# CS305: Computer Architecture

## Caches-V

https://www.cse.iitb.ac.in/~biswa/courses/CS305/main.html
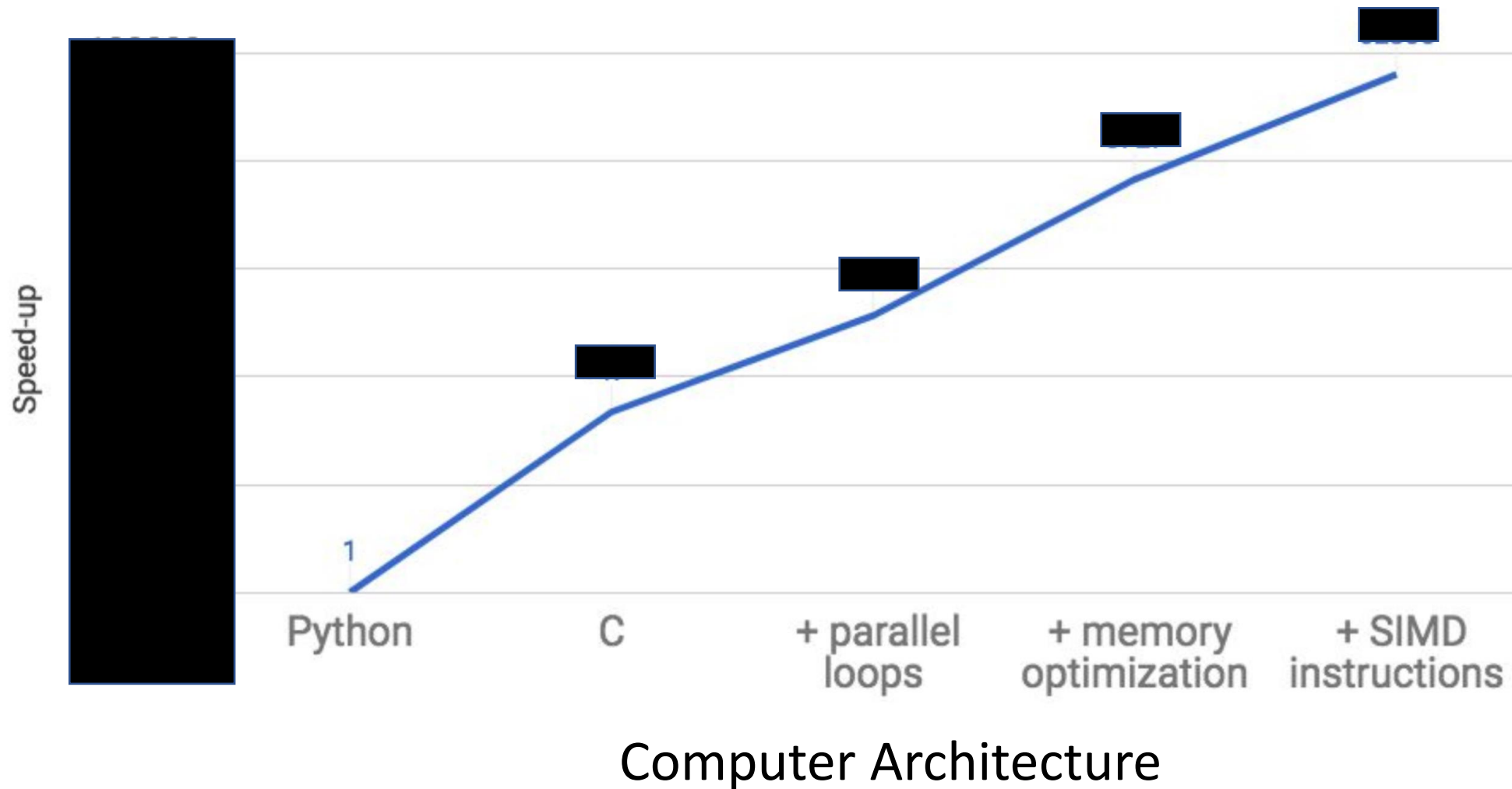
https://www.cse.iitb.ac.in/~biswa/

# Memory Mountain



Think about it, before you write your program

# Does programming languages matter?



Matrix Multiply Speedup Over Native Python

Computer Architecture

# Seriously?

Matrix Multiply Speedup Over Native Python



Speed-up

47

1

Python          C          + parallel          + memory          + SIMD
                           loops            optimization      instructions

