# Cache Friendly Programming

# Oğuzhan KATLI

https://www.linkedin.com/in/ogzhnktl

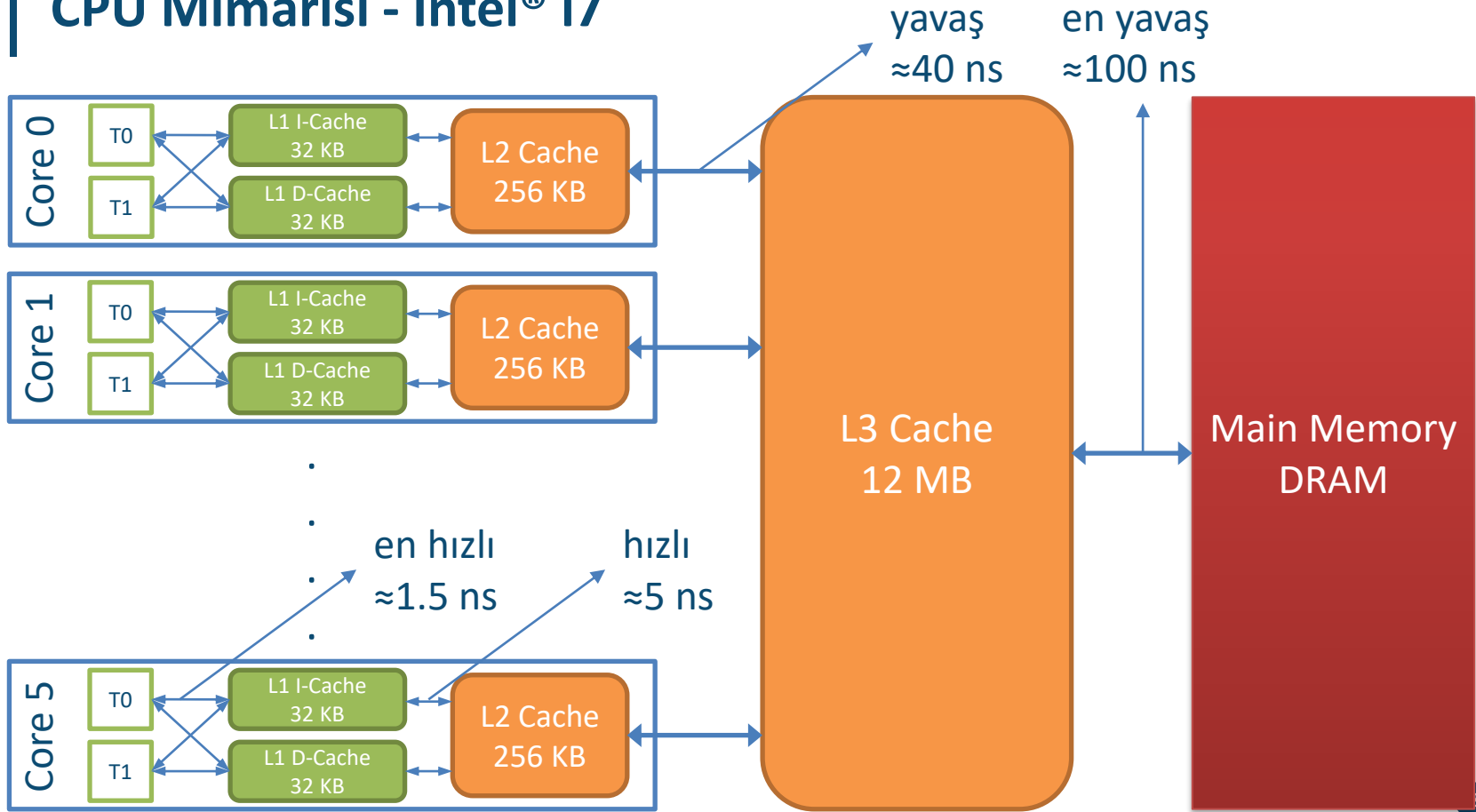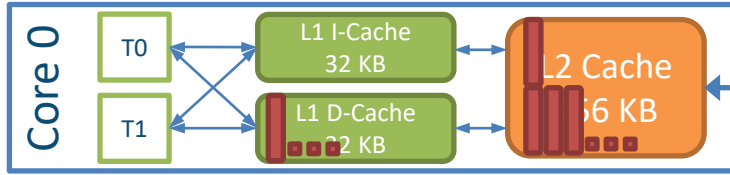https://github.com/nixiz

# Benchmark Setup

- Intel® Core™ i7-9750H İşlemci (12M Önbellek, 4,50 GHz'e kadar)
  - 6 Core – 12 Threads
  - 32 KB      L1 Cache
  - 256 KB    L2, 12 MB L3 Cache
- Windows 10 Home
  - Visual Studio 2019 Version 16.7.2 – MSVC 19.27
- Ubuntu 20.04 LTS
  - GCC 9.3.0
  - Clang 10.0
- Google Benchmark: https://github.com/google/benchmark
- Benchmark Kodları:
  https://gist.github.com/nixiz/c3bb8f6029b64f9281ac44cbc119209a
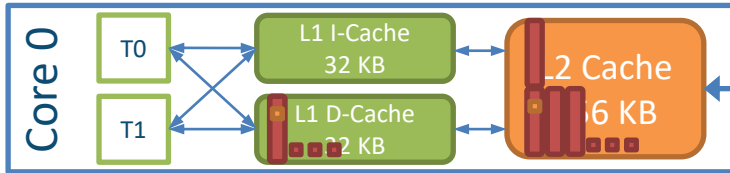
# CPU Mimarisi - Intel® i7

