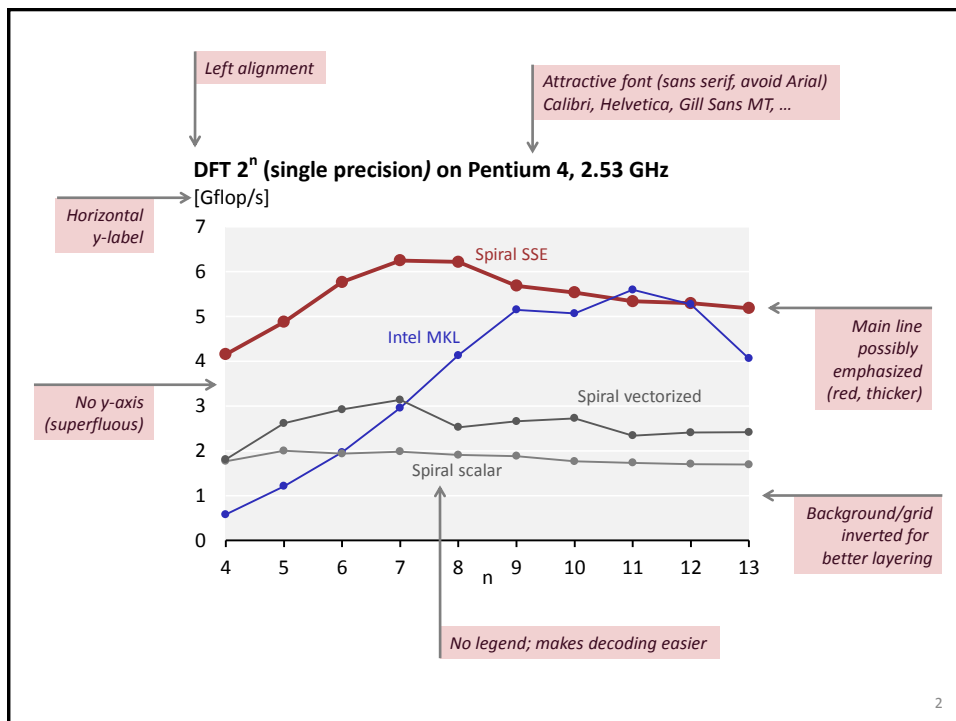# How to Write Fast Numerical Code

Spring 2013
*Lecture:* Memory hierarchy, locality, caches

**Instructor:** Markus Püschel

**TA:** Georg Ofenbeck & Daniele Spampinato

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

---

*Left alignment*

*Attractive font (sans serif, avoid Arial)*
*Calibri, Helvetica, Gill Sans MT, …*

**DFT $2^n$ (single precision) on Pentium 4, 2.53 GHz**

*Horizontal y-label*

[Gflop/s]

Spiral SSE

Intel MKL

Spiral vectorized

*No y-axis (superfluous)*

Spiral scalar

n

*Main line possibly emphasized (red, thicker)*

*Background/grid inverted for better layering*

*No legend; makes decoding easier*

2

*© Markus Püschel*
*Computer Science*
**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2013*

# Organization

- **Temporal and spatial locality**
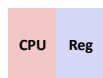- **Memory hierarchy**
- **Caches**

*Chapter 5 in **Computer Systems: A Programmer's Perspective**, 2ⁿᵈ edition,*
*Randal E. Bryant and David R. O'Hallaron, Addison Wesley 2010*
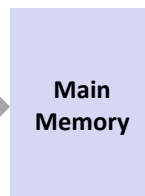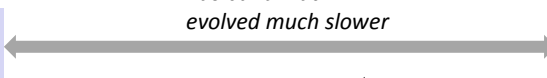*Part of these slides are adapted from the course associated with this book*

