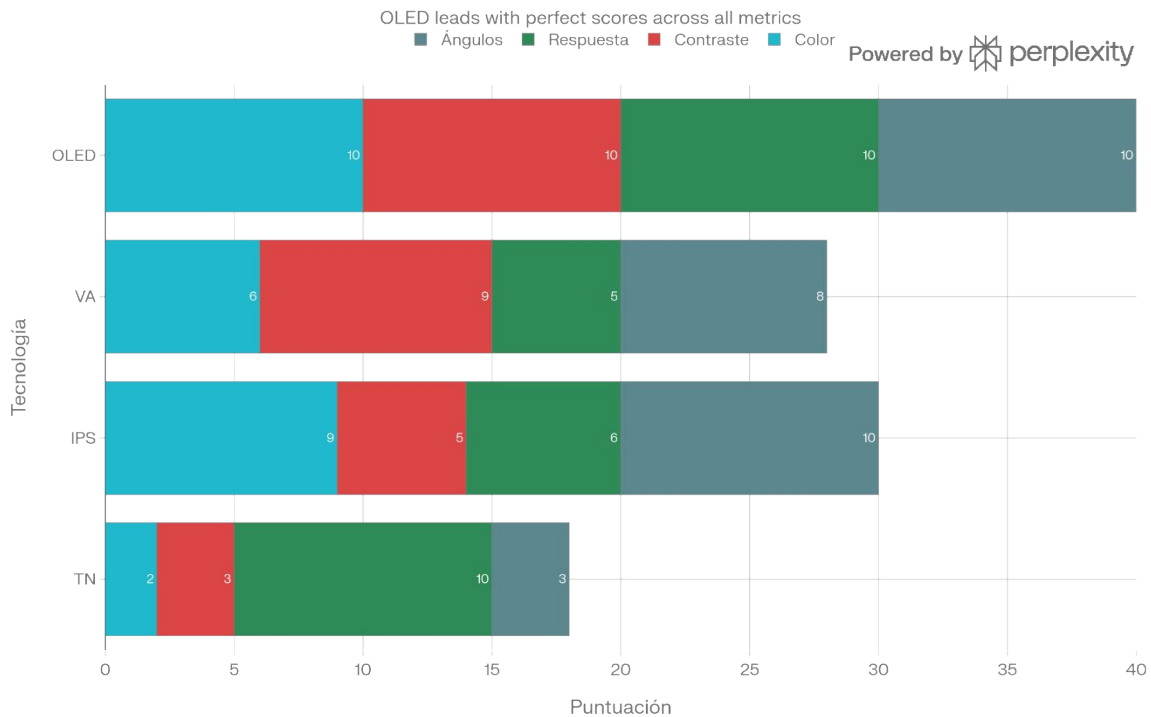
### **OLED (Organic Light Emitting Diode): La Excelencia Autoemisiva**

Los paneles OLED representan la tecnología más avanzada: cada píxel emite su propia luz, lo que permite negros perfectos, contraste prácticamente infinito y tiempos de respuesta extremadamente bajos. Su principal desventaja es el coste y la posible degradación a largo plazo (burn-in).

El contraste infinito se logra apagando físicamente el píxel para representar el negro. Esto elimina el tiempo de respuesta tradicional, situándolo en niveles de 0.03 ms, una cifra órdenes de magnitud superior a cualquier panel LCD.

La gestión térmica es vital para evitar el *burn-in* (marcado estático), lo que requiere que las GPUs modernas utilicen técnicas de desplazamiento de píxeles o limitadores de brillo (ABL).

