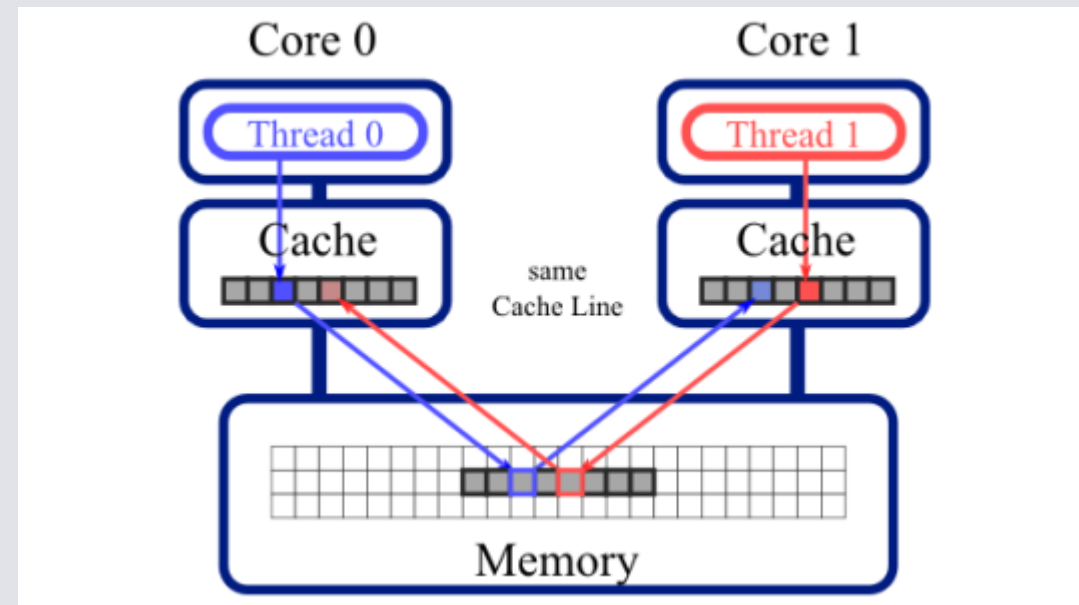# Analyzing False Sharing

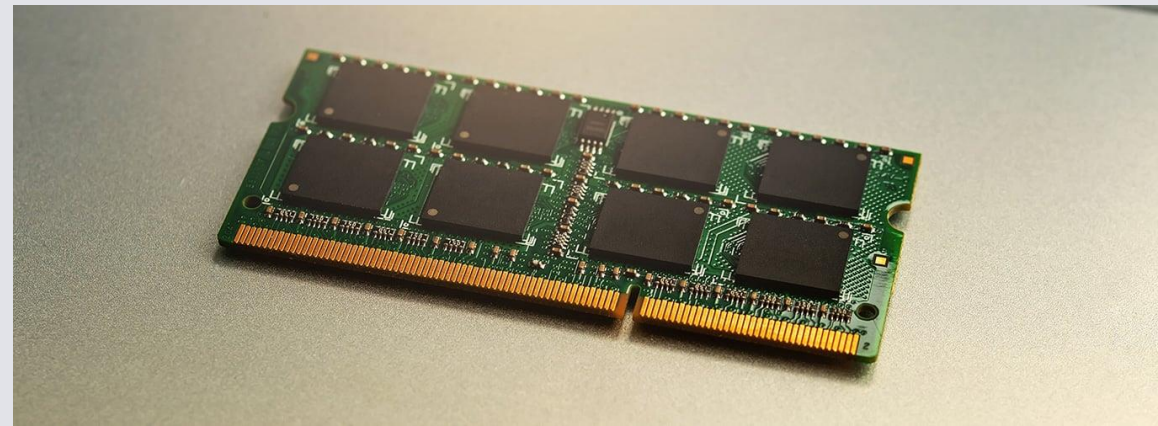Team Members: Alon Raigorodetsky, Puja Girishkumar Kakani