---

# Problem: Processor-Memory Bottleneck

*Processor performance*
*doubled about*
*every 18 months*

*Bus bandwidth*
*evolved much slower*

| CPU | Reg |

**Main Memory**

*Core 2 Duo:*
**Peak performance:**
**2 SSE two operand ops/cycles**
consumes up to 64 Bytes/cycle

*Core 2 Duo:*
**Bandwidth**
2 Bytes/cycle

*Solution: Caches/Memory hierarchy*

# Typical Memory Hierarchy

**Smaller, faster, costlier per byte**

**Larger, slower, cheaper per byte**

**L0:** registers — *CPU registers hold words retrieved from L1 cache*

**L1:** on-chip L1 cache (SRAM) — *L1 cache holds cache lines retrieved from L2 cache*

**L2:** on-chip L2 cache (SRAM) — *L2 cache holds cache lines retrieved from main memory*

**L3:** main memory (DRAM) — *Main memory holds disk blocks retrieved from local disks*

**L4:** local secondary storage (local disks) — *Local disks hold files retrieved from disks on remote network servers*

**L5:** remote secondary storage (tapes, distributed file systems, Web servers)

5

---

**Abstracted Microarchitecture: Example Core 2 (2008) and** *Core i7 Sandybridge (2011)*
Throughput (tp) is measured in doubles/cycle. For example: 2 *(4)*
Latency (lat) is measured in cycles
1 double floating point (FP) = 8 bytes
Rectangles not to scale

Core 2 | *Core i7*

**Memory hierarchy:**
- Registers
- L1 cache
- L2 cache
- Main memory
- Hard disk

**double FP:**
*scalar tp:*
- 1 add/cycle
- 1 mult/cycle

*vector (SSE) tp*
- 1 vadd/cycle = 2 adds/cycle
- 1 vmult/cycle = 2 mults/cycle

*ISA*

fadd
fmul
ALU
load
store

**execution units**

out of order execution superscalar

issue 6 μops/cycle

RISC μops

CISC ops

**instruction pool (up to 96 *(168)* "in flight")**

**instruction decoder (up to 5 ops/cycle)**

**internal registers**

**16 FP register**

lat: 3 *(4)*
tp: 2 *(4)*

**L1 Dcache**

**L1 Icache**

*both:*
**32 KB 8-way 64B CB**

lat: 14 *(12)*
tp: 1 *(4)*

**L2 cache 4 MB 16-way 64B CB**

CB = cache block

lat: 100
tp: 1/4

*depends on platform*

**Main Memory (RAM) 4 GB**

lat: millions
tp: ~1/250
*(~1/100)*

*depends on platform*

**Hard disk ≥ 0.5 TB**

**1 Core**

**Core 2 Duo:** *on die*
Core #1
Core #2
L2 — RAM

*Core i7 Sandy Bridge:*
*256 KB L2 cache*
*2–8MB L3 cache: lat 26-31, tp 4*
*RAM: tp 1*
*vector (AVX) tp*
- 1 vadd/cycle = 4 adds/cycle
- 1 vmult/cycle = 4 mults/cycle

*on die*
Core #1 — L2
Core #2 — L2
Core #3 — L2
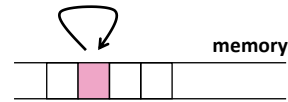Core #4 — L2
L3 — RAM

Source: Intel manual (chapter 2)

6

# Why Caches Work: Locality

- *Locality:* **Programs tend to use data and instructions with addresses near or equal to those they have used recently**
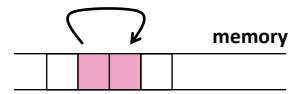  *History of locality*

- *Temporal locality:*
  Recently referenced items are likely
  to be referenced again in the near future

  memory

- *Spatial locality:*
  Items with nearby addresses tend
  to be referenced close together in time

  memory

7

# Example: Locality?

```
sum = 0;
for (i = 0; i < n; i++)
  sum += a[i];
return sum;
```

- **Data:**
  - Temporal: **sum** referenced in each iteration
  - Spatial: array **a[ ]** accessed in stride-1 pattern
- **Instructions:**
  - Temporal: loops cycle through the same instructions
  - Spatial: instructions referenced in sequence

