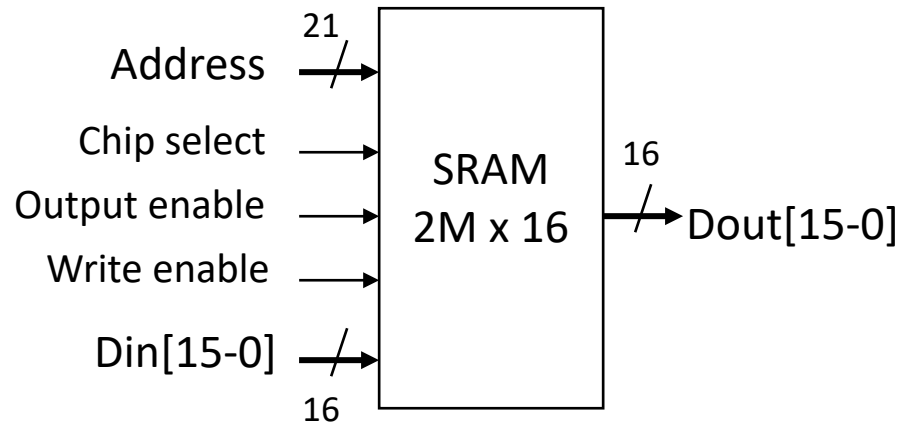# Lecture 4.
# Memory Hierarchy (2)
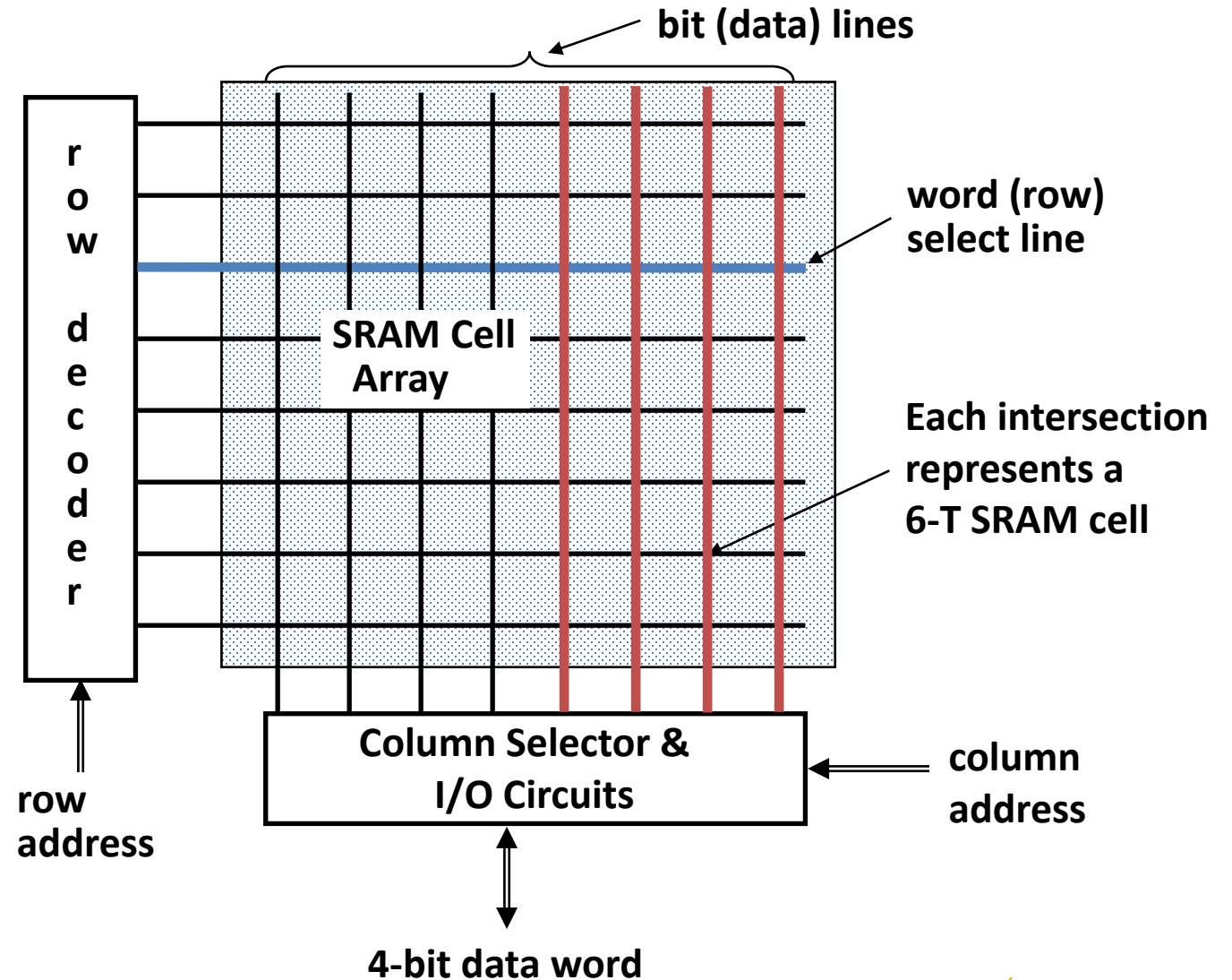
Haonan Wang

