



### INTRODUCCIÓN AL PARALELISMO

#### 1. DEFINICIÓN DE COMPUTACIÓN PARALELA

La computación paralela es un paradigma que busca mejorar el rendimiento y reducir los tiempos de cálculo mediante la ejecución simultánea de múltiples tareas. Este enfoque es fundamental en la actualidad debido al crecimiento exponencial de los datos y la necesidad de procesarlos de manera eficiente.

En términos generales, la computación paralela puede verse como un conjunto de técnicas y metodologías que permiten distribuir la carga de trabajo entre múltiples unidades de procesamiento, aprovechando la capacidad de cómputo de hardware especializado, como procesadores multinúcleo, GPUs y arquitecturas distribuidas.

##### 1.1 Computación Secuencial vs. Computación Paralela

- **Computación Secuencial:** Un único procesador ejecuta una serie de instrucciones en un orden predefinido, sin posibilidad de ejecución simultánea.
- **Computación Paralela:** Varias unidades de procesamiento trabajan en conjunto para ejecutar múltiples instrucciones de manera simultánea, reduciendo el tiempo total de ejecución de un programa.

Ventajas y Desventajas

Característica	Computación Secuencial	Computación Paralela
Rendimiento	Bajo en tareas intensivas	Alto en tareas paralelizables
Complejidad	Baja	Alta
Consumo energético	Moderado	Puede ser alto dependiendo de la implementación
Aplicabilidad	Problemas lineales	Problemas escalables y distribuidos

##### 1.2 Motivación del Paralelismo

El uso de computación paralela está motivado por varios factores, entre ellos:

- **Demanda de procesamiento masivo:** En aplicaciones como inteligencia artificial, procesamiento de imágenes y simulaciones científicas, la cantidad de datos a procesar es demasiado grande para un enfoque secuencial.



- **Estancamiento en la mejora del hardware secuencial:** El crecimiento en la velocidad de los procesadores se ha desacelerado, lo que ha llevado a un aumento en la cantidad de núcleos en lugar de una mayor frecuencia de reloj.
- **Eficiencia energética:** La ejecución de código en móviles introduce un problema: la tecnología de baterías no evoluciona tan rápidamente como la de los procesadores. Aunque la computación paralela puede consumir más energía, en muchos casos optimiza el uso del hardware y permite una distribución más eficiente de los recursos.

### 1.3 Aplicaciones Actuales del Paralelismo

- **Inteligencia Artificial (IA):** La aceleración mediante GPUs permite entrenar redes neuronales profundas en menos tiempo. Por ejemplo, en el entrenamiento de redes neuronales convolucionales (CNNs), cada imagen de entrenamiento puede dividirse en múltiples fragmentos procesados en paralelo.
- **Big Data y Análisis de Datos:** Herramientas como Apache Spark permiten el procesamiento distribuido de grandes volúmenes de información. Un ejemplo es la simulación de modelos climáticos, donde se dividen regiones geográficas entre procesadores.
- **Simulaciones Científicas:** Se aplican en meteorología, modelado de proteínas y otras disciplinas que requieren cálculos masivos.
- **Renderizado de Videojuegos:** Permite gráficos más realistas y física avanzada en entornos interactivos.

---

## 2. TIPOS DE PARALELISMO

Existen diversos tipos de paralelismo, cada uno con aplicaciones específicas dependiendo del problema a resolver.

### 2.1 Paralelismo a Nivel de Datos

Este enfoque permite ejecutar la misma operación sobre diferentes datos de manera simultánea. Se emplea en:

- **Procesamiento Vectorial y SIMD** (Single Instruction Multiple Data): En arquitecturas SIMD, una única instrucción se ejecuta en múltiples elementos de datos en paralelo. Este enfoque es altamente eficiente en aplicaciones como procesamiento de imágenes, multimedia y computación científica.
- **Uso de GPUs para cálculos en paralelo:** Las GPUs están diseñadas para manejar miles de hilos en paralelo, lo que las hace ideales para gráficos computacionales y modelos de aprendizaje profundo.
- **Aplicaciones en Aprendizaje Automático y Gráficos Computacionales:** Muchas técnicas de inteligencia artificial, como redes neuronales convolucionales (CNNs), aprovechan la ejecución paralela en GPUs para reducir los tiempos de entrenamiento.



