# Cache

- SRAM
  - 6 transistors per memory cell.
  - Low density
  - Fast access latency of 0.5 - 5 ns
  - More costly.
  - used in flip-flops, cache.

- DRAM
  - 1 transistor per memory cell
  - High density
  - Slow access latency of 50 - 70 ns
  - less costly
  - used in main memory.

|          | Capacity   | Latency  | Cost/GB |
|----------|------------|----------|---------|
| Register | 100s Bytes | 20 ps    | $$$$    |
| SRAM     | 100s KB    | 0.5-5 ns | $$$     |
| DRAM     | 100s MB    | 50-70 ns | $       |
| Hard Disk| 100s GB    | 5-20 ms  | Cents   |
| **Ideal**| **1 GB**   | **1 ns** | **Cheap** |

L Big and Fast

## Cache

- Idea: keep frequently and recently used data in smaller but faster memory.

- Principle of Locality: Program accesses only a small portion of the memory address space within a small time interval.
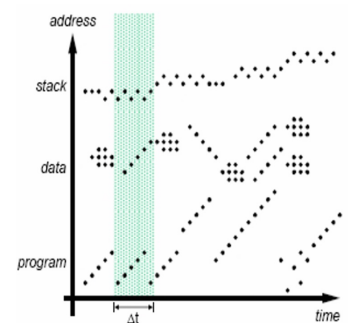
- Temporal Locality (Time)
  - If an item is referenced, it will tend to be referenced again soon.
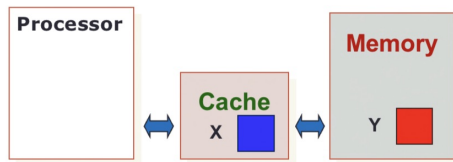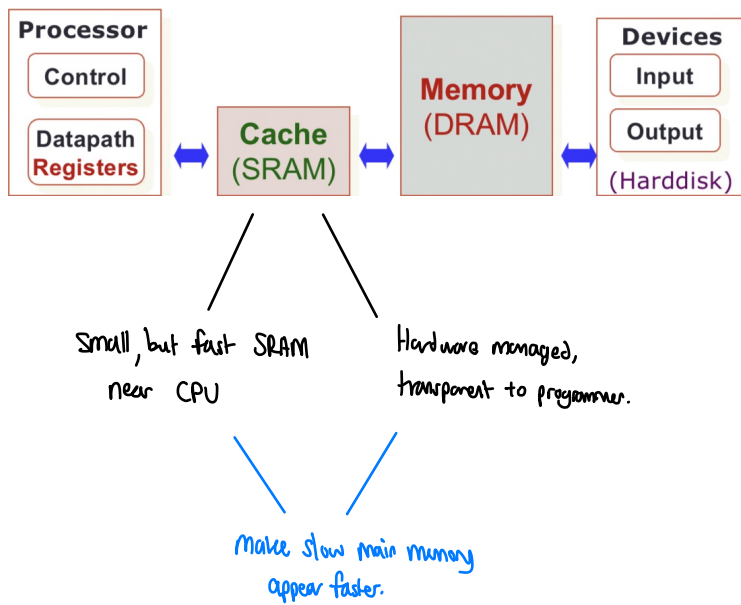
- Spatial Locality (Space)
  - If an item is referenced, nearby items will tend to be referenced soon.

- Different locality for instructions and data.



Capture the working set and keep it in the memory closest to CPU.

# Cache Memory



Small, but fast SRAM near CPU

Hardware managed, transparent to programmer.

Make slow main memory appear faster.



- **Hit**: Data is in cache (e.g., **X**)
  - Hit rate: Fraction of memory accesses that hit
  - Hit time: Time to access cache (very fast)
- **Miss**: Data is not in cache (e.g., **Y**)
  - Miss rate = 1 – Hit rate (opposite of Hit rate)
  - Miss penalty: Time to replace cache block + hit time
- Hit time < Miss penalty

| **Average Access Time** |
| --- |
| = Hit rate x Hit Time + (1-Hit rate) x Miss penalty |

Example:
- Suppose our on-chip SRAM (cache) has **0.8 ns** access time, but the fastest DRAM (main memory) we can get has an access time of **10ns**. **How high a hit rate** do we need to sustain an average access time of **1ns**?