#### **Comparativa Tecnologías Panel (Puntuación 0-10)**



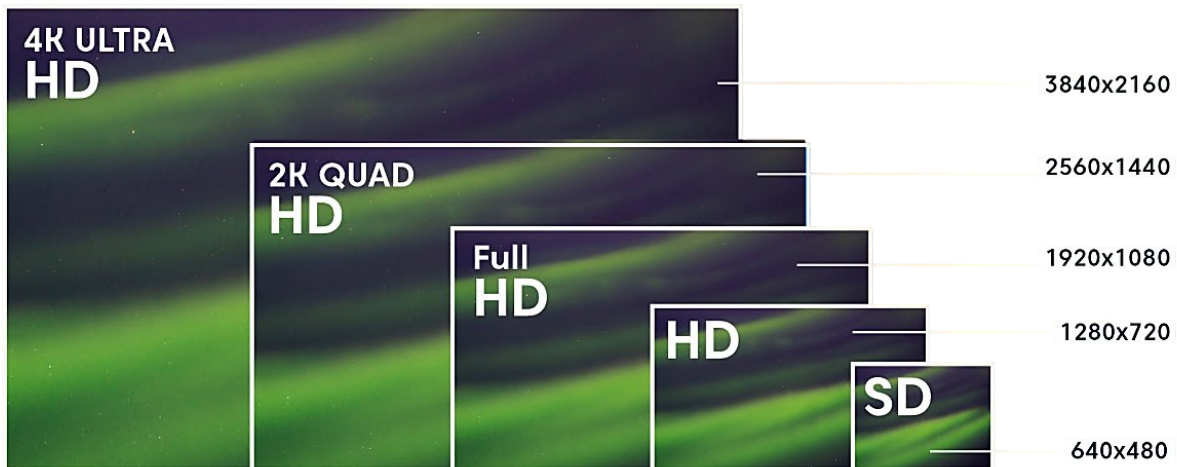
*\* OLED domina total, TN velocidad pura, IPS versátil*

## 10.2 Resolución: cuántos píxeles puede mover tu GPU

La resolución define el número total de píxeles que el monitor puede mostrar simultáneamente. A mayor resolución, mayor nivel de detalle, pero también mayor carga de trabajo para la GPU, que debe calcular y renderizar más información por cada fotograma.

Pasar de 1080p a 1440p o 4K supone un salto significativo tanto en calidad visual como en requisitos gráficos. Por ello, la elección de la resolución debe ir siempre acompañada de una GPU capaz de sostenerla de forma fluida.

### Resolution & Pixels



*\* el salto a 4K u 8K no es incremental, es exponencial.*

Tamaño (")	Aspecto	Resolución	Píxeles Totales	PPI	Dist. Visión Mín (cm)
<b>24</b>	16:9	1920x1080	2.07M	92	94 <a href="#">displayninja</a>
<b>24</b>	16:9	2560x1440	3.69M	122	71
<b>27</b>	16:9	1920x1080	2.07M	82	107
<b>27</b>	16:9	2560x1440	3.69M	109	81
<b>27</b>	16:9	3840x2160	8.29M	163	53
<b>32</b>	16:9	2560x1440	3.69M	93	94
<b>32</b>	16:9	3840x2160	8.29M	140	64
<b>34</b>	21:9 UW	3440x1440	4.95M	110	79 <a href="#">rtings</a>
<b>34</b>	21:9 UW	5120x2160	11.06M	163	48
<b>43</b>	32:9 UW	3840x1080	4.15M	92	94
<b>47</b>	32:9 UW	5120x1440	7.37M	109	81
<b>49</b>	32:9 Super	5120x1440	7.37M	109	81 <a href="#">pcmag</a>
<b>49</b>	32:9 Super	7680x2160	16.59M	163	53

### 10.3 Tasa de refresco y tiempo de respuesta: fluidez percibida

La tasa de refresco, medida en hercios (Hz), indica cuántas veces por segundo el monitor puede actualizar la imagen. Un monitor de 144 Hz puede mostrar hasta 144 fotogramas por segundo, siempre que la GPU sea capaz de generarlos.

El tiempo de respuesta, medido en milisegundos (ms), representa el tiempo que tarda un píxel en cambiar de estado. Valores bajos reducen el desenfoque de movimiento y mejoran la nitidez en escenas rápidas. Ambos parámetros influyen directamente en la sensación de fluidez, especialmente en videojuegos y simulaciones.

Tecnologías como NVIDIA G-Sync o AMD FreeSync eliminan el 'tearing' (corte de imagen) sincronizando los Hz del monitor con los FPS que produce la GPU en tiempo real.