SAN JOSÉ STATE
UNIVERSITY

# Cache Design

- **Caches use *SRAM* for speed and technology compatibility**
  - Low density (6 transistor cells), high power, expensive, fast
  - Static: content will last "forever" (until power  turned off)


- **Example: A (2M X 16 bit) SRAM logic**
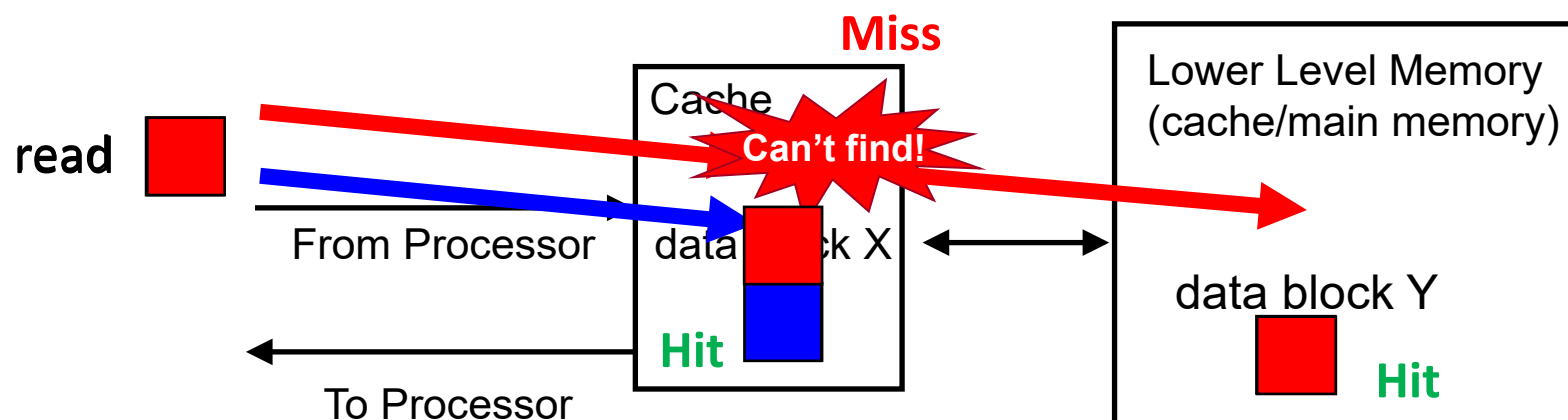
Address $\xrightarrow{\quad 21 \quad}$

Chip select $\longrightarrow$

Output enable $\longrightarrow$ SRAM 2M x 16 $\xrightarrow{\quad 16 \quad}$ Dout[15-0]

