

# Memoria Caché

La memoria caché es un componente esencial en la arquitectura de los sistemas de computación modernos. Sirve como una memoria de alta velocidad que almacena temporalmente los datos más utilizados, mejorando significativamente el rendimiento del sistema al reducir el tiempo de acceso a los datos. En el contexto de la programación paralela y distribuida, la gestión eficiente de la memoria caché es crucial para maximizar el rendimiento y minimizar la latencia de las operaciones.

## Tamaños de caché L1, L2, L3

La memoria caché se organiza en niveles jerárquicos para equilibrar el rendimiento y el costo. Estos niveles son comúnmente conocidos como L1, L2 y L3.

- **Caché L1:** Es la memoria caché más cercana al núcleo del procesador. Generalmente, cada núcleo de CPU tiene su propia caché L1, lo que permite un acceso extremadamente rápido a los datos. La caché L1 suele dividirse en dos partes: una caché de instrucciones (L1i) y una caché de datos (L1d). Los tamaños típicos de la caché L1 varían entre 32KB y 64KB por núcleo. La latencia de la caché L1 es muy baja, lo que permite tiempos de acceso de unos pocos ciclos de reloj.
- **Caché L2:** Esta caché es más grande y ligeramente más lenta que la L1. En muchos diseños modernos, cada núcleo de CPU tiene su propia caché L2, pero en algunos casos, puede ser compartida entre varios núcleos. Los tamaños típicos de la caché L2 oscilan entre 256KB y 512KB por núcleo, con una latencia de acceso mayor que la L1 pero aún bastante baja, generalmente entre 10 y 20 ciclos de reloj.
- **Caché L3:** Es la más grande y la más lenta de las tres. La caché L3 generalmente es compartida entre todos los núcleos de un procesador. Los tamaños de la caché L3 pueden variar significativamente, desde 2MB hasta 32MB o más, dependiendo del diseño del procesador. La latencia de acceso a la caché L3 es considerablemente mayor que la de las cachés L1 y L2, típicamente en el rango de 30 a 50 ciclos de reloj.

## Políticas de reemplazo de Caché

Dado que las memorias caché tienen un tamaño limitado, es esencial contar con políticas de reemplazo eficientes para decidir qué datos mantener y cuáles descartar cuando la caché se llena. Algunas de las políticas de reemplazo más comunes incluyen:

- **Least Recently Used (LRU):** Esta política reemplaza el bloque de caché que no ha sido usado por el mayor tiempo. LRU es popular debido a su simplicidad y eficacia, aunque puede ser costoso en términos de hardware y complejidad computacional.
- **First In, First Out (FIFO):** Reemplaza el bloque más antiguo en la caché, sin considerar cuándo fue accedido por última vez. FIFO es fácil de implementar pero puede no ser tan eficiente como LRU en términos de rendimiento.

- **Least Frequently Used (LFU):** Esta política reemplaza el bloque que ha sido usado con menor frecuencia. LFU puede ser más difícil de implementar debido a la necesidad de mantener un contador para cada bloque de caché.
- **Random Replacement:** Reemplaza un bloque al azar. Esta política es fácil de implementar y, sorprendentemente, puede ser eficaz en ciertos escenarios.

## **Latencia de caché**

La latencia de la caché es el tiempo que tarda el procesador en acceder a los datos almacenados en la caché. La latencia está influenciada por varios factores, incluyendo el tamaño de la caché, su ubicación jerárquica (L1, L2, L3) y la política de reemplazo utilizada. En términos de programación paralela y distribuida, la latencia de la caché puede tener un impacto significativo en el rendimiento de las aplicaciones, especialmente aquellas que requieren un acceso intensivo a los datos.

La baja latencia de la caché L1 permite a los procesadores realizar operaciones rápidas y eficientes, reduciendo el tiempo de espera para acceder a los datos. Sin embargo, a medida que se escala hacia cachés L2 y L3, la latencia aumenta, lo que puede ralentizar las operaciones si los datos necesarios no están presentes en la caché de nivel superior.

## **Gestión de caché en la programación paralela y distribuida**

En entornos de programación paralela y distribuida, la gestión de la caché se vuelve aún más crítica. La concurrencia y la paralelización pueden llevar a situaciones donde múltiples núcleos acceden a los mismos datos, creando posibles conflictos y problemas de coherencia. Algunas de las estrategias y técnicas utilizadas para gestionar eficientemente la caché en estos entornos incluyen:

- **Localidad de referencia:** Optimizar los algoritmos para mejorar la localidad temporal y espacial puede aumentar significativamente la eficiencia de la caché. La localidad temporal se refiere a la reutilización de datos en intervalos cortos de tiempo, mientras que la localidad espacial se refiere al acceso a datos cercanos en la memoria.
- **Afinidad de núcleo:** Asignar tareas a núcleos específicos y mantener la afinidad de los datos a esos núcleos puede reducir la necesidad de transferir datos entre cachés de diferentes núcleos, mejorando así la eficiencia.
- **Coherencia de caché:** En sistemas multinúcleo, mantener la coherencia de caché es crucial. Protocolos de coherencia de caché como MESI (Modified, Exclusive, Shared, Invalid) y MOESI (Modified, Owner, Exclusive, Shared, Invalid) son utilizados para asegurar que todos los núcleos vean una visión consistente de la memoria.
- **Prefetching:** El prefetching es una técnica donde se cargan datos en la caché antes de que sean necesarios, basándose en patrones de acceso previsibles. Esto puede reducir la latencia de caché al tener los datos ya disponibles cuando se requieren.