# CPU Mimarisi – Cache Line

Okuma:



Yazma:

# Veri Erişim Hızları

## Latency Numbers Every Programmer Should Know

■ 1 ns

■ L1 cache reference: 0.5 ns

■ Branch mispredict: 5 ns

■ L2 cache reference: 7 ns

■ Mutex lock/unlock: 25 ns

= 100 ns

■ Main memory reference: 100 ns

= 1 µs

Compress 1 KB with Zippy: 3 µs

= 10 µs

■ Send 1 KB over 1 Gbps network: 10 µs

SSD random read (1Gb/s SSD): 150 µs

Read 1 MB sequentially from memory: 250 µs

Round trip in same datacenter: 500 µs

= 1 ms

■ Read 1 MB sequentially from SSD: 1 ms

Disk seek: 10 ms

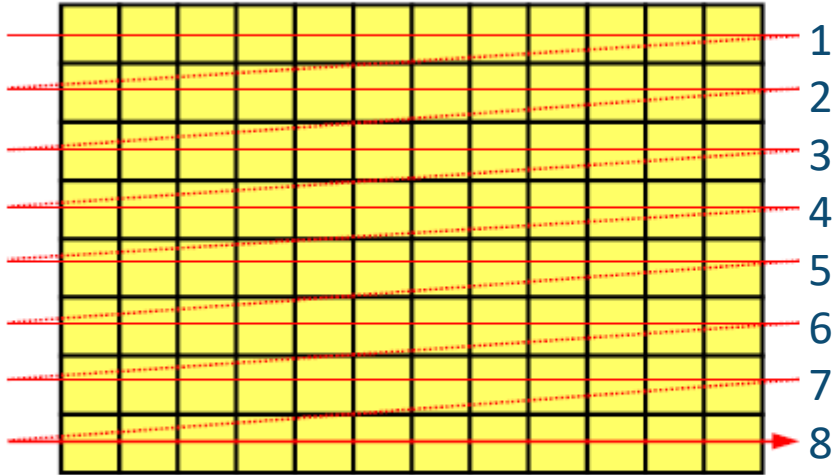Read 1 MB sequentially from disk: 20 ms

Packet roundtrip CA to Netherlands: 150 ms
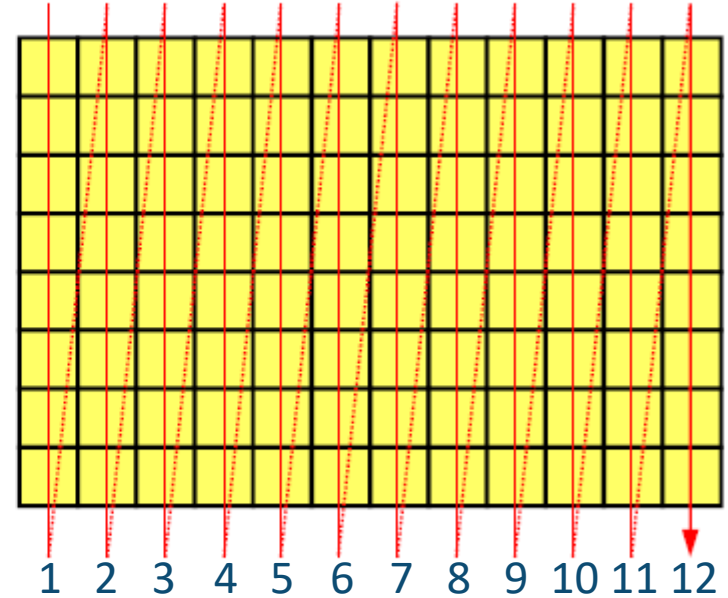
# Cache Friendly Programming