Write enable $\longrightarrow$

Din[15-0] $\xrightarrow{\quad 16 \quad}$

# SRAM Cache Design

- **Each row holds a data block**

- **Column address selects the requested word from block**

bit (data) lines

row decoder

SRAM Cell Array

word (row) select line

Each intersection represents a 6-T SRAM cell

row address

Column Selector & I/O Circuits
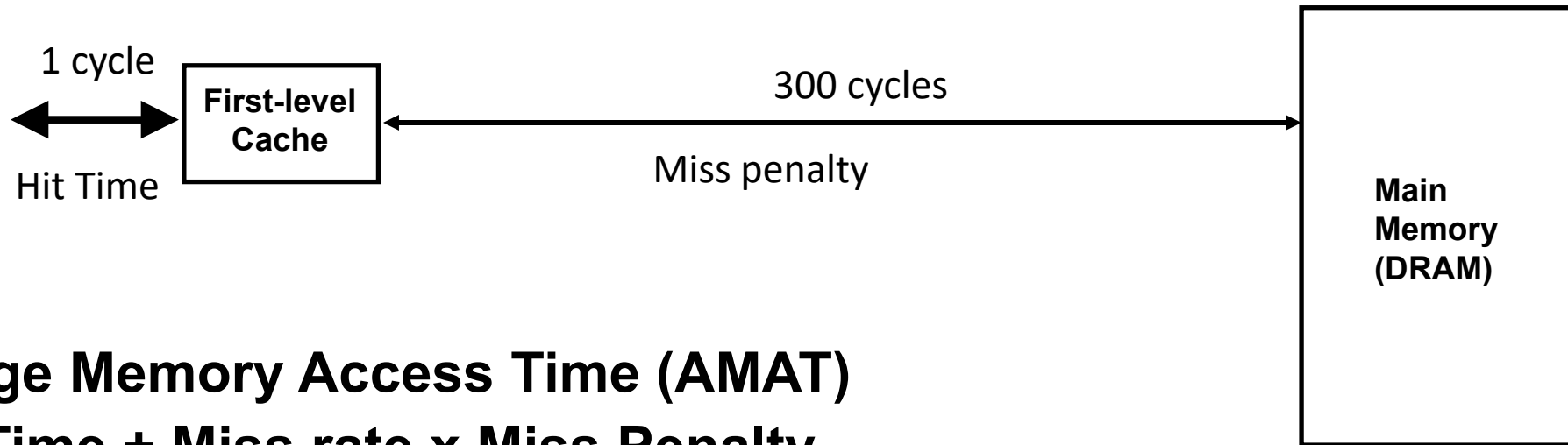
column address

4-bit data word

# Cache Hit and Miss

- **Hit:** Data appears in some block of the cache
  - **Hit Rate:** # hits / total accesses on the cache
  - **Hit Time:** Time to access the cache

- **Miss:** Data needs to be retrieved from the lower level (and stored in cache)
  - **Miss Rate:** 1 - (Hit Rate)
  - **Miss Penalty:** Average delay in the processor caused by each miss

# Memory Hierarchy Performance

```
                  300 cycles
     1 cycle    ┌──────────┐                              ┌──────────┐
  ◄──────────►  │First-level│◄───────────────────────────►│   Main   │
                │  Cache   │                              │  Memory  │
   Hit Time     └──────────┘        Miss penalty          │  (DRAM)  │
                                                          └──────────┘
```

- **Average Memory Access Time (AMAT)**
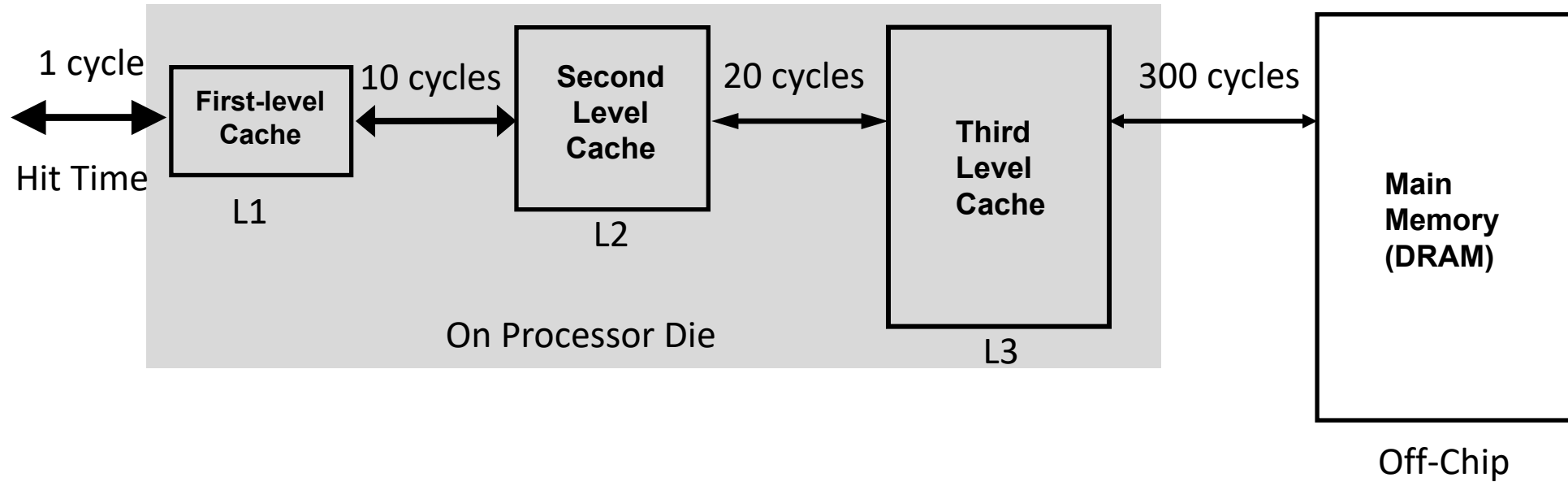  **= Hit Time + Miss rate x Miss Penalty**