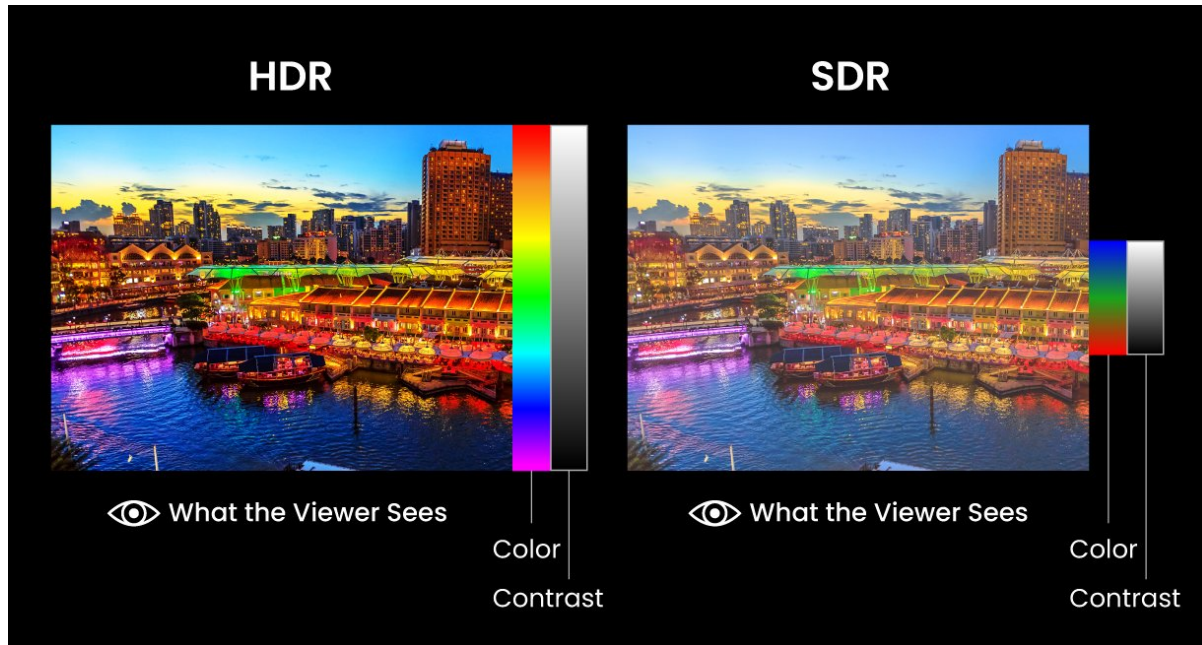


Frecuencia (Hz)	Fluidez percibida	Uso típico recomendado
60 Hz	Baja	Ofimática, vídeo estándar, consolas antiguas
75 Hz	Media	Uso general, gaming casual en PC
144 Hz	Alta	Gaming competitivo 1080p/1440p
240 Hz	Muy alta	eSports, shooters rápidos, alta tasa de FPS

#### 10.4 Brillo, contraste y calidad de imagen

El brillo, medido en nits ( $\text{cd}/\text{m}^2$ ), no solo garantiza la visibilidad en habitaciones muy iluminadas, sino que es el motor fundamental para el HDR (Alto Rango Dinámico); mientras que un monitor estándar ronda los 250-300 nits, para disfrutar de impactos de luz realistas y certificaciones como DisplayHDR se recomiendan paneles que superen los 400 o 600 nits. Por otro lado, el contraste dicta la profundidad de la imagen mediante la relación entre el blanco más brillante y el negro más oscuro. Aquí, la tecnología del panel es clave: mientras los paneles IPS ofrecen un contraste típico de 1000:1, los paneles VA o OLED logran negros mucho más profundos (o infinitos en el caso de OLED), evitando que las escenas oscuras se vean grisáceas o "lavadas".

En el uso diario, la combinación de un buen contraste y un brillo alto suele tener un impacto visual más notable que el simple aumento de resolución; una pantalla 1440p con excelente rango dinámico y cobertura de color amplia (como DCI-P3 para cine y juegos) se verá visualmente superior a un monitor 4K con colores planos. Para creadores de contenido y gamers que buscan inmersión, es crucial priorizar estos valores, ya que revelan detalles ocultos en las sombras y matices en las luces altas que definen la atmósfera de la imagen, elevando la experiencia multimedia muy por encima de la simple nitidez de los píxeles.