# Row vs Column
# Matrix Traversal

# Row vs Column Major
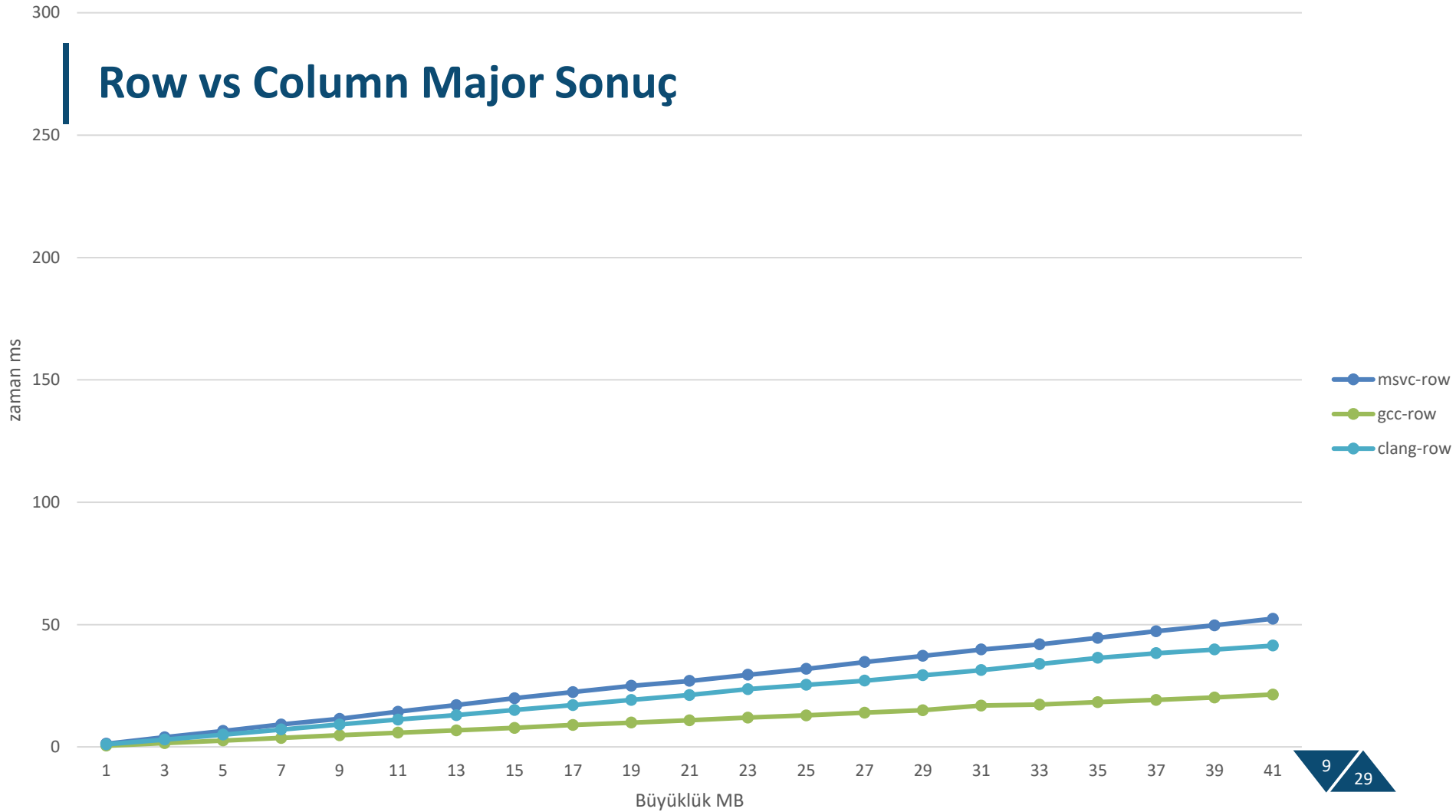


Row Major

Column Major

# Row vs Column Major

```cpp
template <typename T>
auto sumMatrix(const Matrix<T>& m, TraversalOrder order)
{
  unsigned long sum = 0;
  if (order == TraversalOrder::RowMajor)
  {
    for (auto r = 0; r < m.rows(); ++r)
      for (auto c = 0; c < m.columns(); ++c)
        sum += m[r][c];
  }
  else      // TraversalOrder::ColumnMajor
  {
    for (auto c = 0; c < m.columns(); ++c)
      for (auto r = 0; r < m.rows(); ++r)
        sum += m[r][c];
  }
  return sum;
}
```

Row vs Column Major Sonuç

# Row vs Column Major Sonuç

zaman ms

Büyüklük MB

- msvc-row
- msvc-column
- gcc-row
- gcc-column
- clang-row
- clang-column

Row vs Column Major Sonuç

5 MB

mean-row
mean-column
fark

Büyüklük MB

zaman ms

# Cache Friendly Programming

# Branch Prediction

# Branch Prediction

| v: | -1 | -1 | 1 | 1 | -1 | 1 | -1 | -1 | 1 | -1 | -1 | 1 | 1 | -1 | 1 |
|----|----|----|---|---|----|---|----|----|---|----|----|---|---|----|---|