- *Being able to assess the locality of code is a crucial skill for a performance programmer*

8

## Locality Example #1

```c
int sum_array_rows(int a[M][N])
{
  int i, j, sum = 0;

  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      sum += a[i][j];
  return sum;
}
```

## Locality Example #2

```c
int sum_array_cols(int a[M][N])
{
  int i, j, sum = 0;

  for (j = 0; j < N; j++)
    for (i = 0; i < M; i++)
      sum += a[i][j];
  return sum;
}
```

## Locality Example #3

```c
int sum_array_3d(int a[M][N][K])
{
  int i, j, k, sum = 0;

  for (i = 0; i < M; i++)
    for (j = 0; j < N; j++)
      for (k = 0; k < K; k++)
        sum += a[k][i][j];
  return sum;
}
```

**How to improve locality?**

11

## Operational Intensity Again

- **Definition: Given a program P, assume cold (empty) cache**

  *Operational intensity:* $I(n) = \dfrac{W(n)}{Q(n)}$

  — #flops (input size n)

  — #bytes transferred cache $\leftrightarrow$ memory (for input size n)

- **Examples: Determine asymptotic bounds on I(n)**
  - Vector sum: y = x + y                    **O(1)**
  - Matrix-vector product: y = Ax            **O(1)**
  - Fast Fourier transform                   **O(log(n))**
  - Matrix-matrix product: C = AB + C        **O(n)**

12

*© Markus Püschel*
**ETH**
*Computer Science* Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

*How to write fast numerical code*
*Spring 2013*

# Compute/Memory Bound

- **A function/piece of code is:**
  - *Compute bound* if it has high operational intensity
  - *Memory bound* if it has low operational intensity

- **Relationship between operational intensity and locality?**
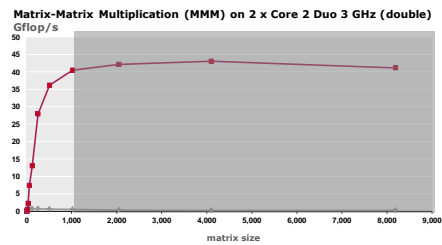  - Operational intensity $\sim$ locality

---

# Effects

**FFT:** $I(n) \leq O(log(n))$

Discrete Fourier Transform (DFT) on 2 x Core 2 Duo 3 GHz (single)
Gflop/s



**MMM:** $I(n) \leq O(n)$

Matrix-Matrix Multiplication (MMM) on 2 x Core 2 Duo 3 GHz (double)
Gflop/s



**Up to 40-50% peak**
**Performance drop outside L2 cache**
*Most time spent transferring data*

**Up to 80-90% peak**
**Performance can be maintained outside L2 cache**
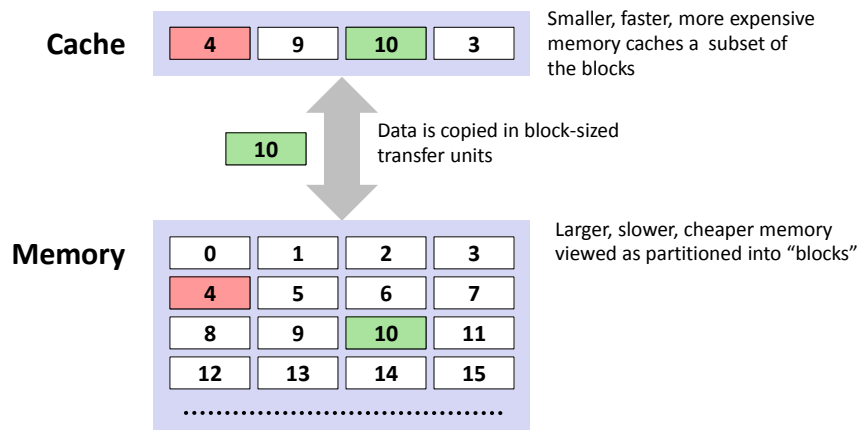*Cache miss time compensated/hidden by computation*