**Ejemplo práctico:** En procesamiento de imágenes, un filtro puede aplicarse a diferentes píxeles de forma simultánea, mejorando significativamente el rendimiento respecto a una ejecución secuencial.

## 2.2 Paralelismo a Nivel de Instrucciones (Resumen)

Se basa en la ejecución simultánea de múltiples instrucciones dentro de una CPU. Se implementa mediante:

- **Pipelining:** Divide la ejecución de instrucciones en múltiples etapas (como búsqueda, decodificación y ejecución).
- **Ejecución Fuera de Orden:** Permite a las CPUs reordenar la ejecución de instrucciones para optimizar el uso de las unidades funcionales.
- **Superescalaridad:** Uso de múltiples unidades de ejecución dentro de un procesador para ejecutar varias instrucciones en un solo ciclo de reloj.

**Nota:** La predicción de ramas y la ejecución especulativa se explicarán en sesiones posteriores cuando se aborden optimizaciones avanzadas en arquitecturas de CPU.

## 2.3 Paralelismo a Nivel de Tareas

Este tipo de paralelismo divide un problema en múltiples tareas independientes que pueden ejecutarse en paralelo. Se emplea en:

- **Computación Distribuida en Clústeres y Supercomputadoras:** En estos sistemas, múltiples nodos trabajan en paralelo para resolver problemas de alto cómputo. Ejemplos incluyen clústeres de alto rendimiento utilizados en simulaciones científicas y modelado climático.
- **Hilos y Procesos en Sistemas Operativos Modernos:** Muchas aplicaciones modernas utilizan hilos y procesos para mejorar la eficiencia. Ejemplo: Un navegador web puede cargar múltiples pestañas en paralelo.
- **Simulaciones Físicas:** En entornos de simulación, cada entidad física puede ser modelada y procesada en paralelo. Ejemplo: En videojuegos, múltiples objetos pueden moverse simultáneamente sin dependencia entre ellos.
- **Tareas Asíncronas y Computación en la Nube:** En entornos distribuidos, servicios en la nube procesan peticiones en paralelo para mejorar el tiempo de respuesta.

**Mención especial:** En sesiones posteriores se detallarán los **modelos de computación paralela** basados en **memoria compartida, memoria distribuida y modelos híbridos**.

---

## 3. LIMITACIONES Y DESAFÍOS DEL PARALELISMO

Aunque la computación paralela ofrece grandes ventajas, presenta desafíos que deben ser abordados para maximizar su eficiencia.



## 3.1 Ley de Amdahl y Limitaciones en la Aceleración

La Ley de Amdahl establece que el beneficio del paralelismo está limitado por la fracción del programa que no puede ejecutarse en paralelo. Su ecuación se expresa como:

$$S(N) = \frac{1}{(1 - P) + \frac{P}{N}}$$

Donde:

- **S(N)** es la aceleración obtenida con **N** procesadores.
- **P** es la fracción del código que puede ejecutarse en paralelo.
- **N** es el número de unidades de procesamiento.

### Ejemplo 1: Código con Fracción Secuencial Significativa

Supongamos que un programa tiene un 70% de su código paralelizable ( $P = 0.7$ ) y el 30% restante debe ejecutarse de forma secuencial.