| v: | 1 | 1 | 1 | 1 | 1 | 1 | 1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 | -1 |
|----|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|

```cpp
std::vector<float> v(65536);
std::generate(std::begin(v), std::end(v), [] {
    return (rand() % 2) ? 1 : -1;
});

std::sort(v.begin(), v.end());

return std::count_if(std::begin(v), std::end(v), [](float x) {
    return x > 0;
});
```
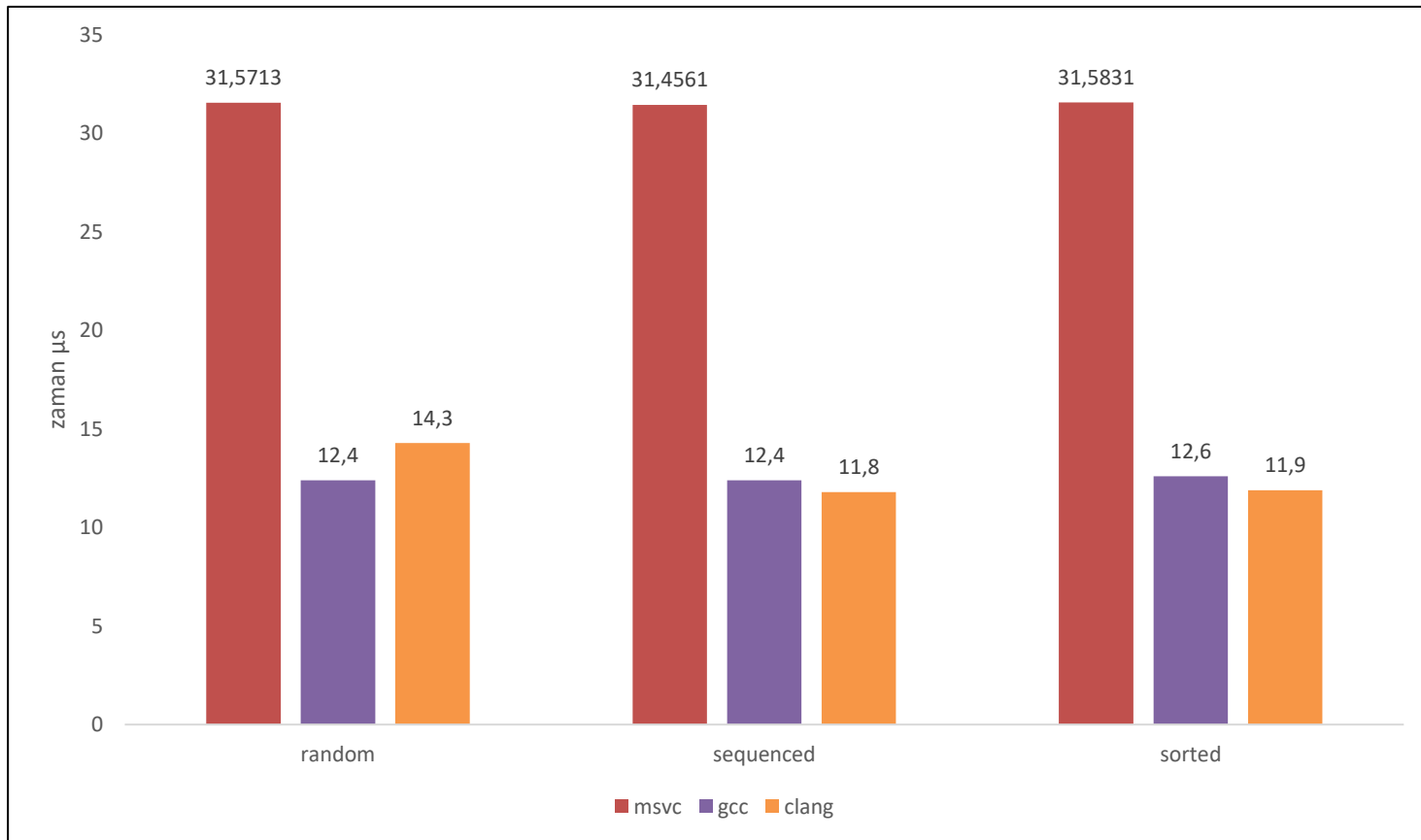
# Branch Prediction – Sequenced If - Else

| v: | 1 | -1 | 1 | -1 | 1 | -1 | 1 | -1 | 1 | -1 | 1 | -1 | 1 | -1 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

```cpp
std::vector<float> v(65536);
std::generate(std::begin(v), std::end(v), [n = 0]() mutable {
    return (++n % 2) ? 1 : -1;
  });


return std::count_if(std::begin(v), std::end(v), [](float x) {