  Let $h$ be the desired hit rate.
  $1 = 0.8h + (1 - h) \times (10 + 0.8)$
  $= 0.8h + 10.8 - 10.8h$
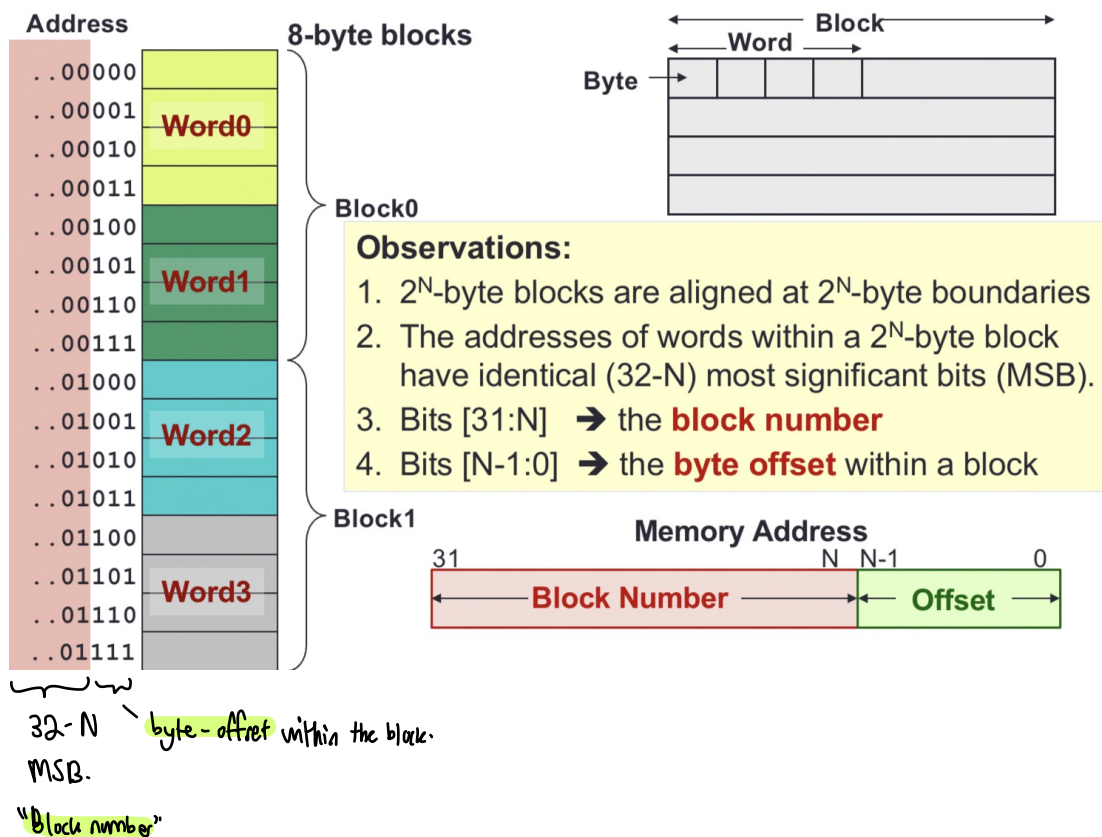  $10h = 9.8 \rightarrow h = 0.98$
  Hence we need a hit rate of **98%**.

# Memory to cache mapping

- Cache Block / Line — unit of transfer between memory and cache.
  - typically one or more words
    - 16-byte block $\approx$ 4-word block.
    - 32-byte block $\approx$ 8-word block.