# Cache

- *Definition:* **Computer memory with short access time used for the storage of frequently or recently used instructions or data**

**CPU** ⟷ **Cache** ⟷ **Main Memory**

- **Naturally supports** *temporal locality*
- *Spatial locality* **is supported by transferring data in blocks**
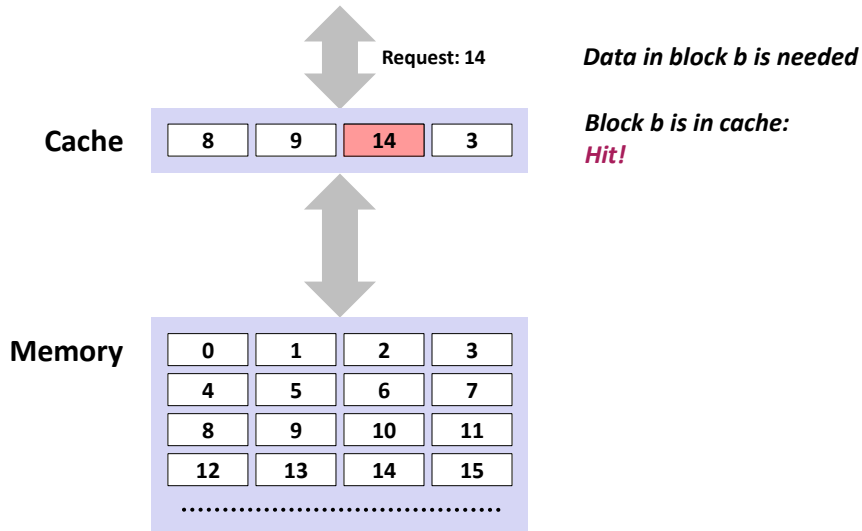  - Core 2: one block = 64 B = 8 doubles
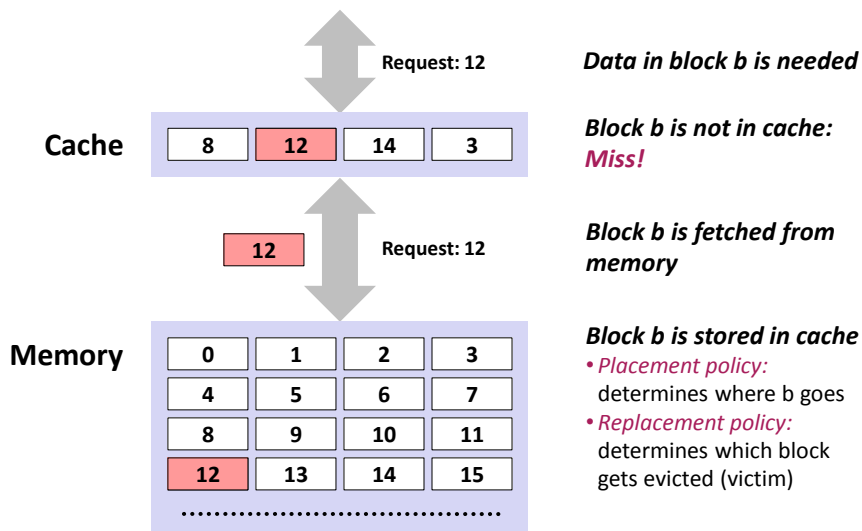
15

# General Cache Mechanics

**Cache**

| 4 | 9 | 10 | 3 |
|---|---|----|---|

Smaller, faster, more expensive memory caches a subset of the blocks

| 10 |
|----|

Data is copied in block-sized transfer units

**Memory**

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Larger, slower, cheaper memory viewed as partitioned into "blocks"