  return x > 0;
});
```

# Count If Sonuç

# Virtual Function Calls

```cpp
struct base_price
{
  virtual ~price() {}
  virtual float getPrice() const noexcept { return 0.0; }
};

struct cheap : public base_price
{
  float getPrice() const noexcept override { return 2.0; }
};

struct expensive : public base_price
{
  float getPrice() const noexcept override { return 3.14159; }
};
```
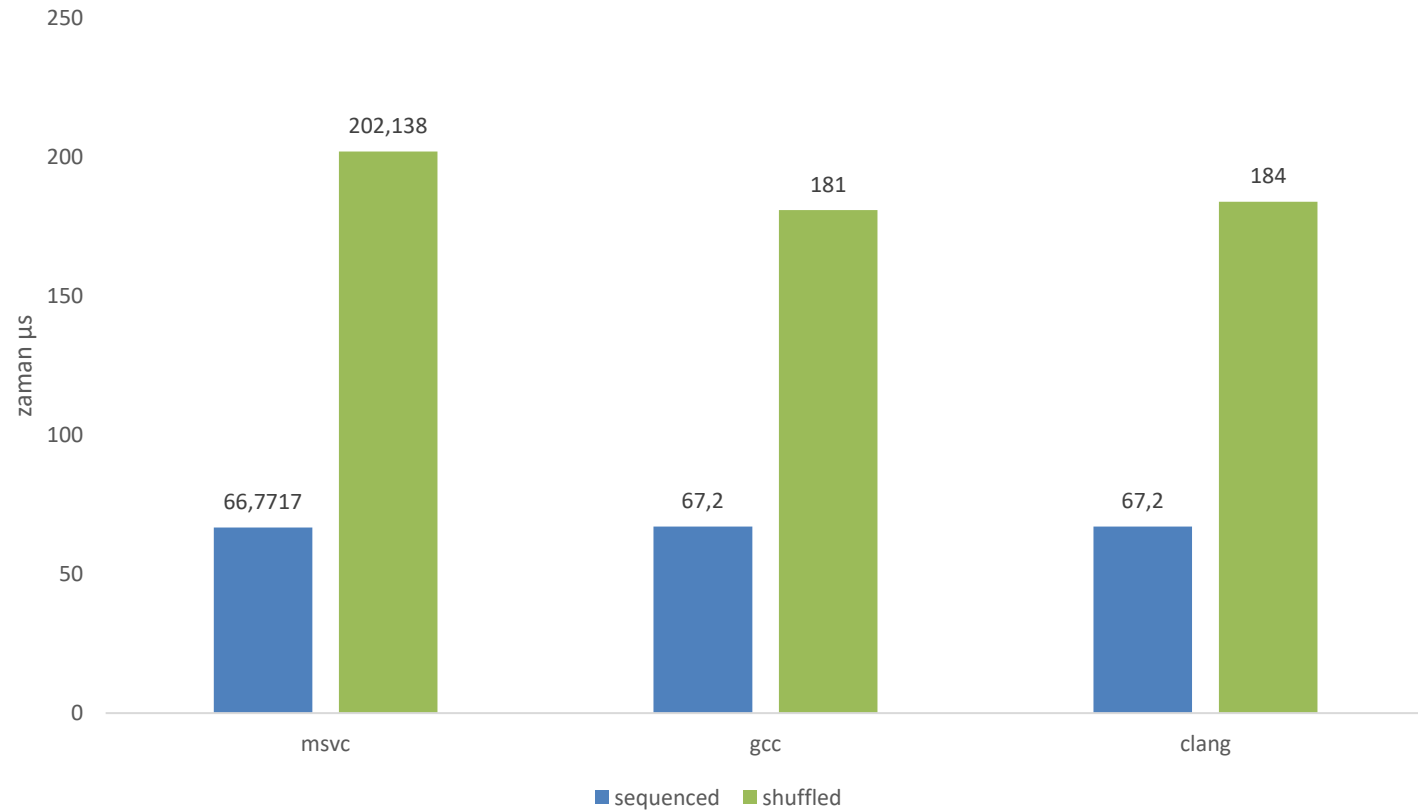
# Örnek

```cpp
std::vector<base_price*> pricelist;
std::fill_n(std::back_inserter(pricelist), 10000, new base_price);
std::fill_n(std::back_inserter(pricelist), 10000, new cheap);
std::fill_n(std::back_inserter(pricelist), 10000, new expensive);

std::random_shuffle(pricelist.begin(), pricelist.end());

float sum = 0;
for (auto *p : pricelist)
{
  sum += p->getPrice();
}
```