# What is False Sharing?

- 2 threads accessing the same data

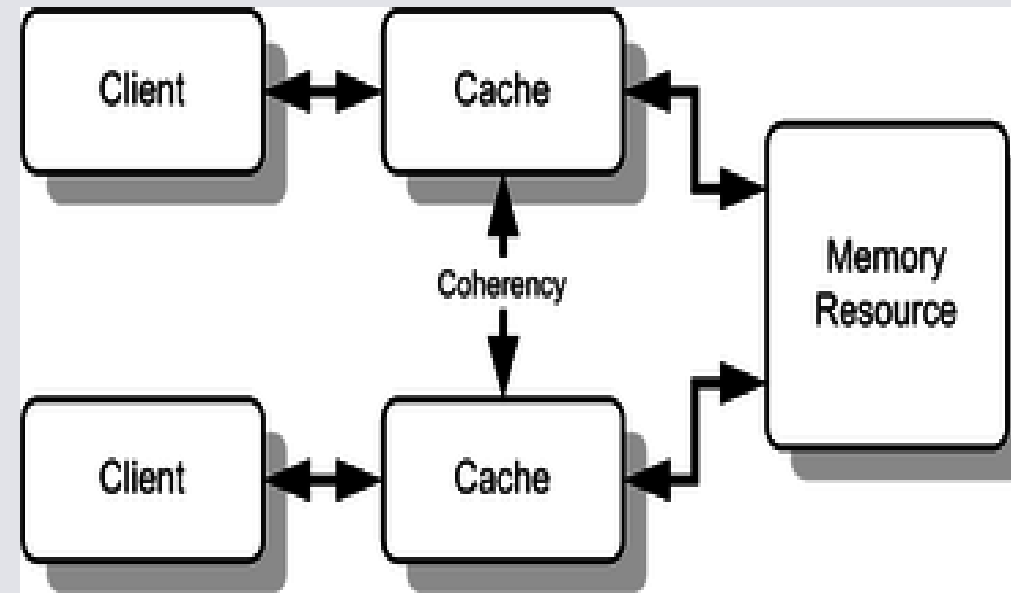- Modification of the data that is being accessed.

# What is Cache?

- A cache is a high speed data storage layer that stores frequently accessed data which is used by the operating system

- Cache can reduce time or resources that take to perform and action but sometimes it can cause a stale if not updated properly
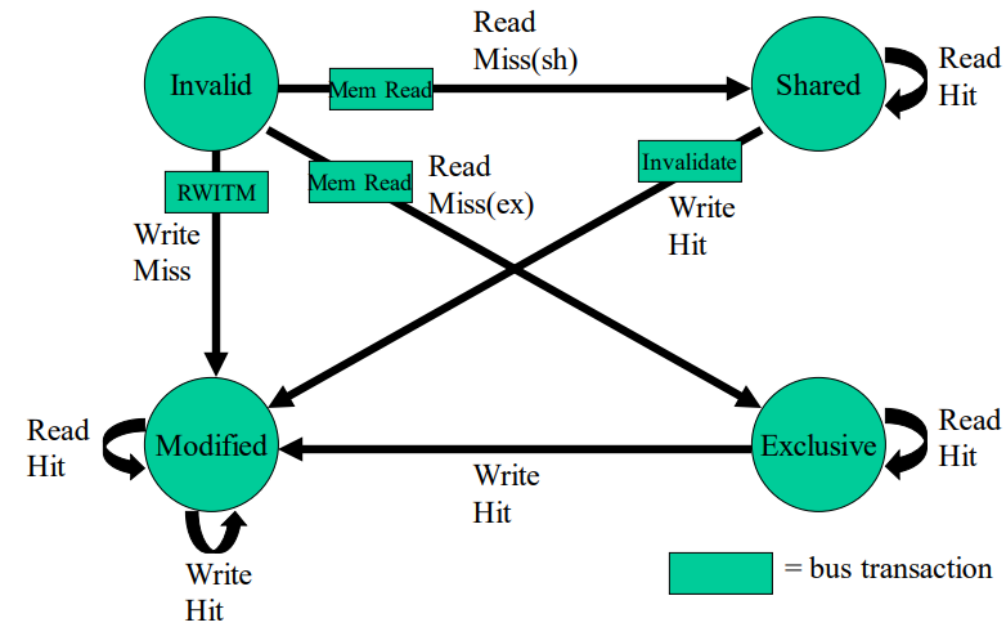
# What is Cache Consistency?