16

# General Cache Concepts: Hit

**Request: 14**

*Data in block b is needed*

**Cache**

| 8 | 9 | 14 | 3 |

*Block b is in cache:*
*Hit!*

**Memory**

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

---

# General Cache Concepts: Miss

**Request: 12**

*Data in block b is needed*

**Cache**

| 8 | 12 | 14 | 3 |

*Block b is not in cache:*
*Miss!*

| 12 |    **Request: 12**

*Block b is fetched from memory*

**Memory**

| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

*Block b is stored in cache*
• *Placement policy:*
  determines where b goes
• *Replacement policy:*
  determines which block
  gets evicted (victim)

# Types of Cache Misses (The 3 C's)

- *Compulsory (cold)* **miss**
  Occurs on first access to a block

- *Capacity* **miss**
  Occurs when working set is larger than the cache

- *Conflict* **miss**
  Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot

- **Not a clean classification but still useful**

# Cache Performance Metrics

- **Miss Rate**
  - Fraction of memory references not found in cache: misses / accesses
    = 1 – hit rate

- **Hit Time**
  - Time to deliver a block in the cache to the processor
  - Core 2:
    3 clock cycles for L1
    14 clock cycles for L2

- **Miss Penalty**
  - Additional time required because of a miss
  - Core 2: about 100 cycles for L2 miss

# Cache Structure

- **Draw a direct mapped cache (E = 1, B = 4 doubles, S = 8)**

- **Show how blocks are mapped into cache**

21

# Example (S=8, E=1)

*Ignore the variables sum, i, j*

**assume: cold (empty) cache, a[0][0] goes here**

```
int sum_array_rows(double a[16][16])
{
  int i, j;
  double sum = 0;

  for (i = 0; i < 16; i++)
    for (j = 0; j < 16; j++)
      sum += a[i][j];
  return sum;
}
```
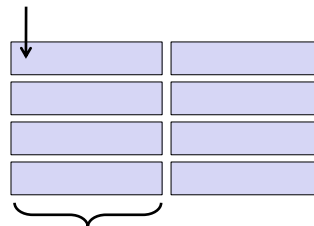
```
int sum_array_cols(double a[16][16])
{
  int i, j;
  double sum = 0;

  for (j = 0; j < 16; i++)
    for (i = 0; i < 16; j++)
      sum += a[i][j];
  return sum;
}
```

**B = 32 byte = 4 doubles**

**blackboard**

22

# Cache Structure

- **Add associativity (E = 2, B = 4 doubles, S = 8)**
- **Show how elements are mapped into cache**

# Example (S=4, E=2)

*Ignore the variables sum, i, j*

```
int sum_array_rows(double a[16][16])
{
  int i, j;
  double sum = 0;

  for (i = 0; i < 16; i++)
    for (j = 0; j < 16; j++)
      sum += a[i][j];
  return sum;
}
```

```
int sum_array_cols(double a[16][16])
{
  int i, j;
  double sum = 0;

  for (j = 0; j < 16; i++)
    for (i = 0; i < 16; j++)
      sum += a[i][j];
  return sum;
}
```