Parámetro	SDR (Estándar)	HDR (Dolby Vision/HDR10+)	Mejora Principal
<b>Brillo pico</b>	100-400 nits	400-4000 nits (HDR2000)	Highlights realistas (sol) <a href="#">newegg</a>
<b>Contraste dinámico</b>	1000:1	1.000.000:1+	Negros absolutos
<b>Rango dinámico</b>	6-8 stops	14-20 stops	Escenas día/noche fluidas
<b>Profundidad color</b>	8-10 bit	10-12 bit	Sin banding, gradientes perfectos
<b>Espacio de color</b>	sRGB (~70% DCI-P3)	98%+ DCI-P3/Rec.2020	Colores Hollywood vívidos
<b>Local Dimming</b>	No	Sí (hasta 1000+ zonas)	Contraste preciso por escena
<b>Metadatos</b>	Estáticos	Dinámicos escena-por-escena	Optimización contenido
<b>Gaming específico</b>	Básico	RT + AutoHDR/DLSS	Reflejos/global illum real-time

El HDR revoluciona la experiencia visual al ofrecer un realismo inigualable mediante picos de brillo de hasta 4000 nits y un contraste dinámico de 1.000.000:1, superando ampliamente las limitaciones del estándar SDR. Al utilizar una profundidad de color de 10-12 bits y cubrir más del 98% del espacio DCI-P3, logra gradientes perfectos y colores vibrantes de nivel cinematográfico sin el molesto efecto de banding. Gracias a los metadatos dinámicos y el Local Dimming, la imagen se optimiza escena por escena, garantizando negros absolutos y una iluminación global precisa que transforma tanto el cine como el gaming de nueva generación.

### **10.5 Sincronización adaptativa: GPU y monitor trabajando juntos**

Cuando la GPU y el monitor no están sincronizados, pueden aparecer defectos visuales como el tearing. Las tecnologías de sincronización adaptativa, como FreeSync (AMD) y G-Sync (NVIDIA), ajustan dinámicamente la tasa de refresco del monitor al número de fotogramas generados por la GPU.

Estas tecnologías mejoran la fluidez, reducen el stuttering y permiten una experiencia más estable incluso cuando los FPS varían.

#### ***El conflicto entre FPS y Hercios (Hz)***

El problema de comunicación surge porque, tradicionalmente, los monitores funcionan a una tasa de refresco fija (por ejemplo, 60Hz), mientras que la GPU genera fotogramas a una velocidad variable dependiendo de la complejidad de la escena. Cuando la GPU envía un nuevo fotograma antes de que el monitor termine de dibujar el anterior, ocurre el tearing (una ruptura horizontal en la imagen). Por el contrario, si la GPU se retrasa, el monitor repite el fotograma anterior, causando stuttering o pequeños tirones que rompen la sensación de fluidez. La sincronización adaptativa elimina esta rigidez, permitiendo que el monitor "espere" activamente a la GPU antes de refrescar la pantalla.

#### ***Ecosistemas de Sincronización: G-Sync, FreeSync y VRR***

Aunque el objetivo es el mismo, la implementación varía según el fabricante. NVIDIA G-Sync utiliza un módulo de hardware propietario dentro del monitor para asegurar una

latencia mínima y un rango de operación más amplio. Por otro lado, AMD FreeSync aprovecha el estándar abierto DisplayPort Adaptive-Sync, lo que permite que sea más accesible y económico al no requerir hardware dedicado. Además, ha surgido el HDMI 2.1 VRR (Variable Refresh Rate), una tecnología estándar adoptada por consolas de nueva generación y televisores modernos. Estas soluciones no solo eliminan artefactos visuales, sino que también reducen el input lag en comparación con el antiguo V-Sync (Sincronización Vertical), ofreciendo una respuesta inmediata a los comandos del usuario.