- **Particionado de caché:** Dividir la caché en segmentos asignados a diferentes núcleos o procesos puede reducir la interferencia y mejorar el rendimiento, especialmente en sistemas donde las tareas tienen diferentes patrones de acceso a los datos.
- **Optimización de algoritmos:** Adaptar los algoritmos para ser más "cache-friendly" es una técnica clave. Esto puede implicar el reordenamiento de bucles, el uso de estructuras de datos adecuadas y la minimización de accesos a memoria dispersos.

El impacto de la gestión de caché en el rendimiento de las aplicaciones paralelas y distribuidas no puede ser subestimado. En aplicaciones de alto rendimiento, como la computación científica y la simulación, la eficiencia de la caché puede determinar la viabilidad del procesamiento en tiempo real. La optimización del uso de la caché puede llevar a mejoras significativas en el rendimiento, reduciendo la latencia de acceso a datos y mejorando el throughput general del sistema.

## Ejercicios

1. Analiza este ejercicio te permitirá medir la latencia de acceso a diferentes tamaños de datos en Python.

```
import time
import numpy as np

def measure_latency(array_size):
    array = np.zeros(array_size)
    start_time = time.perf_counter()

    # Acceder a todos los elementos del array
    for i in range(array_size):
        array[i] += 1

    end_time = time.perf_counter()
    latency = end_time - start_time
    return latency

sizes = [10**3, 10**4, 10**5, 10**6, 10**7]
latencies = []

for size in sizes:
    latency = measure_latency(size)
    latencies.append(latency)
    print(f"Tamaño del array: {size}, Latencia: {latency:.6f} segundos")

# Graficar los resultados
import matplotlib.pyplot as plt

plt.plot(sizes, latencies, marker='o')
plt.xlabel('Tamaño del Array')
plt.ylabel('Latencia (segundos)')
```

```
plt.xscale('log')
plt.yscale('log')
plt.title('Latencia de Acceso a Datos en Diferentes Tamaños de Array')
plt.show()
```

2. Analiza la importancia de la localidad de referencia en el acceso a datos. Compara el tiempo de acceso cuando se recorren los datos de manera secuencial versus de manera aleatoria.

```
import random

def measure_access_time(array, indices):
    start_time = time.perf_counter()

    for i in indices:
        array[i] += 1

    end_time = time.perf_counter()
    access_time = end_time - start_time
    return access_time

size = 10**6
array = np.zeros(size)
sequential_indices = list(range(size))
random_indices = sequential_indices.copy()
random.shuffle(random_indices)

sequential_time = measure_access_time(array, sequential_indices)
random_time = measure_access_time(array, random_indices)

print(f"Tiempo de acceso secuencial: {sequential_time:.6f} segundos")
print(f"Tiempo de acceso aleatorio: {random_time:.6f} segundos")
```

3. Analiza cómo la coherencia de caché afecta el rendimiento en un entorno multihilo.

```
import threading

def worker(shared_data, start, end):
    for i in range(start, end):
        shared_data[i] += 1

size = 10**6
shared_data = np.zeros(size)
num_threads = 4
threads = []

chunk_size = size // num_threads
for i in range(num_threads):
    start = i * chunk_size
    end = (i + 1) * chunk_size
```

```

        thread = threading.Thread(target=worker, args=(shared_data, start,
end))
        threads.append(thread)

start_time = time.perf_counter()

for thread in threads:
    thread.start()

for thread in threads:
    thread.join()

end_time = time.perf_counter()
total_time = end_time - start_time

print(f"Tiempo total con {num_threads} hilos: {total_time:.6f}
segundos")

```

4. Analiza el código que implementa una forma sencilla de prefetching de datos, cargando datos en la caché antes de que sean necesarios.

```

def prefetch_data(array, indices, prefetch_distance):
    for i in range(len(indices) - prefetch_distance):
        _ = array[indices[i + prefetch_distance]] # Prefetch
        array[indices[i]] += 1 # Access

size = 10**6
array = np.zeros(size)
indices = list(range(size))
random.shuffle(indices)
prefetch_distance = 10

start_time = time.perf_counter()
prefetch_data(array, indices, prefetch_distance)
end_time = time.perf_counter()

total_time = end_time - start_time
print(f"Tiempo con prefetching: {total_time:.6f} segundos")

```

5. Analiza cómo la afinidad de núcleo puede influir en el rendimiento.

```

import multiprocessing

def core_affinity_worker(core_id, shared_data, start, end):
    p = multiprocessing.current_process()
    p.cpu_affinity([core_id]) # Asigna el núcleo específico
    for i in range(start, end):
        shared_data[i] += 1

size = 10**6

```

```

shared_data = np.zeros(size)
num_processes = 4
processes = []

chunk_size = size // num_processes
for i in range(num_processes):
    start = i * chunk_size
    end = (i + 1) * chunk_size
    process = multiprocessing.Process(target=core_affinity_worker,
    args=(i, shared_data, start, end))
    processes.append(process)

start_time = time.perf_counter()

for process in processes:
    process.start()

for process in processes:
    process.join()

end_time = time.perf_counter()
total_time = end_time - start_time

print(f"Tiempo total con afinidad de núcleo y {num_processes}
procesos: {total_time:.6f} segundos")

```

6 . Implementa una política de reemplazo LRU para una caché y simula el acceso a los datos.

```

from collections import OrderedDict

class LRUCache:
    def __init__(self, capacity):
        self.capacity = capacity
        self.cache = OrderedDict()

    def get(self, key):
        if key not in self.cache:
            return -1
        self.cache.move_to_end(key)
        return self.cache[key]

    def put(self, key, value):
        if key in self.cache:
            self.cache.move_to_end(key)
        self.cache[key] = value
        if len(self.cache) > self.capacity:
            self.cache.popitem(last=False)

# Simulación de acceso
cache = LRUCache(5)