**Address**

| | 8-byte blocks |
|---|---|
| ..00000 | |
| ..00001 | **Word0** |
| ..00010 | |
| ..00011 | |
| ..00100 | |
| ..00101 | **Word1** |
| ..00110 | |
| ..00111 | |
| ..01000 | |
| ..01001 | **Word2** |
| ..01010 | |
| ..01011 | |
| ..01100 | |
| ..01101 | **Word3** |
| ..01110 | |
| ..01111 | |

Block0

Block1

Block

Word

Byte →

## Observations:
1. $2^N$-byte blocks are aligned at $2^N$-byte boundaries
2. The addresses of words within a $2^N$-byte block have identical (32-N) most significant bits (MSB).
3. Bits [31:N] → the **block number**
4. Bits [N-1:0] → the **byte offset** within a block

**Memory Address**

| 31 | N  N-1 | 0 |
|---|---|---|
| **Block Number** → | | ← **Offset** → |

32-N — byte-offset within the block.
MSB.
"Block number"

- **Direct Mapped Cache.**

**Block Number** (Not Address!)

| | |
|---|---|
| ..00000 | One block |
| Index ..00001 | |
| ..00010 | |
| ..00011 | |
| ..00100 | |
| ..00101 | |
| ..00110 | |
| ..00111 | |
| ..01000 | |
| ..01001 | |
| ..01010 | |
| ..01011 | |
| ..01100 | |
| ..01101 | |
| ..01110 | |
| .01111 | |

**Memory**

**Cache Index**

| | |
|---|---|
| One block | 00 |
| | 01 |
| | 10 |
| | 11 |

**Cache**

— at any one point in time, can only have 1 red block in the location in cache.

**Mapping Function:**
**Cache Index**
**= (BlockNumber) modulo**
**(NumberOfCacheBlocks)**

**Observation:**
If Number of Cache Blocks = $2^M$
➔ the last M bits of the block number is the **cache index**
**Example:**
Cache has $2^2 = 4$ blocks
➔ last 2 bits of the block number is the cache index.

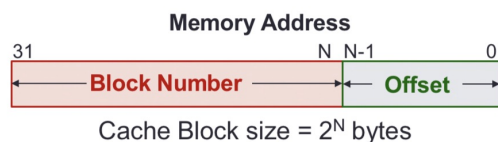2 LSBs → 4 slots in cache.

Multiple memory blocks can map to the same cache block
→ same cache index.
However, they have unique tag number.
Tag = Block number / Number of Cache Block.

**Memory Address**

| 31 | N N-1 | 0 |
|---|---|---|
| ← Block Number → | ← Offset → | |

Cache Block size = $2^N$ bytes

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**Memory Address**

| 31 | N+M-1 | N N-1 | 0 |
|---|---|---|---|
| ← Tag → | ← Index → | ← Offset → | |

Cache Block size = $2^N$ bytes
Number of cache blocks = $2^M$
**Offset = N bits**
**Index = M bits**
**Tag = 32 – (N + M) bits**

# Cache Structure

– 0 or 1

| Valid | Tag | | Data | Index |
|-------|-----|---|------|-------|
| | | | | 00 |
| | | | | 01 |
| | | | | 10 |
| | | | | 11 |

Cache

Along with a data block (line), cache also contains the following administrative information (overheads):
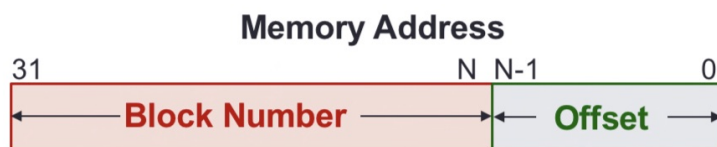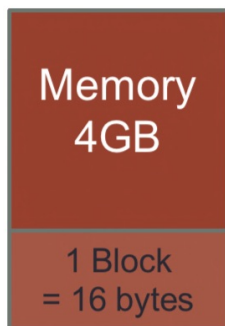1. **Tag** of the memory block
2. **Valid bit** indicating whether the cache line contains valid data

**When is there a cache hit?**
( Valid[index] = TRUE ) **AND**     is the slot filled ?
( Tag[ index ] = Tag[ memory address ] )

## Eg.

**Memory**
**4GB**

1 Block
= 16 bytes

**Cache**
**16KB**

1 Block
= 16 bytes

**Memory Address**

31 ............................... N  N-1 ........... 0

| Block Number | Offset |

Offset, **N** = 4 bits
**Block Number** = 32 − 4 = **28 bits**
Check: Number of Blocks = $2^{28}$

31   18        N+M-1  10    N  N-1    4      0

| Tag | Index | Offset |

(i)