Computer Architecture                                    4

# What?

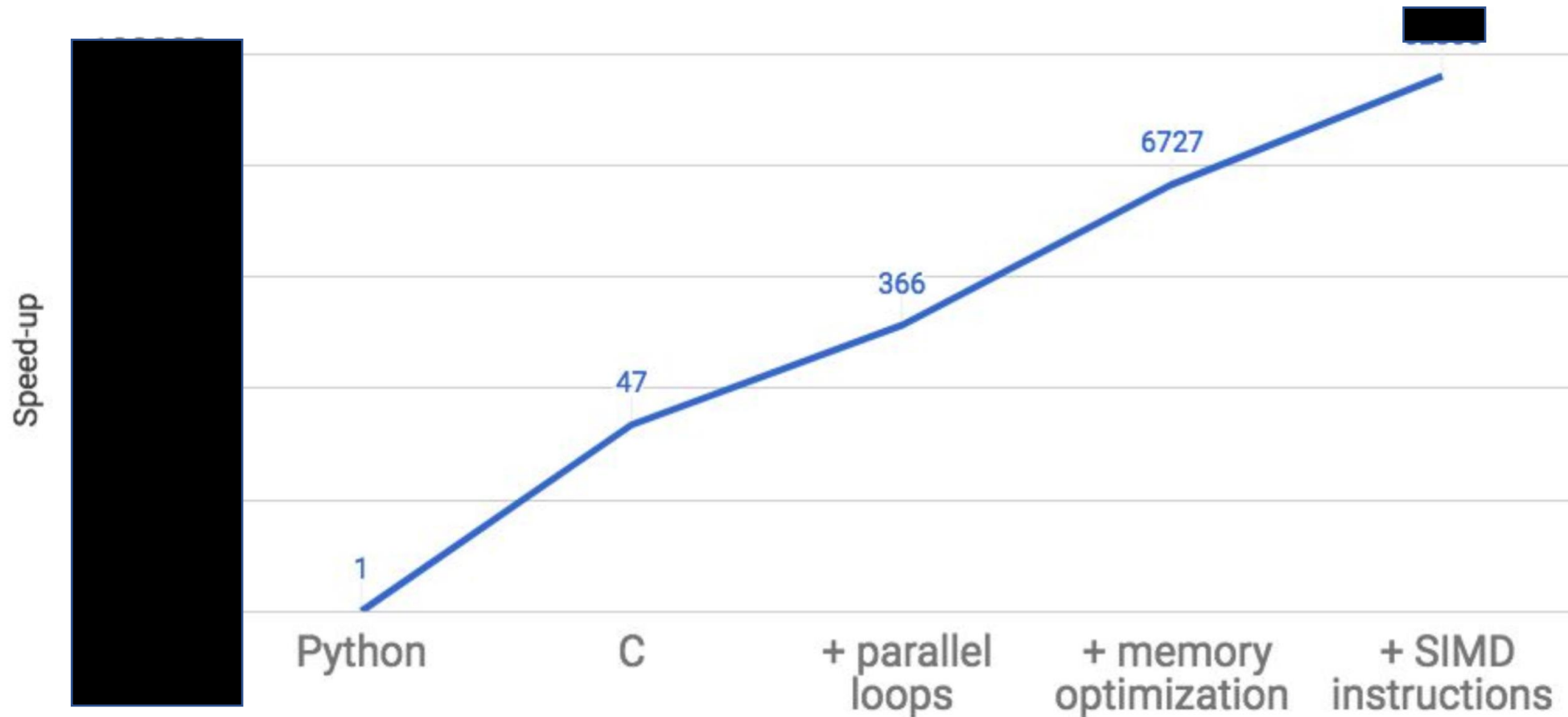

Matrix Multiply Speedup Over Native Python

Computer Architecture

# Insane



Matrix Multiply Speedup Over Native Python

Speed-up

1 — Python
47 — C
366 — + parallel loops
6727 — + memory optimization
+ SIMD instructions

Computer Architecture                                        6

# Still?

Matrix Multiply Speedup Over Native Python



Speed-up

| | 1 | 47 | 366 | 6727 | |
|---|---|---|---|---|---|
| Python | C | + parallel loops | + memory optimization | + SIMD instructions | |