```

```

operations = [("put", 1, "data1"), ("put", 2, "data2"), ("get", 1),
("put", 3, "data3"),
                ("put", 4, "data4"), ("put", 5, "data5"), ("get", 2),
("put", 6, "data6"),
                ("get", 3), ("get", 1)]

for op in operations:
    if op[0] == "put":
        cache.put(op[1], op[2])
        print(f"Put: {op[1]} -> {op[2]}")
    elif op[0] == "get":
        result = cache.get(op[1])
        print(f"Get: {op[1]} -> {result}")
    print(f"Estado de la caché: {list(cache.cache.items())}")

```

7. Implementa una política de reemplazo FIFO para una caché y simula el acceso a los datos.

```

from collections import deque

class FIFOCache:
    def __init__(self, capacity):
        self.capacity = capacity
        self.cache = {}
        self.order = deque()

    def get(self, key):
        return self.cache.get(key, -1)

    def put(self, key, value):
        if key not in self.cache:
            if len(self.cache) >= self.capacity:
                oldest = self.order.popleft()
                del self.cache[oldest]
            self.cache[key] = value
            self.order.append(key)
        else:
            self.cache[key] = value

# Simulación de acceso
cache = FIFOCache(3)
operations = [("put", 1, "data1"), ("put", 2, "data2"), ("put", 3,
"data3"), ("get", 1),
                ("put", 4, "data4"), ("get", 2), ("put", 5, "data5"),
("get", 3), ("get", 1)]

for op in operations:
    if op[0] == "put":
        cache.put(op[1], op[2])
        print(f"Put: {op[1]} -> {op[2]}")
    elif op[0] == "get":

```

```

    result = cache.get(op[1])
    print(f"Get: {op[1]} -> {result}")
    print(f"Estado de la caché: {list(cache.cache.items())}")

```

8. Implementa una política de reemplazo LFU (Reemplazo Least Frequently Used) para una caché y simula el acceso a los datos.

```

from collections import defaultdict, OrderedDict

class LFUCache:
    def __init__(self, capacity):
        self.capacity = capacity
        self.cache = {}
        self.freq = defaultdict(OrderedDict)
        self.min_freq = 0

    def get(self, key):
        if key not in self.cache:
            return -1
        freq = self.cache[key][1]
        self.cache[key][1] += 1
        value = self.cache[key][0]
        del self.freq[freq][key]
        if not self.freq[freq]:
            del self.freq[freq]
            if self.min_freq == freq:
                self.min_freq += 1
        self.freq[freq + 1][key] = None
        return value

    def put(self, key, value):
        if self.capacity == 0:
            return
        if key in self.cache:
            self.cache[key][0] = value
            self.get(key)
            return
        if len(self.cache) >= self.capacity:
            evict = next(iter(self.freq[self.min_freq]))
            del self.freq[self.min_freq][evict]
            del self.cache[evict]
        self.cache[key] = [value, 1]
        self.freq[1][key] = None
        self.min_freq = 1

# Simulación de acceso
cache = LFUCache(3)
operations = [("put", 1, "data1"), ("put", 2, "data2"), ("put", 3, "data3"), ("get", 1),

```



```

        ("put", 4, "data4"), ("get", 2), ("put", 5, "data5"),
        ("get", 3), ("get", 1)]

for op in operations:
    if op[0] == "put":
        cache.put(op[1], op[2])
        print(f"Put: {op[1]} -> {op[2]}")
    elif op[0] == "get":
        result = cache.get(op[1])
        print(f"Get: {op[1]} -> {result}")
    print(f"Estado de la caché: {list(cache.cache.items())}")

```

9 . Implementa una política de reemplazo aleatorio para una caché y simula el acceso a los datos.

```

import random

class RandomCache:
    def __init__(self, capacity):
        self.capacity = capacity
        self.cache = {}

    def get(self, key):
        return self.cache.get(key, -1)

    def put(self, key, value):
        if key not in self.cache:
            if len(self.cache) >= self.capacity:
                evict = random.choice(list(self.cache.keys()))
                del self.cache[evict]
            self.cache[key] = value

# Simulación de acceso
cache = RandomCache(3)
operations = [("put", 1, "data1"), ("put", 2, "data2"), ("put", 3,
"data3"), ("get", 1),
        ("put", 4, "data4"), ("get", 2), ("put", 5, "data5"),
        ("get", 3), ("get", 1)]

for op in operations:
    if op[0] == "put":
        cache.put(op[1], op[2])
        print(f"Put: {op[1]} -> {op[2]}")
    elif op[0] == "get":
        result = cache.get(op[1])
        print(f"Get: {op[1]} -> {result}")
    print(f"Estado de la caché: {list(cache.cache.items())}")

```

10 . Analiza la estrategia de prefetching de datos en un entorno paralelo

```

from concurrent.futures import ThreadPoolExecutor

def prefetching_worker(array, indices, prefetch_distance):
    for i in range(len(indices) - prefetch_distance):
        _ = array[indices[i + prefetch_distance]] # Prefetch
        array[indices[i]] += 1 # Access

size = 10**6
array = np.zeros(size)
indices = list(range(size))
random.shuffle(indices)
prefetch_distance = 10

with ThreadPoolExecutor(max_workers=4) as executor:
    futures = [executor.submit(prefetching_worker, array,
indices[i::4], prefetch_distance) for i in range(4)]

    for future in futures:
        future.result()

print(f"Prefetching completado")

```

11. Analiza este algoritmo para mejorar la localidad temporal y espacial.

```

def locality_optimized_algorithm(matrix):
    n = len(matrix)
    result = np.zeros((n, n))

    for i in range(n):
        for j in range(n):
            result[i][j] = matrix[i][j] + matrix[j][i]

    return result

n = 1000
matrix = np.random.rand(n, n)

start_time = time.perf_counter()
result = locality_optimized_algorithm(matrix)
end_time = time.perf_counter()

print(f"Tiempo de ejecución: {end_time - start_time:.6f} segundos")