# Virtual Call Sonuç

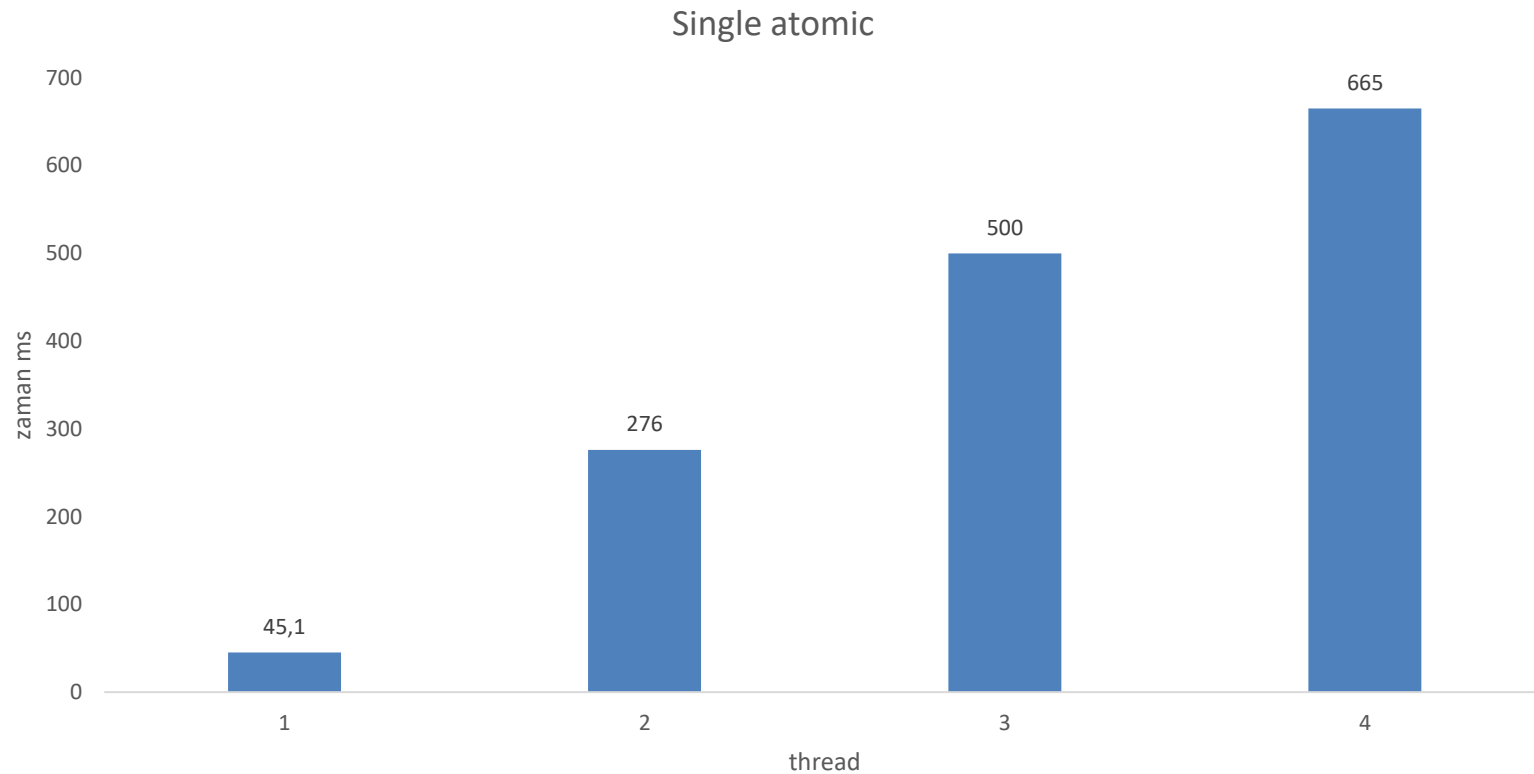# Cache Friendly Programming

# False Sharing

# False Sharing

```cpp
template <typename input_t>
void work(input_t& a)
{
  for (int i = 0; i < 10,000,000; ++i)
    a++;
}

int test_with_4_threads()
{
  std::atomic<int> a; a = 0;

  std::thread t1([&] { work(a); });
  std::thread t2([&] { work(a); });
  std::thread t3([&] { work(a); });
  std::thread t4([&] { work(a); });

  t1.join(); t2.join(); t3.join(); t4.join();
  return a;
}
```
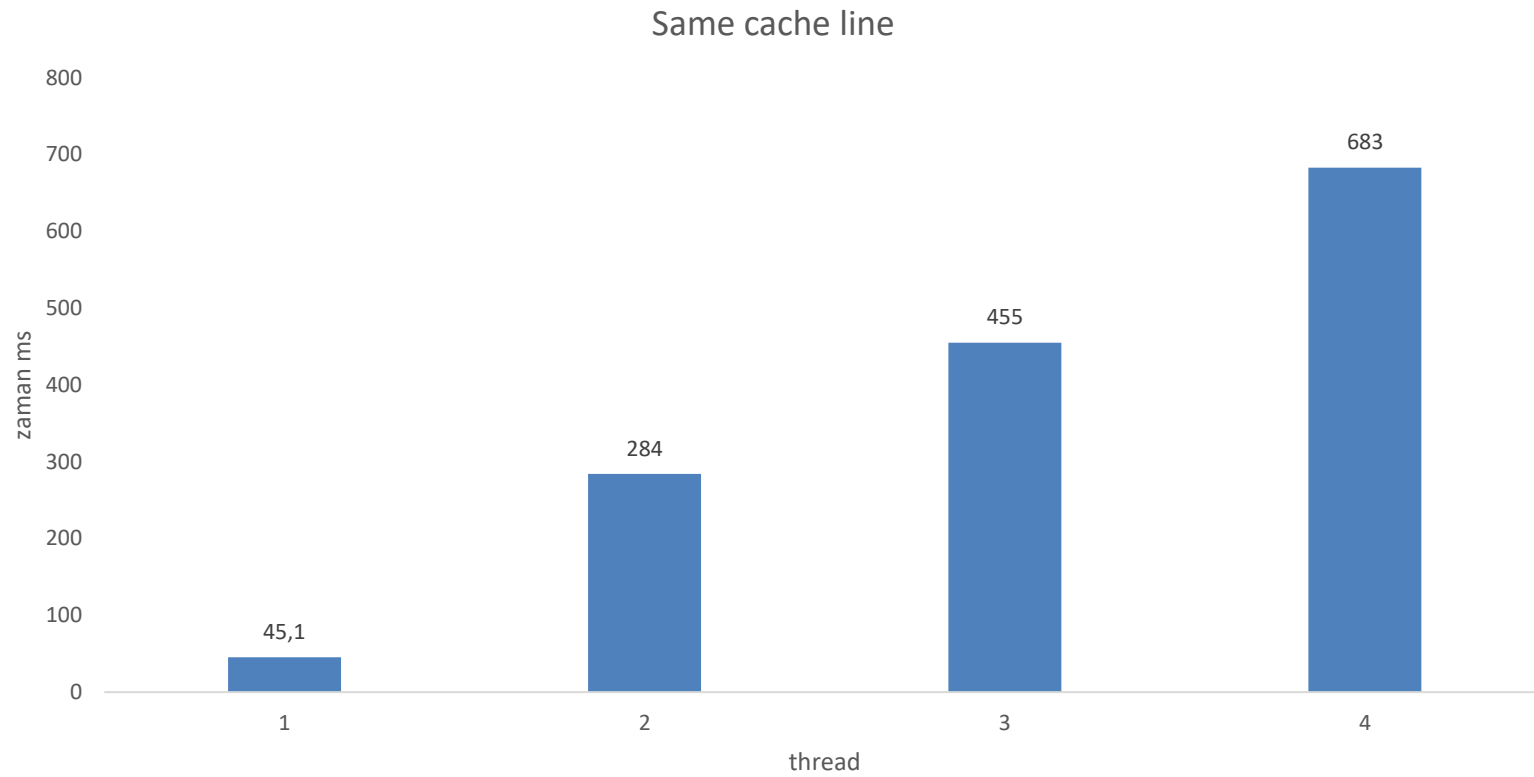
# Sonuç

Single atomic

# Farklı Atomik Nesneler

```cpp
int test()
{
  std::atomic<int> a; a = 0;
  std::atomic<int> b; b = 0;
  std::atomic<int> c; c = 0;
  std::atomic<int> d; d = 0;

  std::thread t1([&] { work(a); });
  std::thread t2([&] { work(b); });
  std::thread t3([&] { work(c); });
  std::thread t4([&] { work(d); });

  t1.join(); t2.join(); t3.join(); t4.join();
  return a + b + c + d;
}
```