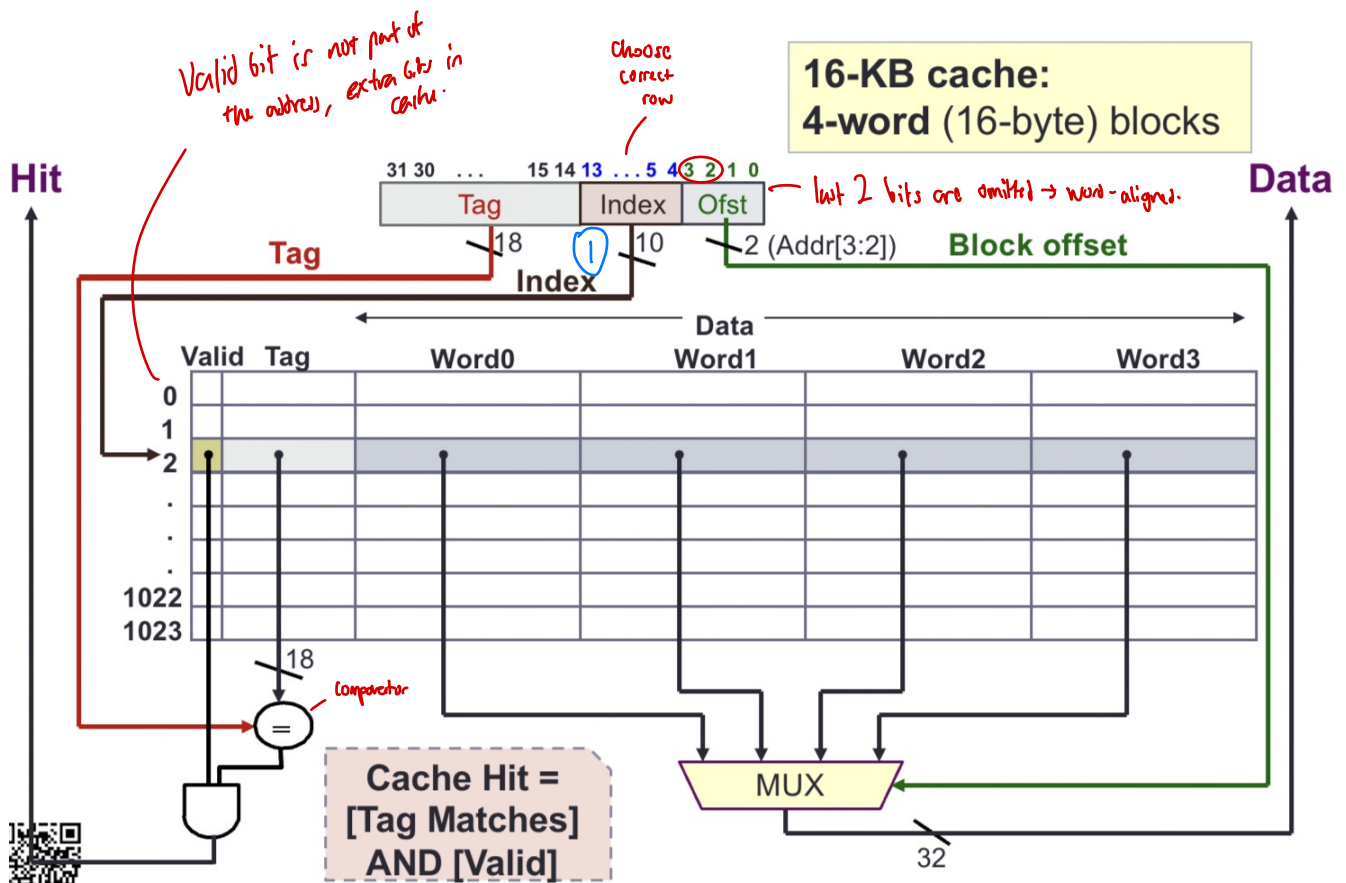**Number of Cache Blocks**
= 16KB / 16bytes = 1024 = $2^{10}$
**Cache Index, M** = 10bits
**Cache Tag** = 32 − 10 − 4 = **18 bits**

# Cache Diagram

**Hit** | **Data**

**16-KB cache:**
**4-word** (16-byte) blocks

Valid bit is not part of the address, extra bits in cache.