```

12. Analiza la simulación de coherencia de caché en un entorno multinúcleo.

```

import multiprocessing
import numpy as np

def coherence_worker(shared_array, start, end):

```

```

        for i in range(start, end):
            shared_array[i] += 1

size = 10**6
shared_array = multiprocessing.Array('d', size)
num_processes = 4
chunk_size = size // num_processes
processes = []

for i in range(num_processes):
    start = i * chunk_size
    end = (i + 1) * chunk_size
    process = multiprocessing.Process(target=coherence_worker,
    args=(shared_array, start, end))
    processes.append(process)

start_time = time.perf_counter()

for process in processes:
    process.start()

for process in processes:
    process.join()

end_time = time.perf_counter()
total_time = end_time - start_time

print(f"Tiempo total con coherencia de caché y {num_processes}
procesos: {total_time:.6f} segundos")

## Tus respuestas

```

## Coherencia de caché

La coherencia de caché es un concepto crucial en sistemas de memoria compartida, especialmente en entornos de multiprocesadores donde múltiples procesadores o núcleos acceden a una memoria común. La coherencia de caché asegura que todas las copias de una misma ubicación de memoria, almacenadas en las cachés de distintos procesadores, mantengan valores consistentes, evitando así problemas de inconsistencia que pueden llevar a errores en los cálculos y resultados inesperados en las aplicaciones.

En sistemas de multiprocesadores, cada procesador generalmente tiene su propia caché para reducir la latencia de acceso a datos y mejorar el rendimiento general del sistema. Sin embargo, esta mejora en el rendimiento viene con el desafío de mantener la coherencia entre las diferentes copias de datos almacenadas en las cachés de los procesadores. Si no se maneja adecuadamente, diferentes procesadores podrían trabajar con datos desactualizados, llevando a errores de sincronización y resultados incorrectos.

### Protocolos de coherencia de caché

Los protocolos de coherencia de caché son mecanismos diseñados para gestionar la consistencia de datos en cachés de múltiples procesadores. Estos protocolos aseguran que cualquier cambio en los datos de una caché se refleje adecuadamente en las demás cachés que contengan una copia de esos datos. Existen varios protocolos de coherencia de caché, cada uno con sus propias estrategias y mecanismos para mantener la coherencia.

- El Protocolo MESI, que es un acrónimo de los cuatro estados posibles de una línea de caché: Modified (Modificado), Exclusive (Exclusivo), Shared (Compartido), e Invalid (Inválido). En el estado Modificado, una línea de caché ha sido alterada y no es coherente con la memoria principal. En el estado Exclusivo, una línea de caché es la única copia en todas las cachés y es coherente con la memoria principal. En el estado Compartido, una línea de caché puede estar presente en múltiples cachés y es coherente con la memoria principal. En el estado Inválido, la línea de caché no es válida y no debe ser utilizada.

El protocolo MESI funciona con base en dos tipos de transacciones: las de lectura y las de escritura. Cuando un procesador desea leer un dato, verifica si la línea de caché está en estado Modificado, Exclusivo o Compartido. Si está en estado Inválido, debe obtener el dato de la memoria principal o de otra caché. Cuando un procesador desea escribir un dato, si la línea de caché está en estado Modificado o Exclusivo, puede proceder con la escritura. Si está en estado Compartido o Inválido, debe invalidar las otras copias antes de escribir.

- El Protocolo MOESI, que es una extensión del protocolo MESI e introduce un estado adicional: Owned (Propietario). En el estado Owned, una línea de caché es la única copia modificada, pero otras cachés pueden tener copias compartidas. El procesador propietario es responsable de actualizar la memoria principal cuando la línea es reemplazada. Este estado adicional permite optimizaciones en la comunicación entre cachés y reduce la necesidad de acceder a la memoria principal para mantener la coherencia.
- El protocolo MSI es una simplificación del protocolo MESI, con solo tres estados: Modified, Shared e Invalid. El estado Exclusivo no está presente en este protocolo. Aunque es más simple de implementar, puede ser menos eficiente en ciertos escenarios debido a la ausencia del estado Exclusivo, lo que podría llevar a más invalidaciones y tráfico de memoria.
- El protocolo Dragon es otro ejemplo de un protocolo de coherencia de caché utilizado en sistemas de multiprocesadores. Introduce estados adicionales como Shared-clean y Shared-modified, que permiten una mayor flexibilidad en la gestión de las líneas de caché compartidas y modificadas. Este protocolo se utiliza para minimizar el tráfico de coherencia y mejorar el rendimiento en sistemas con alta concurrencia.

Los protocolos de coherencia de caché no solo se diferencian por los estados que manejan, sino también por las políticas de invalidación y actualización que implementan. Las políticas de invalidación aseguran que cuando un procesador modifica un dato, las demás cachés invalidan sus copias de ese dato. Esto garantiza que los procesadores no trabajen con datos desactualizados. Las políticas de actualización, por otro lado, actualizan automáticamente las copias de los datos en otras cachés cuando un procesador realiza una modificación. Esto puede

reducir la latencia de acceso a datos actualizados, pero también puede aumentar el tráfico de coherencia.

La coherencia de caché también se puede lograr mediante el uso de **directorios de coherencia**, que son estructuras de datos centralizadas o distribuidas que mantienen información sobre las copias de datos en las cachés de los procesadores. Los directorios de coherencia pueden ser efectivos para reducir el tráfico de coherencia y mejorar el rendimiento en sistemas con un gran número de procesadores.