- Cache consistency refers to the synchronization of data between different caches in a computer system, ensuring that all caches have the most up-to-date version of data.

- Without Cache Consistency, different processors may have different versions of the same data, which can lead to data inconsistencies and errors in the system.

# The MESI Protocol

- **Modified** - The only cached duplicate is the modified, distinct from main memory cache line.
- **Exclusive** - The modified, separate from main memory cache line is the only cached copy.
- **Shared** - identical to main memory, but there might be duplicates in other caches.
- **Invalid** - Line data is not valid.

# How the code works?

```cpp
1  #include <iostream>
2  #include <thread>
3  #include <chrono>
4
5  using namespace std;
6  using namespace chrono;
7
8  // Constants to control the program
9  const int NUM_THREADS = 2;
10 const int NUM_ITERATIONS = 10000000000000;
11 const int CACHE_LINE_SIZE = 64;
12
13 // Define a counter struct to ensure cache line alignment
14 struct alignas(CACHE_LINE_SIZE) Counter {
15     volatile long long value;  // the counter value
16     char padding[CACHE_LINE_SIZE - sizeof(long long)];  // padding to align the struct to cache line size
17 };
18
19 // Define two counter variables
20 Counter counter1, counter2;
21
22 // Function to increment counter 1
23 void increment1() {
24     for (int i = 0; i < NUM_ITERATIONS; i++) {
25         counter1.value++;
26     }
27 }
28
29 // Function to increment counter 2
30 void increment2() {
31     for (int i = 0; i < NUM_ITERATIONS; i++) {
32         counter2.value++;
33     }
34 }
35
```

# How the Code Works?

```cpp
35
36      // Function to check the runtime of the program
37    ⊟void checkRuntime(high_resolution_clock::time_point start_time, high_resolution_clock::time_point end_time) {
38          auto duration = duration_cast<milliseconds>(end_time - start_time).count();
39          cout << "Runtime: " << duration << " ms" << endl;
40    └}
41
42    ⊟int main() {
43          // Print the cache line size
44          cout << "Cache line size: " << CACHE_LINE_SIZE << " bytes" << endl;
45          // Run the program using a single thread
46          cout << "Running program using a single thread" << endl;
47
48          high_resolution_clock::time_point start_time = high_resolution_clock::now();
49    ⊟     for (int i = 0; i < NUM_ITERATIONS; i++) {
50              counter1.value++;
51              counter2.value++;
52          }
53          high_resolution_clock::time_point end_time = high_resolution_clock::now();
54          checkRuntime(start_time, end_time);
55          cout << "Counter 1: " << counter1.value << endl;
56          cout << "Counter 2: " << counter2.value << endl;
57
58          // Run the program using multiple threads
59          cout << "Running program using " << NUM_THREADS << " threads" << endl;
60          counter1.value = 0;
61          counter2.value = 0;
62
63          start_time = high_resolution_clock::now();
64
65          thread t1(increment1);
66          thread t2(increment2);
67          t1.join();
68          t2.join();
69
70          end_time = high_resolution_clock::now();
71          checkRuntime(start_time, end_time);
72          cout << "Counter 1: " << counter1.value << endl;
73          cout << "Counter 2: " << counter2.value << endl;
74
75          return 0;
76    └}
```