# Sonuç

## Same cache line



Chart showing "zaman ms" (y-axis) versus "thread" (x-axis):
- Thread 1: 45,1
- Thread 2: 284
- Thread 3: 455
- Thread 4: 683

# Farklı Atomik Nesneler

Adres blokları <u>aynı</u>

```cpp
int test()
{
  std::atomic<int> a; a = 0;  // &a: 0x...b2f7c0
  std::atomic<int> b; b = 0;  // &b: 0x...b2f7c4
  std::atomic<int> c; c = 0;  // &c: 0x...b2f7c8
  std::atomic<int> d; d = 0;  // &d: 0x...b2f7cc

  std::thread t1([&] { work(a); });
  std::thread t2([&] { work(b); });
  std::thread t3([&] { work(c); });
  std::thread t4([&] { work(d); });

  t1.join(); t2.join(); t3.join(); t4.join();
  return a + b + c + d;
}
```

# Çözüm

Adres blokları **Farklı**

```cpp
struct alignas(64) aligned_type
{
  std::atomic<int> val;
};
```

```cpp
int test()
{
    aligned_type a; a.val = 0; // &a: 0x...4ff240
    aligned_type b; b.val = 0; // &b: 0x...4ff280
    aligned_type c; c.val = 0; // &c: 0x...4ff2c0
    aligned_type d; d.val = 0; // &d: 0x...4ff300

    std::thread t1([&a] { work(a.val); });
    std::thread t2([&b] { work(b.val); });