- **Example:**
  - Cache Hit = 1 cycle
  - Miss rate = 10% = 0.1
  - Miss penalty = 300 cycles
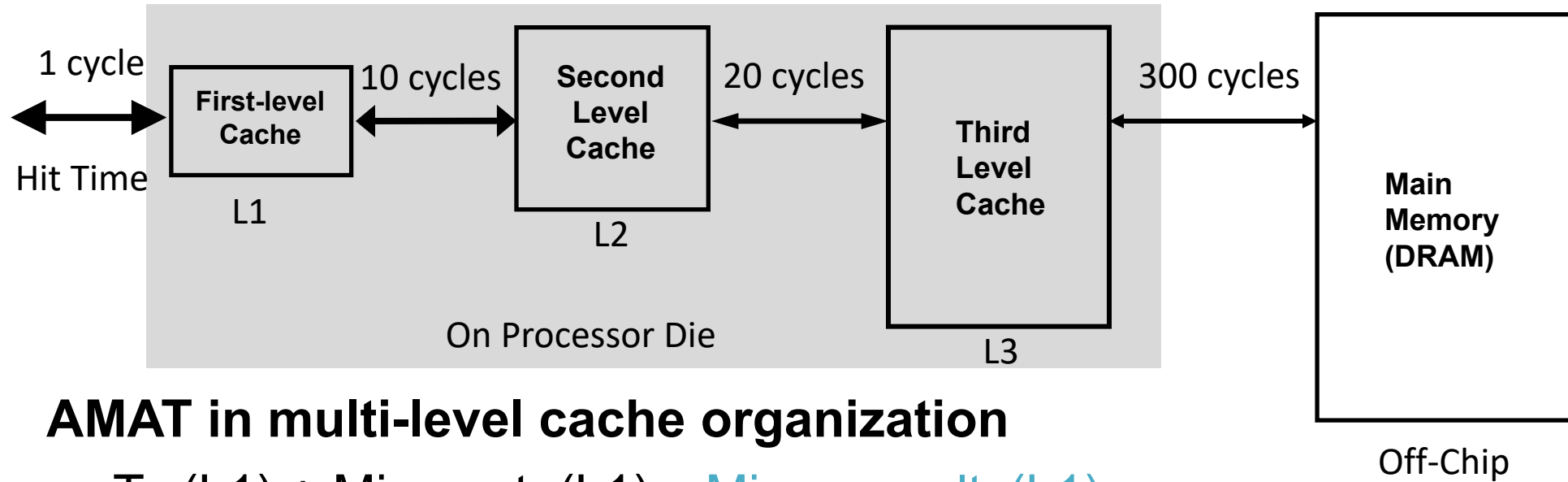  - AMAT = $T_{hit}$(L1) + Miss_rate(L1) x T(Memory) = 1 + 0.1 x 300 = 31 cycles

Due to long-latency memory, AMAT is 30 cycles longer than cache latency.
Can we reduce the overhead?