# Ohhhh!!



Matrix Multiply Speedup Over Native Python

# Can Compilers/programmers exploit locality?

# Matrix Multiplication: 101

```
/* ijk */
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

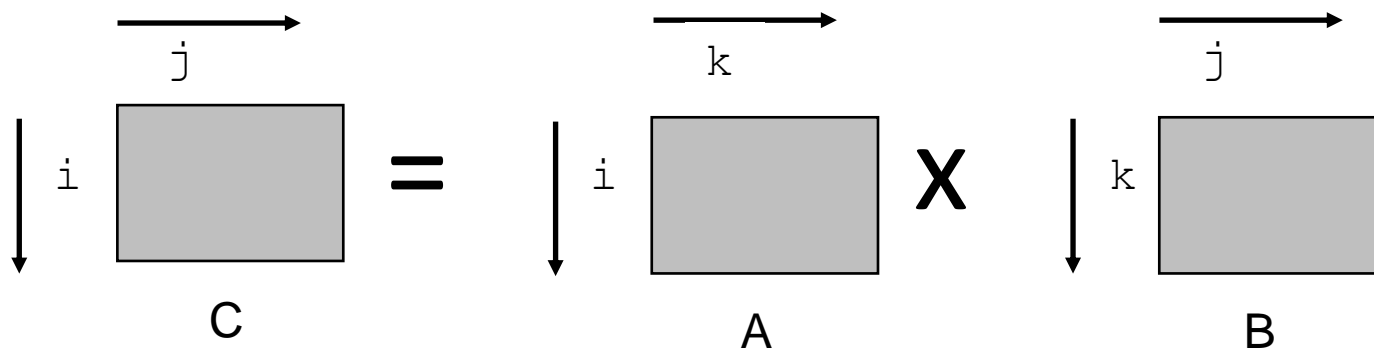$$4 \times 3 + 2 \times 2 + 7 \times 5 = 51$$

| 4 | 2 | 7 |
|---|---|---|
| 1 | 8 | 2 |
| 6 | 0 | 1 |

$\times$

| 3 | 0 | 1 |
|---|---|---|
| 2 | 4 | 5 |
| 5 | 9 | 1 |

$=$

| 51 | | |
|---|---|---|
| | | |
| | | |

Computer Architecture

# Miss Rate analysis

- Assume:
  - Block size = 32B (big enough for four doubles)
  - Matrix dimension (N) is very large
    - Approximate 1/N as 0.0
  - Cache is not even big enough to hold multiple rows

- Analysis Method:
  - Look at access pattern of inner loop



C = A X B

# Effect of Cache Layout

| C arrays allocated in row-major order | Stepping through columns in one row: | Stepping through rows in one column: |
|---|---|---|
| • each row in contiguous memory locations | • **for (i = 0; i < N; i++) sum += a[0][i];**<br>• accesses successive elements<br>• if block size (B) > sizeof($a_{ij}$) bytes, exploit spatial locality<br>   • miss rate = sizeof($a_{ij}$) / B | • **for (i = 0; i < N; i++) sum += a[i][0];**<br>• accesses distant elements<br>• no spatial locality!<br>   • miss rate = 1 (i.e. 100%) |

# Effect of loop order (ijk)

```
/* ijk */
for (i=0; i<n; i++)  {
   for (j=0; j<n; j++) {
      sum = 0.0;
      for (k=0; k<n; k++)
         sum += a[i][k] *
b[k][j];
      c[i][j] = sum;
   }
}
```
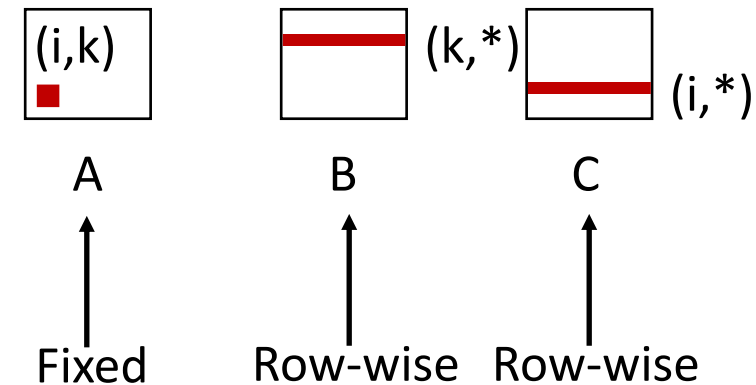
Inner loop:



A — Row-wise

B — Column-wise

C — Fixed

Miss rate for inner loop iterations:

| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

# Effect of loops (kij)