En entornos de programación paralela y distribuida, la coherencia de caché es fundamental para asegurar la consistencia de los datos y el correcto funcionamiento de los programas. Los desarrolladores deben estar conscientes de las implicaciones de la coherencia de caché y utilizar técnicas y herramientas que les permitan gestionar adecuadamente la coherencia en sus aplicaciones. Esto incluye el uso de bibliotecas y marcos de trabajo que implementen protocolos de coherencia de caché, así como el diseño de algoritmos y estructuras de datos que minimicen la necesidad de acceso concurrente a datos compartidos.

Además, la coherencia de caché puede tener un impacto significativo en el rendimiento de las aplicaciones paralelas y distribuidas. El tráfico de coherencia y las latencias asociadas pueden convertirse en cuellos de botella, limitando la escalabilidad y eficiencia de las aplicaciones. Por lo tanto, es esencial optimizar el acceso a la memoria y el uso de la caché para minimizar estos impactos.

El manejo eficiente de la coherencia de caché también implica el uso de técnicas de prefetching y políticas de reemplazo de caché que aprovechen la localidad temporal y espacial de los datos. El prefetching anticipa las necesidades de datos futuros y los carga en la caché antes de que sean necesarios, reduciendo así las latencias de acceso. Las políticas de reemplazo de caché, por otro lado, determinan qué datos deben mantenerse en la caché y cuáles deben ser reemplazados cuando la caché está llena, basándose en patrones de acceso recientes y previstos.

## Consistencia de caché

La consistencia de caché es un problema crítico en sistemas multiprocesadores donde múltiples núcleos de procesamiento acceden a una memoria compartida. La caché se utiliza para mejorar el rendimiento al almacenar datos frecuentemente accedidos cerca del procesador, reduciendo así el tiempo de acceso a la memoria. Sin embargo, en un sistema multiprocesador, mantener la coherencia entre múltiples cachés que pueden almacenar copias de los mismos datos se vuelve un desafío significativo.

### Fundamentos de la consistencia de caché

La memoria caché es una memoria de alta velocidad que almacena copias de datos de la memoria principal. Cada núcleo del procesador puede tener su propia caché, permitiendo un acceso más rápido a los datos frecuentemente usados.

El problema surge cuando múltiples cachés contienen copias de los mismos datos y una de estas copias es modificada. Es crucial asegurar que todas las copias reflejen el cambio para mantener la coherencia de datos. Si una caché modifica un dato y las otras cachés no se actualizan, los procesadores pueden trabajar con datos inconsistentes, llevando a errores en la ejecución del programa.

## Protocolos de consistencia de caché

Para manejar la consistencia de caché, se utilizan varios protocolos. Estos protocolos aseguran que cualquier modificación en los datos se refleje en todas las copias existentes de la caché. Algunos de los protocolos más comunes incluyen:

Protocolo MESI: MESI es uno de los protocolos de coherencia de caché más utilizados y sus siglas representan los cuatro estados posibles de una línea de caché: Modificado, Exclusivo, Compartido e Inválido.

- Modificado (M): La línea de caché ha sido modificada y no coincide con la copia en la memoria principal. Solo esta caché tiene una copia válida.
- Exclusivo (E): La línea de caché es la única copia y coincide con la memoria principal.
- Compartido (S): La línea de caché puede ser compartida por múltiples cachés y coincide con la memoria principal.
- Inválido (I): La línea de caché no es válida.

Este protocolo asegura que cualquier modificación de datos en una caché se refleje en las otras cachés, invalidando las copias antiguas o actualizando las mismas.

Protocolo MOESI: MOESI es una extensión del protocolo MESI que agrega un estado adicional: Owned (O). Este estado permite que una línea de caché modificada se comparta entre varias cachés, con una caché actuando como la "propietaria" que coordina la actualización.

- Owned (O): La línea de caché es una copia modificada que puede ser compartida y es responsable de propagar los cambios a la memoria principal cuando sea necesario.

Protocolo MSI: MSI es una versión simplificada de MESI con solo tres estados: Modificado, Compartido e Inválido. Aunque más simple, puede llevar a más tráfico de bus debido a la falta del estado Exclusivo.

## Mecánicas de consistencia

Snooping: El snooping es una técnica donde todas las cachés en un sistema observan el bus para detectar operaciones que puedan afectar las líneas de caché que almacenan. Si una caché detecta una operación relevante (como una escritura en una línea de caché que contiene), toma medidas para mantener la coherencia.

Directorios: Los sistemas basados en directorios utilizan una estructura centralizada que mantiene un registro del estado de cada línea de caché en todas las cachés del sistema. Este enfoque escala mejor para sistemas con muchos núcleos, pero introduce una latencia adicional debido a la necesidad de consultar el directorio.

### Ejemplo 1:

Imaginemos un sistema con dos núcleos, cada uno con su propia caché L1. Ambos núcleos necesitan acceder a una variable X almacenada en la memoria principal. Inicialmente, ambos cachés cargan X, y ambas copias están en el estado Compartido.

- Si el núcleo 1 modifica X, su caché cambia el estado de la línea de X a Modificado. El núcleo 1 debe invalidar la copia en la caché del núcleo 2.

- Si el núcleo 2 intenta leer X después de la modificación, encontrará su copia inválida y deberá solicitar la nueva copia desde la memoria principal o desde la caché del núcleo 1.

### Ejemplo 2:

Considera un sistema con cuatro núcleos utilizando el protocolo MOESI. El núcleo 1 modifica un dato D, cambiando su estado a Modificado.

- El núcleo 2 solicita D, y el núcleo 1 cambia el estado a Owned y responde con la copia modificada. El núcleo 2 ahora tiene D en estado Compartido.
- Si el núcleo 3 también solicita D, el núcleo 1 o el 2 pueden proporcionar la copia, manteniendo la coherencia sin necesidad de acceder a la memoria principal.