- Si usamos 4 procesadores ( $N = 4$ ), aplicamos la ecuación:

$$S(4) = \frac{1}{(1 - 0.7) + \frac{0.7}{4}}$$
$$S(4) = \frac{1}{0.3 + 0.175} = \frac{1}{0.475} \approx 2.11$$

Esto significa que, **a pesar de cuadruplicar los procesadores**, la aceleración solo es **2.11 veces mayor**, muy por debajo del ideal de **4x**.

### Ejemplo 2: Código Altamente Paralelizable

Ahora supongamos que **95% del código es paralelizable** ( $P = 0.95$ ) y usamos **16 procesadores** ( $N = 16$ ):

$$S(16) = \frac{1}{(1 - 0.95) + \frac{0.95}{16}}$$
$$S(16) = \frac{1}{0.05 + 0.059} = \frac{1}{0.109} \approx 9.17$$



A pesar de usar **16 procesadores**, la **aceleración real solo es 9.17 veces mayor**, mostrando que aún con alto paralelismo, hay límites en la mejora del rendimiento.

### Soluciones a las Limitaciones de Amdahl

- **Incrementar P:** Reescribir código para **hacer más partes paralelizables**.
- **Reducir la fracción secuencial:** **Optimización** de algoritmos secuenciales.
- **Ley de Gustafson:** **En problemas de mayor escala**, la fracción paralelizable tiende a aumentar, mitigando el impacto de Amdahl.

### 3.2 Ley de Gustafson: Un Enfoque Alternativo

La Ley de Gustafson sugiere que el paralelismo es más efectivo en problemas escalables. Conforme aumenta el tamaño del problema, la fracción paralelizable también crece. Su ecuación se expresa como:

$$S(N) = N - \alpha (N - 1)$$

Donde:

- **S(N)** es la aceleración obtenida con **N** procesadores.
- **$\alpha$**  es la fracción secuencial del código.
- **N** es el número de procesadores.

#### Ejemplo de la Ley de Gustafson

Supongamos que un problema tiene una fracción secuencial del **5%** ( **$\alpha = 0.05$** ). Si usamos **16 procesadores (N = 16)**:

$$S(16) = 16 - 0.05 (16 - 1)$$

$$S(16) = 16 - 0.05 (15)$$

$$S(16) = 16 - 0.75$$

$$S(16) = 15.25$$

Este resultado indica que, en lugar de la aceleración de **9.17 veces** obtenida con la Ley de Amdahl, la Ley de Gustafson sugiere una aceleración de **15.25 veces**, lo que refleja un mejor aprovechamiento del paralelismo en problemas que pueden escalar con la cantidad de procesadores.

#### Diferencias entre Amdahl y Gustafson

- **Amdahl** asume un problema de tamaño fijo y muestra que hay un límite en la aceleración que se puede lograr.
- **Gustafson** permite que el problema escale, argumentando que conforme aumenta el número de procesadores, podemos abordar problemas más grandes y complejos.



- **Conclusión:** En aplicaciones prácticas, la **Ley de Gustafson** es más relevante para cargas de trabajo reales en supercomputación y computación de alto rendimiento.

### Comparación práctica: Amdahl vs. Gustafson

Situación práctica	¿Qué ley usar?	Explicación
Optimización en sistemas pequeños	<b>Ley de Amdahl</b>	El problema es fijo y necesitas evaluar si agregar más procesadores mejorará significativamente el rendimiento.
Desarrollo de simulaciones escalables	<b>Ley de Gustafson</b>	El problema puede crecer (más datos, más complejidad) y necesitas aprovechar más procesadores.
Diagnóstico del rendimiento actual	<b>Ley de Amdahl</b>	Para identificar cuánta parte del programa no es paralelizable y cómo limita la aceleración.
Diseño de nuevas aplicaciones	<b>Ley de Gustafson</b>	Cuando planificas un sistema que puede aprovechar recursos adicionales en el futuro.

En resumen, **Amdahl nos indica si vale la pena optimizar el código existente**, mientras que **Gustafson nos muestra cómo escalar la aplicación para aprovechar más recursos**.