Tecnología	Fabricante	Requisito Hardware	Ventaja Principal	Compatibilidad
<b>V-Sync</b>	Genérico	Software (driver)	Elimina tearing pero +lag input	Todos <a href="#">newegg</a>
<b>G-Sync</b>	NVIDIA	Chip dedicado monitor	Máxima precisión, control overdrive, <1ms lag	Monitores certificados
<b>FreeSync</b>	AMD	Estándar Adaptive-Sync	Compatible amplia, económico	90% monitores modernos
<b>HDMI VRR</b>	Estándar	HDMI 2.1	Ideal consolas PS5/Xbox/TV	TVs gaming HDMI 2.1



### 10.6 La sinergia GPU–Monitor: elegir bien el conjunto

La relación entre la GPU y el monitor es una simbiosis donde el rendimiento de una define el límite del otro. Instalar una tarjeta de gama entusiasta, como una RTX 4090, en un monitor 1080p a 60Hz crea un "cuello de botella visual" severo: la GPU podría estar generando 300 fotogramas por segundo, pero el usuario solo verá 60, desperdiciando el 80% del potencial del hardware. Por el contrario, intentar mover un monitor 4K a 144Hz con una GPU de gama entrada resultará en una experiencia frustrante, con tasas de cuadros por segundo tan bajas que harán que el movimiento se vea entrecortado, obligando al usuario a reducir drásticamente la calidad gráfica y perdiendo la nitidez por la que pagó al comprar el monitor.

Para lograr una sinergia perfecta, el usuario debe alinear la resolución del panel con la potencia de cómputo y la tecnología de refresco variable (VRR). En el gaming moderno, el objetivo ideal es que la GPU sea capaz de mantener una tasa de FPS que oscile dentro del rango de hercios del monitor; por ejemplo, una configuración equilibrada para la gama media actual sería un panel 1440p (QHD) a 144Hz combinado con una GPU de la serie RTX 4070 o RX

7800. Esta combinación permite aprovechar la mayor densidad de píxeles sin sacrificar la fluidez, mientras que tecnologías como DLSS o FSR actúan como puentes para mantener esa sinergia en títulos exigentes, asegurando que la inversión en ambos componentes se traduzca en una mejora visual tangible y constante.

Perfil Usuario	Resolución Recomendada	Frecuencia (Hz)	GPU	Panel Sugerido	Presupuesto Aprox.
<b>Competitivo</b>	1080p FHD (1920x1080)	240Hz+	RTX 4060 Ti / RX 7600	TN/OLED	€400
<b>Gaming</b>	1440p QHD (2560x1440)	144-165Hz	RTX 4070 / RX 7800 XT	IPS/VA	€700
<b>Entusiasta/Pro</b>	4K UHD (3840x2160)	120-144Hz	RTX 4080 / 4090	OLED/IPS MiniLED	€1500+

## 11. Conclusiones y perspectivas de futuro

El estudio de las tarjetas gráficas y los monitores permite comprender que el sistema gráfico de un ordenador moderno es el resultado de la interacción coordinada entre múltiples componentes especializados. La GPU ha evolucionado desde un dispositivo dedicado exclusivamente al procesamiento de gráficos hasta convertirse en un acelerador de propósito general, capaz de afrontar tareas complejas como la inteligencia artificial, el renderizado profesional o la simulación científica gracias a su arquitectura altamente paralela.

A lo largo del proyecto se ha analizado cómo factores como la arquitectura interna de la GPU, el paralelismo masivo, la VRAM, el ancho de banda, la alimentación y la refrigeración influyen directamente en el rendimiento real del sistema. Asimismo, se ha puesto de manifiesto que el rendimiento no depende únicamente de la potencia bruta, sino del equilibrio global del sistema, evitando cuellos de botella entre CPU, GPU, memoria y monitor.

En este contexto, la elección de una tarjeta gráfica no debe basarse exclusivamente en cifras comerciales como la cantidad de VRAM o el nombre del modelo, sino en un análisis racional del uso previsto, la resolución objetivo y las características del monitor. Del mismo modo, el monitor deja de ser un elemento secundario para convertirse en una pieza clave que condiciona la experiencia visual final, siendo imprescindible comprender parámetros como la tasa de refresco, la resolución, el tipo de panel o la sincronización adaptativa.