# Output



```
Cache line size: 64 bytes
Running program using a single thread
Runtime: 2543 ms
Counter 1: 1215752192
Counter 2: 1215752192
Running program using 2 threads
Runtime: 2529 ms
Counter 1: 1215752192
Counter 2: 1215752192
```

```
Cache line size: 64 bytes
Running program using a single thread
Runtime: 214 ms
Counter 1: 100000000
Counter 2: 100000000
Running program using 2 threads
Runtime: 229 ms
Counter 1: 100000000
Counter 2: 100000000
```

# Solution

- A performance hit results as a result of each thread invalidating the cache line for the other thread. In order to ensure that each variable is aligned to its own cache line, we can fix this issue by adding padding to the Counter struct.

- By ensuring that each thread updates a different cache line, false sharing is prevented.

# Code for the Solution

```cpp
#include <iostream>
#include <thread>
#include <chrono>
#include <atomic>

using namespace std;
using namespace chrono;

// Constants to control the program
const int NUM_THREADS = 2;
const int NUM_ITERATIONS = 100000000;
const int CACHE_LINE_SIZE = 64;

// Define a counter struct with padding for cache line alignment
struct alignas(CACHE_LINE_SIZE) Counter {
    atomic<long long> value;  // the counter value
    char padding[CACHE_LINE_SIZE - sizeof(atomic<long long>)];  // padding to align the struct to cache line size
};

// Define two counter variables
Counter counter1, counter2;

// Function to increment counter 1
void increment1() {
    for (int i = 0; i < NUM_ITERATIONS; i++) {
        counter1.value++;
    }
}

// Function to increment counter 2
void increment2() {
    for (int i = 0; i < NUM_ITERATIONS; i++) {
        counter2.value++;
    }
}

// Function to check the runtime of the program
void checkRuntime(high_resolution_clock::time_point start_time, high_resolution_clock::time_point end_time) {
```

# Code for the Solution

```cpp
void checkRuntime(high_resolution_clock::time_point start_time, high_resolution_clock::time_point end_time) {
    auto duration = duration_cast<milliseconds>(end_time - start_time).count();
    cout << "Runtime: " << duration << " ms" << endl;
}

int main() {
    // Print the cache line size
    cout << "Cache line size: " << CACHE_LINE_SIZE << " bytes" << endl;

    // Run the program using a single thread
    cout << "Running program using a single thread" << endl;
    high_resolution_clock::time_point start_time = high_resolution_clock::now();
    for (int i = 0; i < NUM_ITERATIONS; i++) {
        counter1.value++;
        counter2.value++;
    }
    high_resolution_clock::time_point end_time = high_resolution_clock::now();
    checkRuntime(start_time, end_time);
    cout << "Counter 1: " << counter1.value << endl;
    cout << "Counter 2: " << counter2.value << endl;

    // Run the program using multiple threads
    cout << "Running program using " << NUM_THREADS << " threads" << endl;
    counter1.value = 0;
    counter2.value = 0;

    start_time = high_resolution_clock::now();
    thread t1(increment1);
    thread t2(increment2);
    t1.join();
    t2.join();
    end_time = high_resolution_clock::now();
    checkRuntime(start_time, end_time);
    cout << "Counter 1: " << counter1.value << endl;
    cout << "Counter 2: " << counter2.value << endl;

    return 0;
}
```

# Code Output

```
Cache line size: 64 bytes
Running program using a single thread
Runtime: 355 ms
Counter 1: 100000000
Counter 2: 100000000
Running program using 2 threads
Runtime: 479 ms
Counter 1: 100000000
Counter 2: 100000000
```

# Summary

- False sharing is a performance problem that can happen in multi-threaded applications when different variables on the same cache line are accessed by different threads at the same time.

- Finding the cache lines that are shared by several threads and then identifying the variables that cause the false sharing are the first steps in false sharing analysis.

- In conclusion, watch out for false sharing because it kills scalability.

- Prevent false sharing by lowering the frequency of updates to the variables that are shared.

# Do you have any questions?