**assume: cold (empty) cache,
a[0][0] goes here**

**B = 32 byte = 4 doubles**
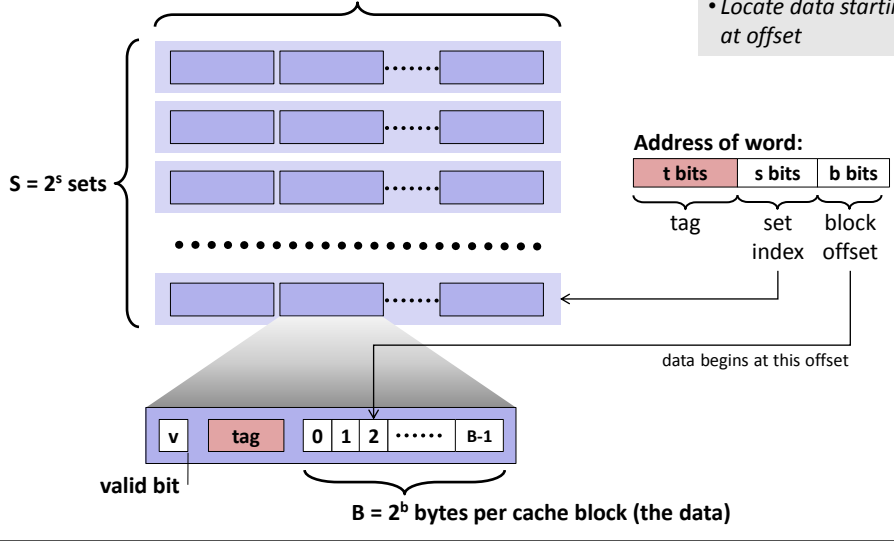
**blackboard**

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set
E = associativity, E=1: direct mapped

set

line

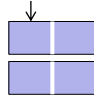$S = 2^s$ sets

**Cache size:**
**S x E x B data bytes**

| v | tag | 0 | 1 | 2 | ····· | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)

25

---

# Cache Read

- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

$E = 2^e$ lines per set
E = associativity, E=1: direct mapped

$S = 2^s$ sets

**Address of word:**

| t bits | s bits | b bits |

tag    set index    block offset

data begins at this offset

| v | tag | 0 | 1 | 2 | ····· | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)

26

# Small Example, Part 1

x[0]
↓

**Cache:**
E = 1 (direct mapped)
S = 2
B = 16 (2 doubles)

**Array (accessed twice in example)**
x = x[0], …, x[7]

```
% Matlab style code
for j = 0:1
  for i = 0:7
    access(x[i])
```

**Access pattern:**  0123456701234567
**Hit/Miss:**  MHMHMHMHMHMHMHMH

**Result:** 8 misses, 8 hits
**Spatial locality:** yes
**Temporal locality:** no

27

# Small Example, Part 2

x[0]
↓

**Cache:**
E = 1 (direct mapped)
S = 2
B = 16 (2 doubles)

**Array (accessed twice in example)**
x = x[0], …, x[7]

```
% Matlab style code
for j = 0:1
  for i = 0:2:7
    access(x[i])
  for i = 1:2:7
    access(x[i])
```

**Access pattern:**  0246135702461357
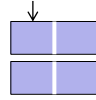**Hit/Miss:**  MMMMMMMMMMMMMMMM

**Result:** 16 misses
**Spatial locality:** no
**Temporal locality:** no

28

## Small Example, Part 3

```
x[0]
 ↓
```

**Cache:**
E = 1 (direct mapped)
S = 2
B = 16 (2 doubles)

**Array (accessed twice in example)**
`x = x[0], …, x[7]`

```
% Matlab style code
for j = 0:1
  for k = 0:1
    for i = 0:3
      access(x[i+4j])
```

**Access pattern:**  0123012345674567
**Hit/Miss:**          MHMHHHHHHMHMHHHHH

**Result:** 4 misses, 12 hits (is optimal, why?)
**Spatial locality:** yes
**Temporal locality:** yes

---
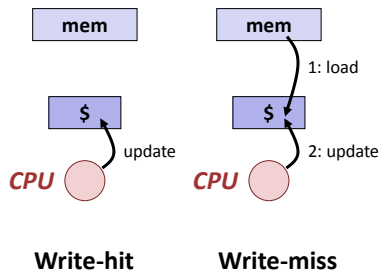
## Terminology

- **Direct mapped cache:**
    - Cache with E = 1
    - Means every block from memory has a unique location in cache