Mirando hacia el futuro, todo apunta a que la GPU continuará ganando protagonismo como pilar fundamental de la computación moderna. El crecimiento de la inteligencia artificial, el aprendizaje automático y los modelos de lenguaje de gran tamaño refuerzan el papel de la GPU como motor de cálculo paralelo. Paralelamente, tecnologías como el ray tracing, el escalado mediante IA y los paneles de nueva generación seguirán elevando el nivel de realismo y calidad visual.

En conclusión, comprender cómo funcionan las tarjetas gráficas y los monitores no solo permite tomar mejores decisiones de compra, sino también entender uno de los pilares tecnológicos que sustentan la sociedad digital actual. La GPU ya no es “solo para jugar”, sino una herramienta clave en la evolución de la informática moderna.

## 12. ANEXO: HORIZONTES EXPANDIDOS (La Era GPGPU)

Si bien a lo largo de este proyecto se ha analizado la GPU como el motor indiscutible de la representación gráfica, su evolución reciente ha trascendido la barrera del píxel. Hoy en día, la tarjeta gráfica se ha consolidado como el cerebro de la computación paralela de propósito general, un fenómeno conocido técnicamente como **GPGPU** (*General-Purpose Computing on Graphics Processing Units*).

Este anexo explora cronológicamente cómo el hardware diseñado para videojuegos acabó resolviendo los problemas más complejos de la ciencia, la economía y la medicina moderna.

### 12.1 Principios de los 2000: La "Alquimia Digital" y los Inicios

A comienzos del nuevo milenio, mucho antes de que existieran estándares como CUDA, un grupo de ingenieros pioneros, entre ellos Mark Harris, se enfrentó a un hardware que estaba "bloqueado" para entender solo gráficos. En esta etapa, conocida como la era de la "alquimia digital", los investigadores tuvieron que recurrir al ingenio para realizar cálculos matemáticos complejos: disfrazaban los datos numéricos como si fueran colores.

El proceso consistía en traducir matrices matemáticas a valores RGB, engañando a la GPU para que creyera que estaba procesando texturas de un videojuego cuando, en realidad, estaba resolviendo ecuaciones físicas. El resultado visual —un píxel de un color específico— se traducía de vuelta a un número, sentando así las bases conceptuales para que la tarjeta gráfica dejara de ser un juguete y se convirtiera en una herramienta científica.

### 12.2 2006 – 2020: De Folding@home a la Lucha Global contra Pandemias

El potencial latente de las GPU se materializó masivamente en 2006 con la actualización del proyecto Folding@home de la Universidad de Stanford. Este proyecto buscaba simular el plegamiento tridimensional de las proteínas, un proceso biológico crítico cuyo fallo provoca enfermedades como el Alzheimer o el cáncer, y que requiere una potencia de cálculo astronómica.

La verdadera demostración de fuerza llegó años más tarde, durante la crisis sanitaria de 2020. Miles de usuarios domésticos cedieron la potencia de sus tarjetas gráficas para investigar la estructura del virus SARS-CoV-2. Gracias a la arquitectura paralela de las GPU de consumo, la red combinada alcanzó los 1.5 exaflops, convirtiéndose temporalmente en la supercomputadora más rápida del planeta y superando la capacidad combinada de las máquinas de la NASA e IBM.

### **12.3 2010: El Cluster Condor y la Supercomputación con Consolas**

En un fascinante giro de eficiencia económica, la Fuerza Aérea de los Estados Unidos demostró en 2010 que el hardware de videojuegos podía tener aplicaciones de defensa crítica. Ante la necesidad de procesar imágenes satelitales de alta resolución en tiempo real, y frente al coste prohibitivo de los servidores tradicionales, la USAF construyó el "Cluster Condor" conectando en red 1.760 consolas PlayStation 3.

El uso del chip *Cell* de la consola, que compartía principios de diseño vectorial con las GPU modernas, permitió construir este superordenador por solo 2 millones de dólares, una fracción del coste habitual de 10 millones. Este hito validó definitivamente que el hardware de consumo masivo, optimizado para la física y los gráficos, ofrecía una relación coste-rendimiento inigualable para tareas de vigilancia y procesamiento de datos.