### Desafíos en la consistencia de caché

- Latencia de Comunicación: Mantener la coherencia de caché introduce latencia adicional debido al tráfico en el bus o las consultas al directorio.
- Escalabilidad: A medida que el número de núcleos en un sistema aumenta, el tráfico de coherencia se incrementa, lo que puede llevar a cuellos de botella. Los sistemas basados en snooping pueden no escalar bien para un gran número de núcleos debido al aumento del tráfico de bus.
- Complejidad del protocolo: Los protocolos más avanzados, como MOESI, introducen complejidad adicional en el diseño de hardware, lo que puede incrementar los costos de desarrollo y fabricación.

Herramientas como GEM5 y SimpleScalar permiten a los investigadores simular diferentes arquitecturas de caché y protocolos de coherencia para evaluar su desempeño.

Los procesadores modernos de Intel y AMD implementan sofisticados protocolos de coherencia de caché, como MESIF y MOESIF, que son variaciones avanzadas de los protocolos MESI y MOESI.

## Ejercicios

1. Explica cada uno de los estados del protocolo MESI (Modified, Exclusive, Shared, Invalid) y dé un ejemplo práctico de cuándo una línea de caché puede cambiar de un estado a otro.
2. Compara y contraste los protocolos de coherencia de caché basados en invalidación (ej. MESI) y basados en actualización (ej. MOESI). ¿Cuáles son las ventajas y desventajas de cada uno?
3. Discute cómo los diferentes protocolos de coherencia de caché pueden impactar el rendimiento de un sistema multiprocesador. Considere factores como el tráfico de red y la latencia de acceso a datos.
4. Explica cómo la coherencia de caché afecta la escalabilidad de sistemas multiprocesadores. ¿Qué desafíos presentan los sistemas con un gran número de procesadores?
5. Explica cómo funciona un directorio de coherencia en un sistema multiprocesador y cuáles son sus ventajas sobre los protocolos basados en bus.

**## Tus respuestas**

## 6 . Implementación de un simulador de Protocolo MESI:

Descripción: Implementa un simulador simple del protocolo MESI en Python. Simula varios procesadores y cachés que acceden y modifican líneas de datos, cambiando los estados de las líneas de caché según las reglas de MESI.

```
# Código de inicio
class CacheLine:
    def __init__(self):
        self.state = 'Invalid'
        self.data = None

class Processor:
    def __init__(self, id):
        self.id = id
        self.cache = {}

    def read(self, address):
        if address in self.cache and self.cache[address].state != 'Invalid':
            return self.cache[address].data
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Shared'
            self.cache[address].data = memory_read(address)
            return self.cache[address].data

    def write(self, address, data):
        if address in self.cache and (self.cache[address].state == 'Shared' or self.cache[address].state == 'Exclusive'):
            self.cache[address].state = 'Modified'
            self.cache[address].data = data
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Modified'
            self.cache[address].data = data
            memory_write(address, data)

    def memory_read(address):
        # Simula una lectura de memoria principal
        return "data_from_memory"

    def memory_write(address, data):
        # Simula una escritura en memoria principal
        pass

# Ejemplo de uso
p1 = Processor(1)
p2 = Processor(2)
```



```
p1.write(0x1A, "data1")
p2.read(0x1A)
```

*## Tu respuesta*

## 7. Simulación de tráfico de coherencia:

Descripción: Extiende el simulador anterior para incluir múltiples procesadores y rastrear el tráfico de coherencia (lecturas y escrituras en el bus de memoria).

a.

```
## Código inicial

traffic = 0

class CacheLine:
    def __init__(self):
        self.state = 'Invalid'
        self.data = None

class Processor:
    def __init__(self, id):
        self.id = id
        self.cache = {}

    def read(self, address):
        global traffic
        if address in self.cache and self.cache[address].state != 'Invalid':
            return self.cache[address].data
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Shared'
            traffic += 1
            self.cache[address].data = memory_read(address)
            return self.cache[address].data

    def write(self, address, data):
        global traffic
        if address in self.cache and (self.cache[address].state == 'Shared' or self.cache[address].state == 'Exclusive'):
            self.cache[address].state = 'Modified'
            self.cache[address].data = data
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Modified'
            traffic += 1
            self.cache[address].data = data
            memory_write(address, data)
```

```

def memory_read(address):
    return "data_from_memory"

def memory_write(address, data):
    pass

# Ejemplo de uso
p1 = Processor(1)
p2 = Processor(2)
p1.write(0x1A, "data1")
p2.read(0x1A)
print(f"Traffic: {traffic} transactions")

# Tu respuesta

```

## 8 . Protocolo de directorio distribuido:

Descripción: Implementa un simulador básico de un protocolo de directorio distribuido. Cada procesador tiene su propio directorio local que mantiene información sobre las líneas de datos.

```

# Codigo de inicio
class Directory:
    def __init__(self):
        self.entries = {}

    def update(self, address, processor_id):
        if address not in self.entries:
            self.entries[address] = set()
        self.entries[address].add(processor_id)

    def invalidate(self, address):
        if address in self.entries:
            self.entries[address] = set()

class Processor:
    def __init__(self, id, directory):
        self.id = id
        self.cache = {}
        self.directory = directory

    def read(self, address):
        if address in self.cache and self.cache[address].state != 'Invalid':
            return self.cache[address].data
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Shared'
            self.directory.update(address, self.id)
            self.cache[address].data = memory_read(address)