- **Fully associative cache**
    - Cache with S = 1 (i.e., maximal E)
    - Means every block from memory can be mapped to any location in cache
    - In practice to expensive to build

- **LRU (least recently used) replacement**
    - when selecting which block should be replaced (happens only for E > 1), the least recently used one is chosen
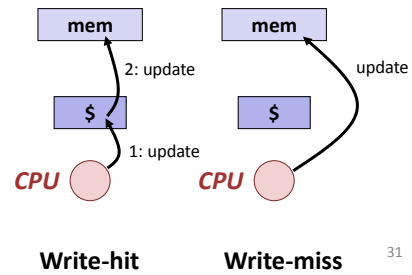
# What about writes?

- **What to do on a write-hit?**
  - *Write-through:* write immediately to memory
  - *Write-back:* defer write to memory until replacement of line
- **What to do on a write-miss?**
  - *Write-allocate:* load into cache, update line in cache
  - *No-write-allocate:* writes immediately to memory

### *Write-back/write-allocate (Core)*

| mem | mem |
|-----|-----|
|  |  1: load ↑ |
| **$** | **$** |
| update ↓ | 2: update ↑ |
| **CPU** ◯ | **CPU** ◯ |
| **Write-hit** | **Write-miss** |

### *Write-through/no-write-allocate*

| mem | mem |
|-----|-----|
| 2: update ↓ | update ↑ |
| **$** | **$** |
| 1: update ↑ |  |
| **CPU** ◯ | **CPU** ◯ |
| **Write-hit** | **Write-miss** |

31

---

# Example: (Blackboard)

- **z = x + y, x, y, z vector of length n**
- **assume they fit jointly in cache + cold cache**
- **memory traffic Q(n)?**
- **operational intensity I(n)?**

32

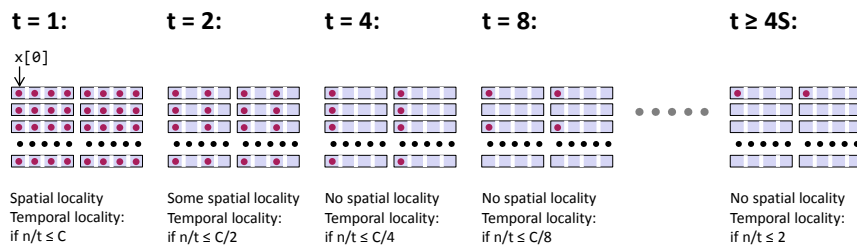# Locality Optimization: Blocking

- **Example: MMM (blackboard)**

33

# The Killer: Two-Power Strided Working Sets

```
% t = 1,2,4,8,… a 2-power
% size of working set: n/t
for (i = 0; i < n; i += t)
  access(x[i])
for (i = 0; i < n; i += t)
  access(x[i])
```

**blackboard**

**Cache: E = 2, B = 4 doubles**

**t = 1:**     **t = 2:**     **t = 4:**     **t = 8:**                    **t ≥ 4S:**

x[0]

• • • • •

Spatial locality          Some spatial locality     No spatial locality       No spatial locality       No spatial locality
Temporal locality:        Temporal locality:        Temporal locality:        Temporal locality:        Temporal locality:
if n/t ≤ C                if n/t ≤ C/2              if n/t ≤ C/4              if n/t ≤ C/8              if n/t ≤ 2

34

# The Killer: Where Can It Occur?

- **Accessing two-power size 2D arrays (e.g., images) columnwise**
    - 2d Transforms
    - Stencil computations
    - Correlations
- **Various transform algorithms**
    - Fast Fourier transform
    - Wavelet transforms
    - Filter banks

# Summary

- **It is important to assess temporal and spatial locality in the code**
- **Cache structure is determined by three parameters**
- **You should be able to roughly simulate a computation on paper**
- **Blocking to improve locality**
- **Two-power strides are problematic (conflict misses)**