```
/* kij */
for (k=0; k<n; k++) {
    for (i=0; i<n; i++) {
        r = a[i][k];
        for (j=0; j<n; j++)
            c[i][j] += r * b[k][j];
    }
}
```

Inner loop:



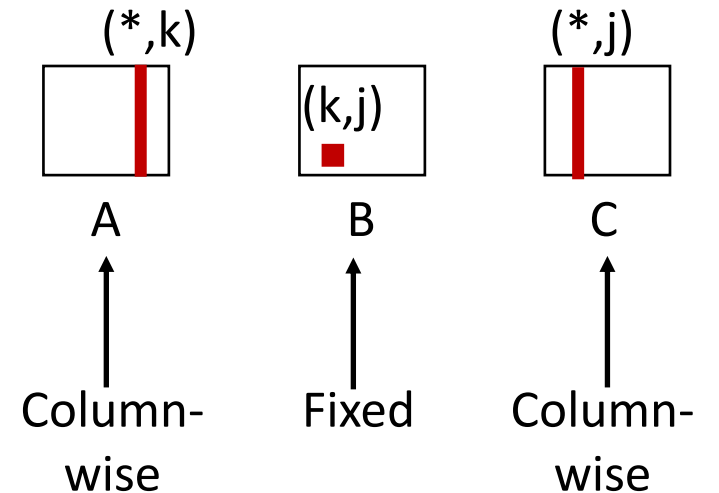A          B          C

Fixed    Row-wise  Row-wise

Miss rate for inner loop iterations:

| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

# Effect of loops (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```
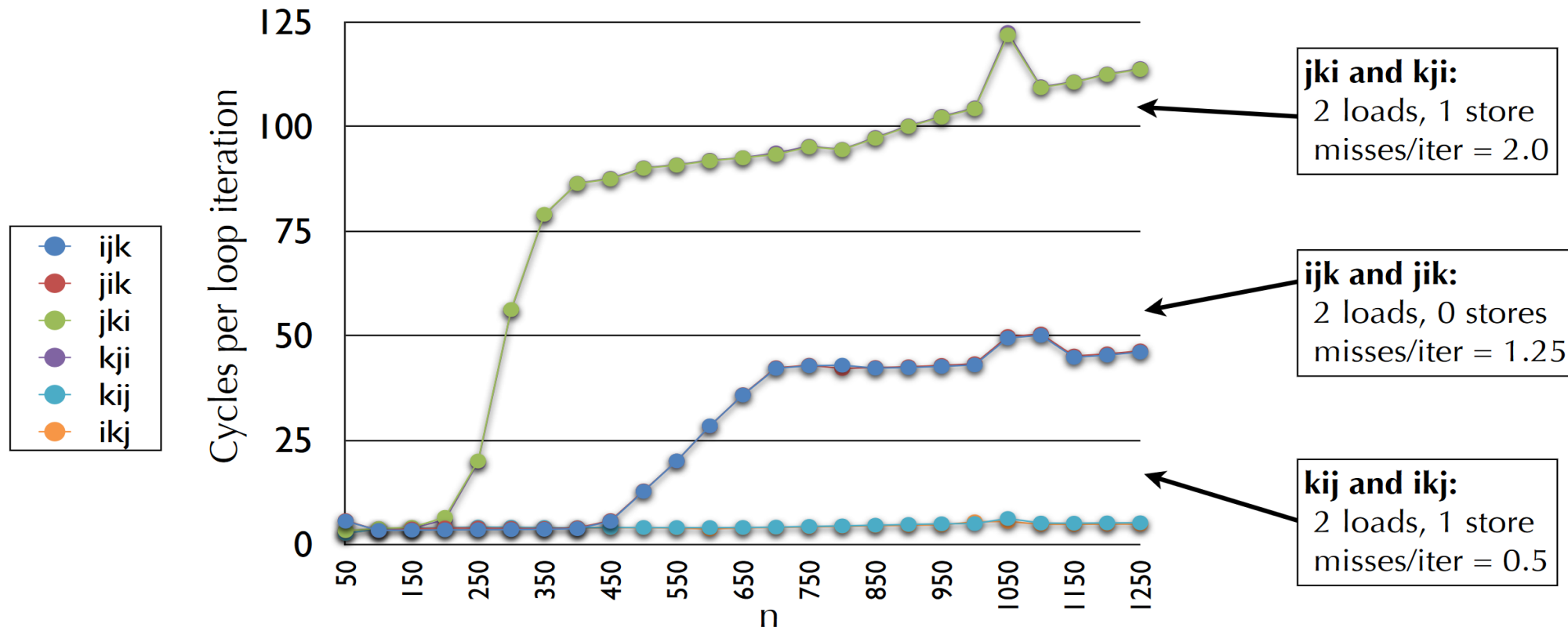
Inner loop:



| | (*,k) | | (k,j) | | (*,j) |
| A | | | B | | C |
| Column-wise | | Fixed | | Column-wise | |

Miss rate for inner loop iterations:

| A | B | C |
|---|---|---|
| 1.0 | 0.0 | 1.0 |

# Effect of loops



**jki and kji:**
2 loads, 1 store
misses/iter = 2.0

**ijk and jik:**
2 loads, 0 stores
misses/iter = 1.25

**kij and ikj:**
2 loads, 1 store
misses/iter = 0.5

- Miss rate better predictor or performance than number of mem. accesses!
- For large N, kij and ikj performance almost constant.
  Due to **hardware prefetching**, able to recognize stride-1 patterns.

# Few Linux commands of interest

perf:
https://perf.wiki.kernel.org/index.php/Tutorial#Counting_with_perf_stat

dmidecode

/proc/cpuinfo

# Vd'aka