SJSU   SAN JOSÉ STATE UNIVERSITY

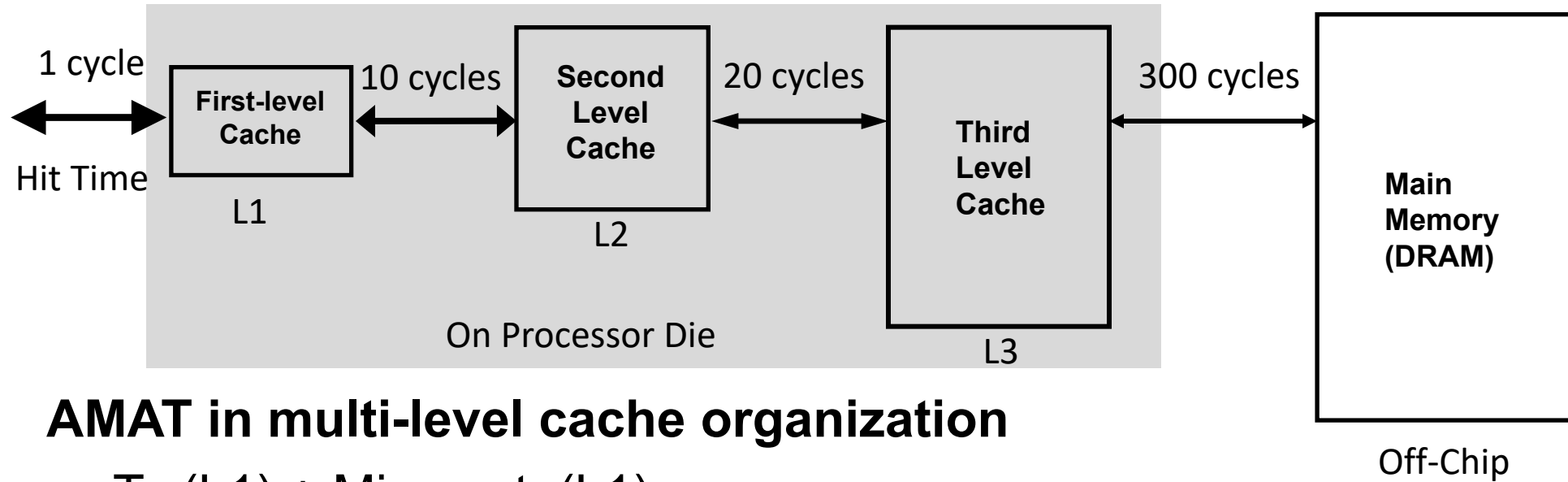# Reducing Penalty: Multi-Level Cache

# Reducing Penalty: Multi-Level Cache



1 cycle    First-level Cache    10 cycles    Second Level Cache    20 cycles    Third Level Cache    300 cycles    Main Memory (DRAM)

Hit Time

L1    L2    L3

On Processor Die

Off-Chip

- **AMAT in multi-level cache organization**

    $= T_{hit}(L1) + Miss\_rate(L1) \times Miss\_penalty(L1)$

SJSU    SAN JOSÉ STATE UNIVERSITY

# Reducing Penalty: Multi-Level Cache



- **AMAT in multi-level cache organization**
  = $T_{hit}(L1)$ + Miss_rate(L1) x

  [ $T_{hit}(L2)$ + Miss_rate(L2) x Miss_penalty(L2) ]