### 3.3 Costes de Comunicación y Sincronización

En los sistemas paralelos, la comunicación y sincronización entre procesadores puede convertirse en un cuello de botella significativo. La eficiencia del sistema puede degradarse debido a los siguientes factores:

#### 3.3.1 Latencia de comunicación

- La latencia es el tiempo necesario para que un mensaje viaje desde un nodo emisor hasta un nodo receptor. Este retraso puede ser causado por la distancia física entre nodos, la tecnología de red o la congestión de la misma.
  - **Ejemplo práctico:** En un clúster de supercomputación, si la latencia es alta, los procesadores podrían estar ociosos esperando datos de otros nodos antes de continuar su ejecución.

#### 3.3.2 Ancho de banda limitado

- Aunque tengas múltiples procesadores trabajando en paralelo, la capacidad de transferencia de datos entre ellos puede estar limitada por el ancho de banda. Un ancho de banda insuficiente puede crear cuellos de botella.
  - **Ejemplo práctico:** Una aplicación de análisis de big data que transfiera grandes volúmenes de datos entre nodos puede experimentar tiempos de espera si el ancho de banda no es adecuado.





### 3.3.3 Costes de acceso a memoria compartida

- En sistemas de memoria compartida, varios procesadores pueden acceder simultáneamente a una misma región de memoria. Este acceso concurrente puede provocar conflictos y esperas.
  - **Ejemplo práctico:** Un sistema multinúcleo donde varios núcleos intentan modificar la misma variable compartida puede experimentar degradación del rendimiento.

### 3.3.4 Costes de sincronización

- La sincronización es necesaria cuando varios hilos o procesos deben coordinarse para completar una tarea. Si esta sincronización no está bien gestionada, puede haber tiempos muertos.
  - **Ejemplo práctico:** En simulaciones físicas, cada paso temporal puede requerir que todos los procesadores esperen hasta que todos hayan terminado su parte del cálculo antes de avanzar.

### 3.3.5 Estrategias para mitigar los costes

- **Agrupación de cálculos:** Reducir la frecuencia de comunicación agrupando más cálculos en cada paso para minimizar los tiempos de espera asociados a la transferencia de datos entre nodos o procesadores.
- **Prefetching de datos:** Traer datos anticipadamente a la caché del procesador para reducir el tiempo de espera asociado al acceso a la memoria. Esta técnica anticipa qué datos necesitará el procesador y los coloca cerca del mismo para su acceso rápido.
- **Optimización del ancho de banda:** Comprimir o reducir los datos transferidos, o diseñar algoritmos más eficientes que minimicen la cantidad de datos enviados, evitando cuellos de botella en la red.
- **Minimizar barreras de sincronización:** Evitar puntos de sincronización innecesarios y utilizar enfoques asíncronos que permitan a los procesadores seguir trabajando mientras otros procesadores completan sus tareas.
- **Particionamiento eficiente de la memoria:** Diseñar estructuras de datos y algoritmos de forma que se minimice la compartición de datos entre procesadores, lo que reduce la sobrecarga de coherencia de caché y evita conflictos de acceso a memoria compartida.



## EJERCICIOS DE REFLEXIÓN

### EJERCICIO PROPUESTO 1: COMPARACIÓN DE LEY DE AMDAHL Y LEY DE GUSTAFSON

#### Escenario:

Una empresa tecnológica está desarrollando una aplicación para procesar grandes volúmenes de datos. Actualmente, el procesamiento es secuencial, pero el equipo ha decidido optimizarlo para ejecutarlo en un sistema paralelo con 16 procesadores.

Se han realizado análisis del programa y se ha determinado lo siguiente:

- **El 70% del programa puede paralelizarse ( $P = 0.7$ ).**
- **El 30% del programa debe ejecutarse secuencialmente.**
- **El equipo está considerando dos escenarios de crecimiento:**
  - Escenario 1: Mantener el tamaño del problema actual.
  - Escenario 2: Escalar el tamaño del problema conforme se añadan más procesadores.