    std::thread t3([&c] { work(c.val); });
    std::thread t4([&d] { work(d.val); });

    t1.join(); t2.join(); t3.join(); t4.join();
    return a.val + b.val + c.val + d.val;
}
```
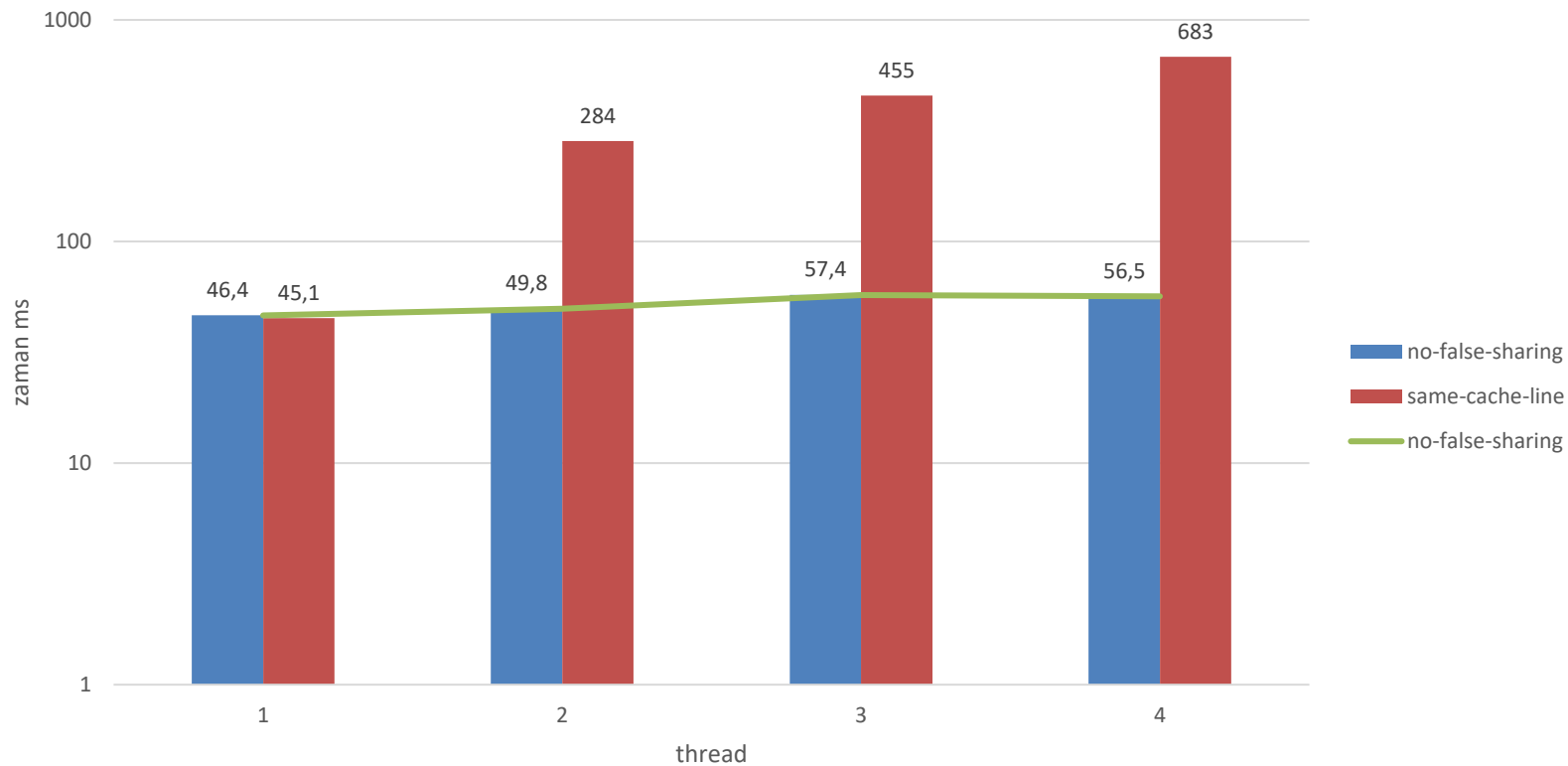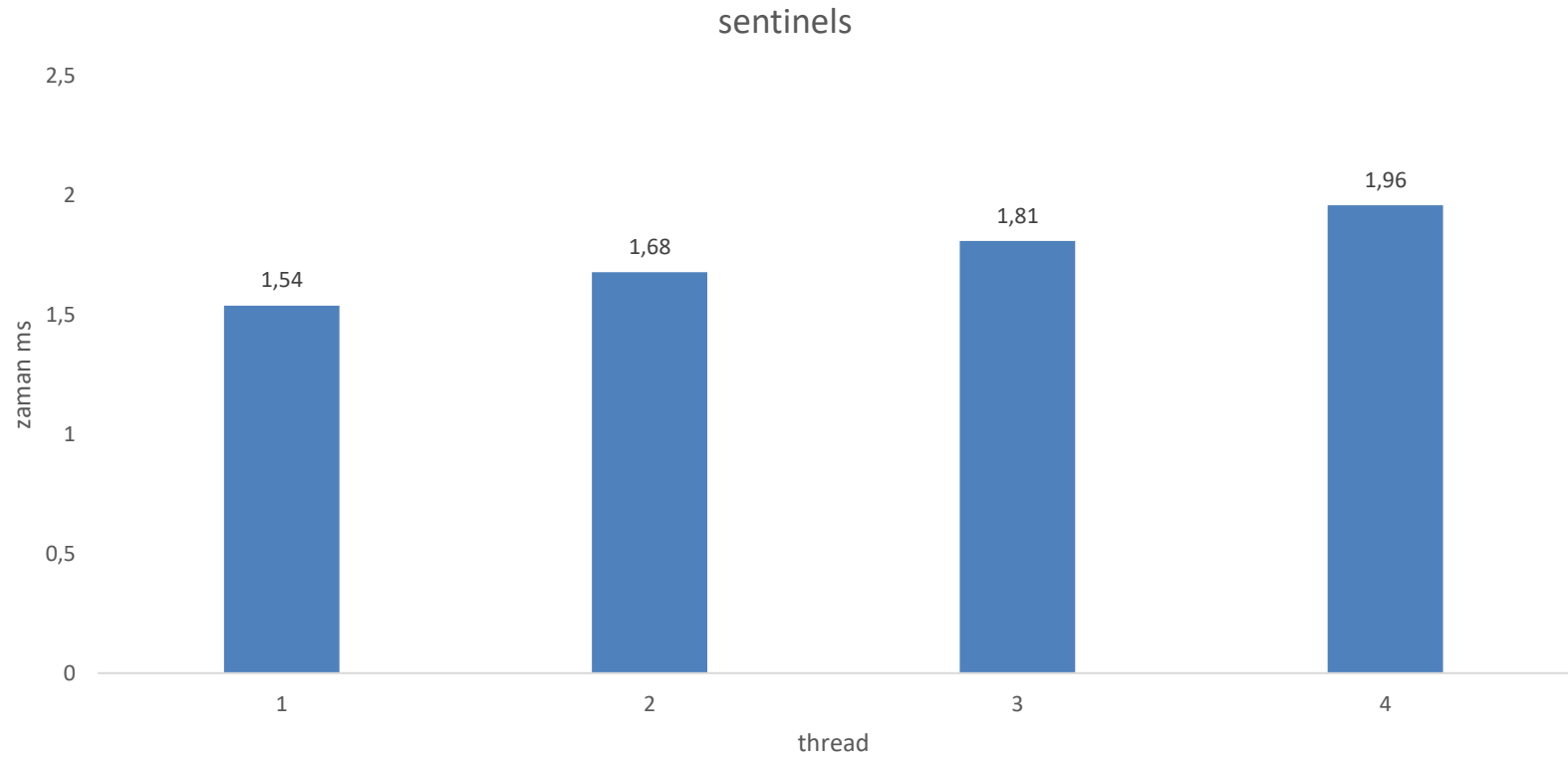
# Sonuç

# False Sharing – Sentinels

```cpp
void work_with_sentinel(std::atomic<int>& a) {
  thread_local unsigned int sentinel = 0;
  for (int i = 0; i < 10,000,000; ++i)
    ++sentinel;
  a += sentinel;
}

int test_with_sentinels() {
  std::atomic<int> a; a = 0;

  std::thread t1([&] { work_with_sentinel(a); });
  std::thread t2([&] { work_with_sentinel(a); });
  std::thread t3([&] { work_with_sentinel(a); });
  std::thread t4([&] { work_with_sentinel(a); });

  t1.join(); t2.join(); t3.join(); t4.join();
  return a;
}
```

# Sonuç



sentinels

Bar chart with x-axis labeled "thread" showing values 1, 2, 3, 4 and y-axis labeled "zaman ms" ranging from 0 to 2,5.

| thread | zaman ms |
|--------|----------|
| 1 | 1,54 |
| 2 | 1,68 |
| 3 | 1,81 |
| 4 | 1,96 |

# Cache Friendly Programming

# Teşekkürler