SJSU    SAN JOSÉ STATE
        UNIVERSITY

# Reducing Penalty: Multi-Level Cache



1 cycle

Hit Time

**First-level Cache**

L1

10 cycles

**Second Level Cache**

L2

20 cycles

**Third Level Cache**

L3

300 cycles

**Main Memory (DRAM)**

Off-Chip

On Processor Die
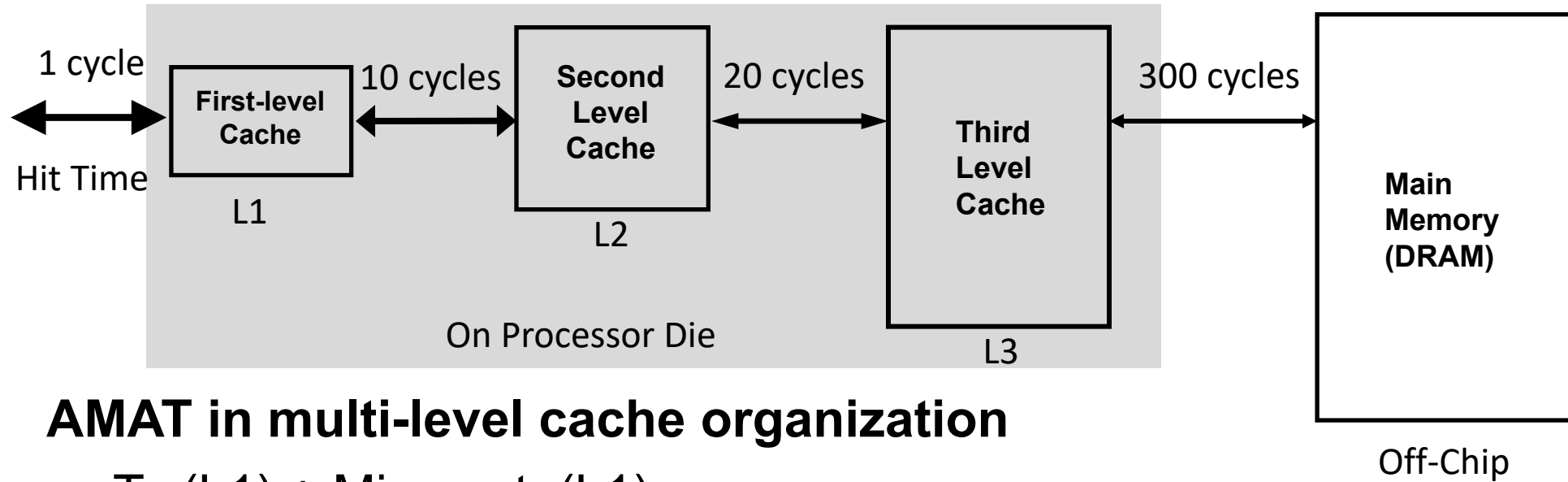
- **AMAT in multi-level cache organization**

$= T_{hit}(L1) + Miss\_rate(L1) \times$

$[\ T_{hit}(L2) + Miss\_rate(L2) \times$

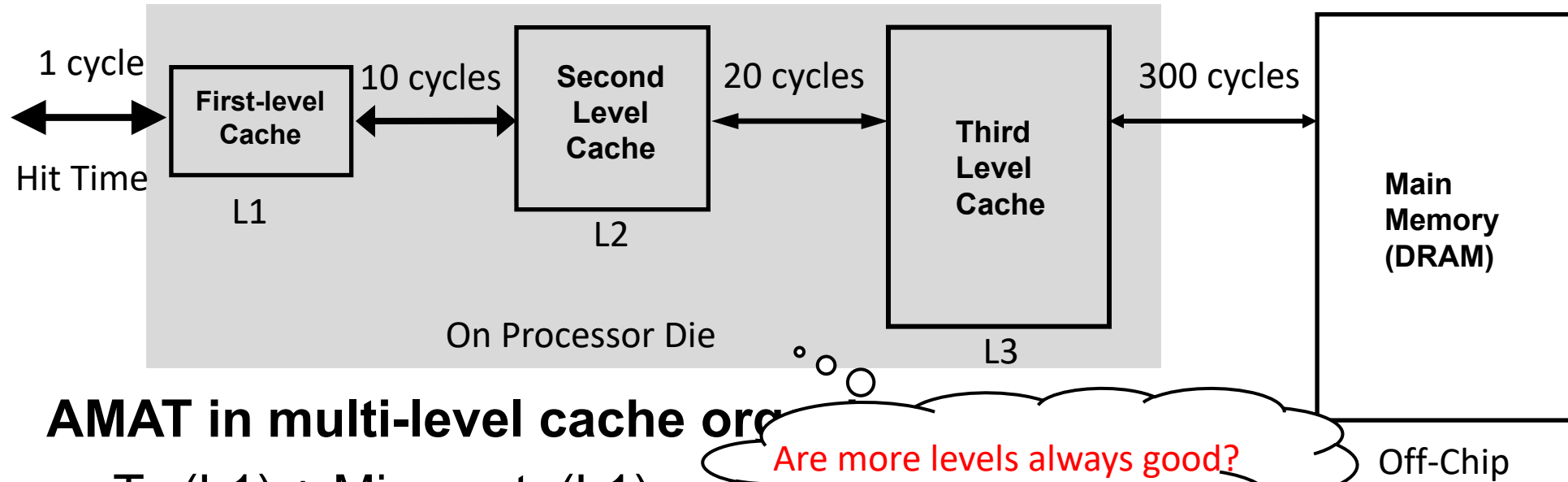$\{\ T_{hit}(L3) + Miss\_rate(L3) \times Miss\_penalty(L3)\ \}\ ]$