### **12.4 2010 – Actualidad: La Revolución del Criptoanálisis**

Paralelamente, a partir de 2010, la economía digital encontró en las GPU su motor de minería ideal. Algoritmos de *hash* como SHA-256 (Bitcoin) o Ethash (Ethereum) requerían resolver problemas matemáticos repetitivos millones de veces por segundo. La arquitectura SIMT de las GPU, diseñada para ejecutar la misma instrucción sobre múltiples datos simultáneamente, resultó ser miles de veces más eficiente que cualquier CPU convencional para esta tarea.

Este descubrimiento transformó el mercado para siempre, provocando crisis de stock globales en 2017 y 2021. Por primera vez en la historia, el valor y la disponibilidad de las

tarjetas gráficas se desvincularon de su rendimiento en videojuegos para ligarse a la volatilidad de los activos financieros, demostrando la versatilidad bruta del silicio gráfico.

### **12.5 2012: El "Momento AlexNet" y el Renacimiento de la IA**

La inteligencia artificial moderna tiene una fecha de nacimiento clara vinculada al hardware. En la competición *ImageNet* de 2012, el equipo de Alex Krizhevsky presentó AlexNet, una red neuronal que aplastó a la competencia reduciendo la tasa de error casi a la mitad. Lo revolucionario no fue solo el algoritmo, sino el método: utilizaron dos tarjetas gráficas NVIDIA GTX 580 para el entrenamiento, programando manualmente las operaciones matemáticas.

Este evento marcó un punto de inflexión donde la comunidad científica comprendió que la IA no estaba limitada por el software, sino por la falta de paralelismo en el cálculo. Desde ese momento, el desarrollo de las GPU y la IA se fusionaron, llevando a los fabricantes a integrar núcleos específicos (Tensor Cores) en sus diseños para acelerar las matrices que hoy dan vida a tecnologías como ChatGPT.

### **12.6 Presente y Futuro: Del Renderizado al Diagnóstico Médico**

Hoy en día, la GPU ha cerrado el círculo en campos como la medicina, pasando de ser una herramienta de visualización a una de análisis. Si antes su función se limitaba a renderizar los datos de una resonancia magnética para mostrarlos en pantalla, ahora utiliza el *Deep Learning* para "entender" esas imágenes.

Mediante el entrenamiento con millones de casos previos, las GPU actuales pueden detectar microfracturas, tumores incipientes o anomalías vasculares invisibles al ojo humano en cuestión de segundos. Esta capacidad de convertir datos crudos en diagnósticos precisos representa la frontera final de la tecnología gráfica, donde la velocidad de procesamiento ya no solo sirve para ganar fluidez en un videojuego, sino para salvar vidas mediante la medicina de precisión.

### 13. Bibliografía y fuentes consultadas

NVIDIA Corporation. (2023). *NVIDIA CUDA C Programming Guide. NVIDIA Developer Documentation.*

NVIDIA Corporation. (2024). *GPU Architecture and Ray Tracing Overview. NVIDIA Technical Whitepapers.*

AMD. (2023). *AMD RDNA Architecture Whitepaper. Advanced Micro Devices.*

Intel Corporation. (2023). *Intel Arc Graphics Architecture. Intel Developer Zone.*

PCI-SIG. (2022). *PCI Express Base Specification. PCI Special Interest Group.*

JEDEC Solid State Technology Association. (2021). *GDDR6 and GDDR6X Memory Standards.*

Khronos Group. (2023). *Vulkan Graphics API Overview. Khronos Documentation.*

Akenine-Möller, T., Haines, E., & Hoffman, N. (2018). *Real-Time Rendering (4th ed.). CRC Press.*

AnandTech. (2023). *GPU Benchmarks and Architecture Analysis. AnandTech Hardware Reviews.*

TechPowerUp. (2024). *GPU Database and Performance Analysis. TechPowerUp.*

NVIDIA Corporation. (2022). *NVIDIA Ada Lovelace Architecture Whitepaper.*

AMD. (2023). *RDNA 3 Instruction Set Architecture Reference Guide.*

Hennessy, J. L., & Patterson, D. A. (2019). *Computer Architecture: A Quantitative Approach.*

Jon Peddie Research. (2024). *GPU Market Share Report Q4.*