#### Datos:

- Número de procesadores: 16
- Fracción paralelizable ( $P$ ): 0.7 (para Ley de Amdahl)
- Fracción secuencial ( $\alpha$ ): 0.3 (para Ley de Gustafson)

#### PREGUNTAS:

##### 1. Ley de Amdahl:

- Calcula la aceleración teórica del programa utilizando la Ley de Amdahl con 16 procesadores.

##### 2. Ley de Gustafson:

- Calcula la aceleración teórica del programa si asumimos que el tamaño del problema escala con el número de procesadores (Ley de Gustafson).

##### 3. Comparación:

- Compara los resultados obtenidos en los dos escenarios.
- ¿Cuál de las dos leyes resulta más favorable y en qué situaciones es mejor aplicar una u otra?

##### 4. Optimización:

- Si la empresa quiere mejorar la eficiencia del programa, sugiere estrategias concretas para reducir la fracción secuencial del código y aumentar la aceleración obtenida.





## EJERCICIO 2: OPTIMIZACIÓN DE UN PROCESO INDUSTRIAL MEDIANTE PARALELIZACIÓN

### Escenario:

Una fábrica de **automóviles eléctricos** está experimentando retrasos en la producción debido a un proceso de ensamblaje de **baterías y componentes electrónicos** que se realiza secuencialmente. El proceso actual implica los siguientes pasos:

1. **Montaje de las celdas de batería.**
2. **Verificación de calidad de cada celda ensamblada.**
3. **Ensamblaje de las celdas en módulos.**
4. **Inspección final y pruebas de carga.**

Actualmente, estos pasos se realizan uno tras otro, y la producción total de baterías es insuficiente para cumplir con la demanda creciente. La fábrica cuenta con un equipo de ingenieros interesados en aplicar la computación paralela y la automatización para optimizar el proceso.

### Preguntas:

1. **Identificación de oportunidades de paralelización:**
  - Analiza cada paso del proceso e identifica cuáles podrían ejecutarse en paralelo.
  - ¿Existen tareas dentro de un mismo paso que podrían dividirse y ejecutarse simultáneamente? Justifica tu respuesta.
2. **Modelo de paralelismo adecuado:**
  - ¿Qué tipo de paralelismo (a nivel de datos, instrucciones o tareas) propondrías para mejorar el proceso y por qué?
  - Por ejemplo, ¿podría la verificación de calidad dividirse entre distintos robots o sensores trabajando en paralelo?
3. **Ley de Amdahl o Gustafson:**
  - Supón que el 60% del proceso puede paralelizarse ( $P = 0.6$ ) y el resto sigue siendo secuencial.
    - Si la fábrica invierte en una línea de producción automatizada con 8 robots que trabajan en paralelo, calcula la aceleración teórica utilizando la **Ley de Amdahl**.
  - Luego, asume que el tamaño de la producción se puede escalar aumentando el número de baterías fabricadas simultáneamente. Calcula la aceleración usando la **Ley de Gustafson**.
4. **Diseño de una solución paralela:**
  - Propón un diseño específico para implementar la paralelización en la fábrica (por ejemplo, distribuir la inspección de celdas entre varios robots, dividir el montaje de módulos en diferentes estaciones, etc.).



- ¿Qué retos crees que podrían surgir al implementar esta paralelización (por ejemplo, sincronización, comunicación entre estaciones)?

### 5. Optimización adicional:

- ¿Qué otras tecnologías o mejoras (como optimización de algoritmos, robots colaborativos o simulaciones) podrías sugerir para complementar la paralelización y maximizar la producción?

### Sintetiza resultados:

Representa tu propuesta de solución en un diagrama de flujo o esquema visual de la línea de producción optimizada, destacando las secciones paralelizables.