SJSU   SAN JOSÉ STATE UNIVERSITY

# Reducing Penalty: Multi-Level Cache



- **AMAT in multi-level cache organization**

$$= T_{hit}(L1) + Miss\_rate(L1) \times$$

$$[\ T_{hit}(L2) + Miss\_rate(L2) \times$$

$$\{\ T_{hit}(L3) + Miss\_rate(L3) \times T(memory)\ \}\ ]$$

# Reducing Penalty: Multi-Level Cache



```
                    10 cycles              20 cycles              300 cycles
  1 cycle    First-level          Second                 Third                Main
             Cache                Level                  Level                Memory
             L1                   Cache                  Cache                (DRAM)
  Hit Time                        L2                     L3
                        On Processor Die                                     Off-Chip
```

Are more levels always good?

- **AMAT in multi-level cache org**
  = $T_{hit}(L1)$ + Miss_rate(L1) x
     [ $T_{hit}(L2)$ + Miss_rate(L2) x
       { $T_{hit}(L3)$ + Miss_rate(L3) x T(memory) } ]

- **Example:**
  – Miss rate of L1, L2, L3 = 10%, 5%, 1%, respectively
  – AMAT = 1 + 0.1 x [ 10 + 0.05 x { 20 + 0.01 x 300 } ] = 2.115 cycles

Vs. 31 cycles
14.7x speedup!

SJSU  SAN JOSÉ STATE UNIVERSITY

# Discussion

**Background: The cache space is limited. We can only keep a subset of data in cache.**

**Question: What happens when the cache is full?**

**Question: What to keep in the cache (which block should be evicted)?**

SJSU    SAN JOSÉ STATE
UNIVERSITY

# What to Keep in Caches (1)?

- **It depends on the cache organization and replacing policy.**
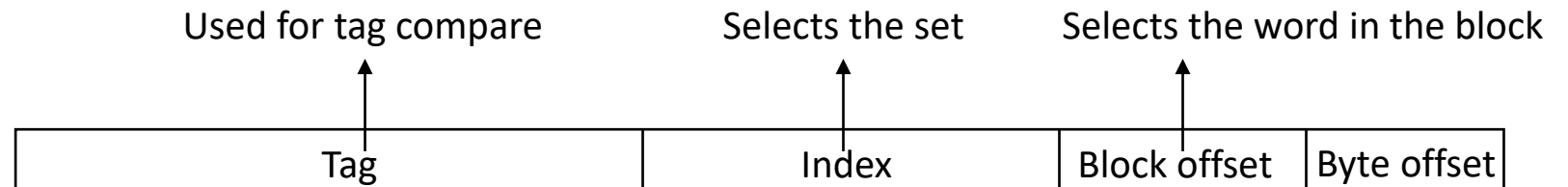
|  | Way 0 | Way 1 ... |
|---|---|---|
| Set 0 | block 0 | block 2 |
| Set 1 | block 1 | block 3 |
⋮

- **Cache organization:**

  - **Cache line (block):** The basic unit of data replacement. A longer cache line fetches and replaces more data per miss.

  - **Set:** An entry that one cache line is mapped to according to certain bits of its address.

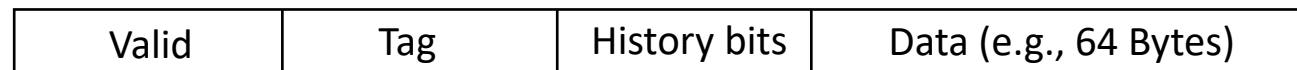  - **Way:** A slot within a cache set to hold one history cache line.