```

```

        return self.cache[address].data

    def write(self, address, data):
        if address in self.cache and (self.cache[address].state ==
'Shared' or self.cache[address].state == 'Exclusive'):
            self.cache[address].state = 'Modified'
            self.cache[address].data = data
            self.directory.invalidate(address)
        else:
            self.cache[address] = CacheLine()
            self.cache[address].state = 'Modified'
            self.cache[address].data = data
            self.directory.invalidate(address)
            memory_write(address, data)

def memory_read(address):
    return "data_from_memory"

def memory_write(address, data):
    pass

# Ejemplo de uso
directory = Directory()
p1 = Processor(1, directory)
p2 = Processor(2, directory)
p1.write(0x1A, "data1")
p2.read(0x1A)
print(directory.entries)

# Tu respuesta

```

## Problemas de coherencia y consistencia

El problema de la coherencia y consistencia de caché es un desafío fundamental en la arquitectura de sistemas multiprocesadores. La necesidad de mantener una vista coherente y consistente de la memoria compartida entre múltiples procesadores es crucial para asegurar la correcta ejecución de aplicaciones paralelas y distribuidas. A medida que los sistemas de computación evolucionan y la demanda de rendimiento crece, se hace imperativo implementar mecanismos efectivos para gestionar la coherencia y consistencia de caché. Uno de los enfoques más comunes para abordar este problema es el uso de protocolos de coherencia de caché basados en snooping y sistemas de bus compartido.

En un sistema multiprocesador, cada procesador tiene su propia caché local para acelerar el acceso a datos frecuentes y reducir la latencia de acceso a memoria. Sin embargo, cuando múltiples procesadores trabajan en la misma región de memoria, pueden ocurrir inconsistencias si no se implementan mecanismos adecuados para mantener la coherencia de caché.

Uno de los métodos más comunes para mantener la coherencia de caché en sistemas multiprocesadores es el uso de un **protocolo de snooping**. En este enfoque, todas las cachés monitorean (o "husmean") las transacciones en un bus compartido para detectar accesos a

líneas de memoria que pueden estar en caché localmente. Cuando un procesador realiza una operación de lectura o escritura en una línea de memoria, emite una transacción en el bus compartido. Los demás procesadores observan esta transacción y, si tienen una copia de la línea de memoria en sus cachés, actúan en consecuencia para mantener la coherencia.

El protocolo MESI (Modified, Exclusive, Shared, Invalid) y MOESI son uno de los protocolos de snooping más utilizados.

Cuando un procesador desea leer o escribir una línea de caché, realiza una transacción en el bus. Si otro procesador tiene una copia de la línea en estado Modified, debe escribir los datos de vuelta en la memoria principal (write-back) antes de que el primer procesador pueda proceder. Este mecanismo asegura que todas las copias de una línea de caché estén sincronizadas y que las modificaciones sean visibles para todos los procesadores.

MOESI puede reducir el tráfico de coherencia al permitir que las líneas de caché modificadas se compartan sin necesidad de escribir de vuelta a la memoria principal inmediatamente.

Además de los protocolos basados en snooping, existen protocolos de coherencia basados en directorios, que son más adecuados para sistemas con un gran número de procesadores. En estos sistemas, el uso de un bus compartido puede convertirse en un cuello de botella debido al elevado tráfico de coherencia. Los **protocolos de directorio** utilizan una estructura de datos centralizada o distribuida que mantiene un registro de las líneas de caché y sus estados en cada procesador. Cuando un procesador desea acceder a una línea de memoria, consulta el directorio para determinar si la línea está presente en otras cachés y en qué estado se encuentra. El directorio coordina las actualizaciones y las invalidaciones, reduciendo la necesidad de tráfico en el bus compartido.

La consistencia de caché es otro aspecto crítico que debe gestionarse en sistemas multiprocesadores. Los modelos de consistencia definen el orden en que las operaciones de memoria deben ser visibles a los procesadores. El modelo de consistencia más estricto es la **consistencia secuencial**, que asegura que las operaciones de memoria se vean en el mismo orden en que se ejecutaron. Sin embargo, la consistencia secuencial puede ser costosa en términos de rendimiento, ya que requiere una fuerte sincronización entre los procesadores.

Para mejorar el rendimiento, se utilizan modelos de consistencia más relajados, como la **consistencia débil** y la **consistencia de memoria liberada** (release consistency). En la consistencia débil, las operaciones de memoria pueden ejecutarse en cualquier orden, siempre y cuando se sincronicen en puntos específicos del programa. La consistencia de memoria liberada permite que las operaciones de lectura y escritura se realicen en cualquier orden, pero requiere que las operaciones de sincronización (como las barreras de memoria) se ejecuten en un orden específico para asegurar la correcta ejecución del programa.

El problema de la coherencia y consistencia de caché es particularmente desafiante en sistemas distribuidos, donde los nodos pueden estar ubicados en diferentes ubicaciones geográficas y comunicarse a través de redes de alta latencia. En estos sistemas, mantener una vista coherente y consistente de la memoria compartida requiere algoritmos sofisticados y técnicas de sincronización eficientes. Por ejemplo, los **algoritmos de consenso** como Paxos y Raft se utilizan para asegurar que los datos sean consistentes en todos los nodos, incluso en presencia de fallos de red y nodos defectuosos.

Además, en sistemas distribuidos modernos, se emplean técnicas como el uso de **cachés distribuidas** y **replicación de datos** para mejorar el rendimiento y la disponibilidad. Las cachés distribuidas permiten que los datos se almacenen en múltiples nodos, reduciendo la latencia de acceso al aprovechar la cercanía geográfica de los datos. Sin embargo, esto introduce nuevos desafíos para mantener la coherencia y consistencia de los datos replicados. Las políticas de coherencia deben asegurarse de que cualquier actualización realizada en una copia de los datos se propague a todas las demás copias, y las políticas de consistencia deben definir cómo se manejan los conflictos de actualización y las fallas de red.