SJSU SAN JOSÉ STATE UNIVERSITY

# What to Keep in Caches (2)?

- **It depends on the cache organization and replacing policy**

- **Cache line replacing policy:** which history line should be replaced?
  - In a set entry, each cache line can be identified with a **Tag** (a portion of its address).

  - Least recently used (LRU), First-in first-out (FIFO), Random, etc.

Used for tag compare      Selects the set      Selects the word in the block

**Address decoding:**

| Tag | Index | Block offset | Byte offset |
|-----|-------|--------------|-------------|

**Cache line:**

| Valid | Tag | History bits | Data (e.g., 64 Bytes) |
|-------|-----|--------------|-----------------------|

# Cache Indexing

**Example: a simple 2-set x 2-way cache**

|  | **Way 0** | **Way 1** |
|---|---|---|
| **Set 0** | block 0 | block 2 |
| **Set 1** | block 1 | block 3 |

**Address decoding:**

Used for tag compare                    Selects the set         Selects the word in the block

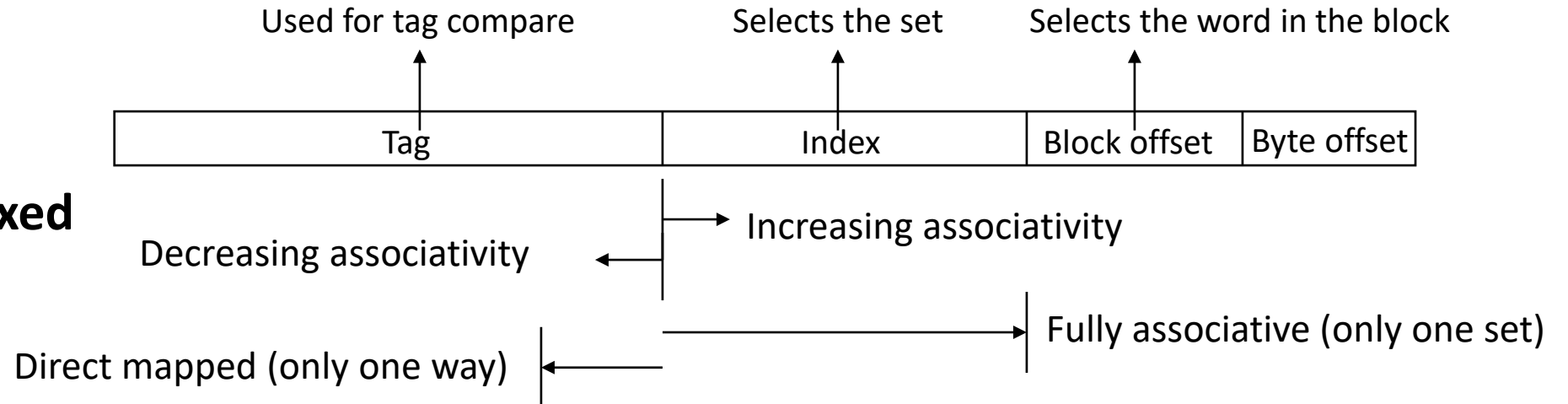| Tag | Index | Block offset | Byte offset |
|---|---|---|---|

# Cache Types

- **N-way Set-Associative:** Number of ways > 1 & Number of sets > 1
  - Slightly complex searching mechanism

- **Direct Mapped:** Number of ways = 1
  - Fast indexing mechanism

- **Fully-Associative:** Number of sets = 1
  - Extensive hardware resources required to search

|  | Way 0 | Way 1 ⋯ |
|---|---|---|
| **Set 0** | block 0 | block 2 |
| **Set 1** | block 1 | block 3 |

⋮

# Cache Types

|  | Way 0 | Way 1 ... |
|---|---|---|
| **Set 0** | block 0 | block 2 |
| **Set 1** | block 1 | block 3 |

⋮

Used for tag compare  Selects the set  Selects the word in the block

| Tag | Index | Block offset | Byte offset |
|---|---|---|---|

**Assuming fixed sized cache:**

Decreasing associativity  Increasing associativity

Direct mapped (only one way)  Fully associative (only one set)

17

SJSU SAN JOSÉ STATE UNIVERSITY

SAN JOSÉ STATE UNIVERSITY *powering* SILICON VALLEY