## Ejercicios

- Explica el problema de la coherencia de caché en un sistema multiprocesador. ¿Por qué es importante mantener la coherencia de la caché?
- Describe el protocolo MESI. ¿Cuáles son los estados involucrados y qué significan? Explica cómo el protocolo MESI maneja una operación de lectura y escritura en un entorno multiprocesador.
- Compara y contraste los protocolos MSI y MESI. ¿Cuáles son las diferencias clave entre estos protocolos de coherencia de caché?
- Explica cómo funcionan los protocolos de directorio para mantener la coherencia de la caché. ¿Qué ventajas ofrecen estos protocolos sobre los basados en snooping?
- Describe las diferencias entre un protocolo de directorio centralizado y uno distribuido. ¿Cuáles son los beneficios y desventajas de cada enfoque?
- Analiza un escenario en el que un protocolo de directorio distribuido podría ser más eficiente que uno centralizado. ¿Qué factores influirían en esta decisión?
- Explica el propósito de los algoritmos de consenso en sistemas distribuidos. ¿Por qué es importante alcanzar el consenso entre nodos en un sistema distribuido?
- Describe el algoritmo Paxos. ¿Cuáles son los roles principales (proposer, acceptor, learner) y cómo interactúan entre sí para alcanzar el consenso?
- Compare y contraste Paxos con el algoritmo Raft. ¿Cuáles son las diferencias clave en su funcionamiento y uso?
- Explica el concepto de una caché distribuida. ¿Cuáles son los principales desafíos en la implementación de una caché distribuida?
- Describe una estrategia de consistencia eventual en cachés distribuidas. ¿Cómo se maneja la replicación de datos y qué garantías de consistencia se ofrecen?
- Analiza un escenario en el que la consistencia eventual podría ser insuficiente. ¿Qué tipo de aplicaciones requieren una consistencia más fuerte?
- Explica los diferentes modelos de consistencia en la replicación de datos. ¿Cuáles son las diferencias entre consistencia fuerte, consistencia eventual y consistencia causal?
- Describe un sistema de replicación de datos con consistencia eventual. ¿Cómo se manejan las actualizaciones y cómo se resuelven los conflictos?
- Analice las ventajas y desventajas de la replicación de datos con consistencia eventual frente a la consistencia fuerte. ¿En qué tipo de aplicaciones sería más adecuada cada una?
- Supongamos un sistema multiprocesador con cuatro núcleos (P1, P2, P3 y P4) que utilizan un protocolo MESI. Describa paso a paso cómo se manejaría una operación de escritura en una dirección de memoria que inicialmente no está en la caché de ningún procesador.

- Considera un sistema distribuido con nodos que implementan el algoritmo de consenso Paxos. Describe el proceso completo de una propuesta que pasa por los estados de promesa y aceptación hasta alcanzar el consenso.
- Analiza un escenario en el que un sistema de directorio distribuido podría fallar al mantener la coherencia de caché. Plantea soluciones para mitigar estos problemas y mejorar la robustez del sistema.

### # Tus respuestas

1. Implementa una versión simplificada del algoritmo de consenso Paxos en Python.

```
# Código base
import random

class Proposer:
    def __init__(self, id, acceptors):
        self.id = id
        self.acceptors = acceptors

    def propose(self, value):
        proposal_number = random.randint(1, 100)
        promises = 0
        for acceptor in self.acceptors:
            if acceptor.promise(proposal_number):
                promises += 1

        if promises > len(self.acceptors) // 2:
            for acceptor in self.acceptors:
                acceptor.accept(proposal_number, value)
            return True
        return False

class Acceptor:
    def __init__(self, id):
        self.id = id
        self.promised_proposal = 0
        self.accepted_proposal = 0
        self.accepted_value = None

    def promise(self, proposal_number):
        if proposal_number > self.promised_proposal:
            self.promised_proposal = proposal_number
            return True
        return False

    def accept(self, proposal_number, value):
        if proposal_number >= self.promised_proposal:
            self.accepted_proposal = proposal_number
            self.accepted_value = value
```

```

# Ejemplo de uso
acceptors = [Acceptor(i) for i in range(3)]
proposer = Proposer(1, acceptors)

if proposer.propose("value1"):
    print("Proposal accepted")
else:
    print("Proposal rejected")

## Tu respuesta

```

2. Implementa una caché distribuida simple en Python utilizando un mecanismo básico de replicación de datos.

```

class DistributedCache:
    def __init__(self):
        self.nodes = {}

    def add_node(self, node_id):
        self.nodes[node_id] = {}

    def put(self, self, key, value):
        for node in self.nodes.values():
            node[key] = value

    def get(self, self, node_id, key):
        return self.nodes[node_id].get(key, None)

# Ejemplo de uso
cache = DistributedCache()
cache.add_node('node1')
cache.add_node('node2')

cache.put('key1', 'value1')
print(cache.get('node1', 'key1'))
print(cache.get('node2', 'key1'))

## Tu respuesta

```

3. Implementa un sistema de replicación de datos con consistencia eventual en Python.

```

class Node:
    def __init__(self, id):
        self.id = id
        self.data = {}
        self.log = []

    def put(self, self, key, value):
        self.data[key] = value
        self.log.append((key, value))

```

```
def get(self, key):
    return self.data.get(key, None)

def replicate(self, other_node):
    for entry in self.log:
        other_node.data[entry[0]] = entry[1]
    other_node.log.extend(self.log)

# Ejemplo de uso
node1 = Node('node1')
node2 = Node('node2')

node1.put('key1', 'value1')
node1.put('key2', 'value2')

node1.replicate(node2)

print(node2.get('key1'))
print(node2.get('key2'))

## Tu respuesta
```