Choose correct row

31 30 ... 15 14 13 ... 5 4 3 2 1 0

| Tag | Index | Ofst |

last 2 bits are omitted → word-aligned.

**Tag** — 18
**Index** — 10 — 1

**Block offset** — 2 (Addr[3:2])

**Data**

| Valid | Tag | Word0 | Word1 | Word2 | Word3 |

0
1
2
.
.
.
1022
1023

18

Comparator =

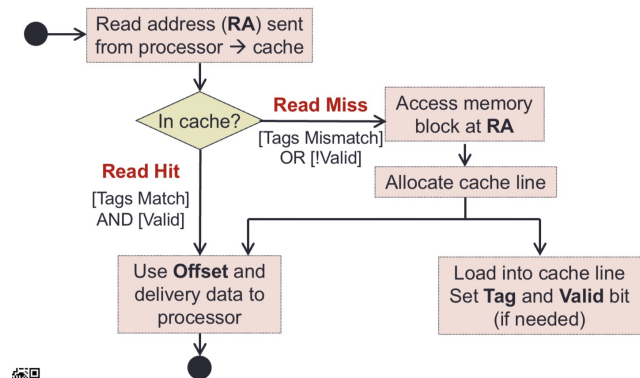**Cache Hit = [Tag Matches] AND [Valid]**

MUX

32

---

**\* Initially cache is empty**

Steps: (Load)

1. Check Cache Block at specified index

2. Data in specified block is **invalid** (Cold / Compulsory Miss)

3. Load 16 bytes from memory; Set **Tag** and **Valid** bit.
   ↳ To populate the "words" slot.     ↳ To 1.

4. Return (offset **2 MSB**) Word to Register

**(2nd time)**
1. Check Cache Block at specified index
2. [Cache Block is Valid] AND [Tags match] → Cache hit!
3. Look at offset to get the specified Word. [Spatial Locality].

**(overriding)**
1. Check Cache Block at specified index
2. Cache block is Valid but **Tags mismatched** [Cold miss]
3. **Replace** block 1 with new data; Set Tag.
4. Check offset, return specified Word to Register.

---

Read address (**RA**) sent from processor → cache

In cache?

**Read Miss** [Tags Mismatch] OR [!Valid] → Access memory block at **RA** → Allocate cache line

**Read Hit** [Tags Match] AND [Valid]

Use **Offset** and delivery data to processor

Load into cache line Set **Tag** and **Valid** bit (if needed)

---

hit if block in mem == block in cache

# Cache Misses

1. **Compulsory / Cold Misses**
   - On first access to a block; the block must be brought into the cache.
   - Aka Cold Start misses or first reference misses.

2. **Conflict Misses**
   - Occur in direct mapped cache / set associative cache, when several blocks are mapped to the same block/set. *(Some other blocks)*
   - Aka Collision misses / interference misses.

3. **Capacity Misses**
   - Occur when blocks are discarded from cache as cache cannot contain all blocks needed.

# Handling Cache Misses

- On **Read Miss**:
  - Data loaded into cache and then load from there to register.

- **Write Miss** option 1: **Write allocate**
  - Load the complete block into cache
  - Change only the required word in the cache
  - Write to main memory depends on write policy. (WB/WT)

- **Write Miss** option 2: **Write around**
  - Do not load the block to cache
  - Write directly to main memory only.

---

# Steps: (Write)

1. Check index, offset
2. Check if Cache is Valid and Tags match → Cache hit!
3. Replace chosen offset with new Word. ( Changed data in Cache only )
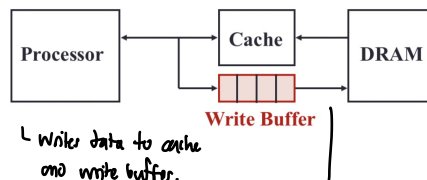
⇓

**Write Policy**

- Cache and main memory are inconsistent.
- Solution 1: **Write-through** Cache
  - Write data both to cache and main memory.
    - └ Slow — will only operate at the speed of main memory.
      - └ Buffer btw Cache and main memory ( program could carry on, no need wait)



└ Writes data to cache and write buffer.

- Memory Controller: write contents of buffer to memory.

- Solution 2: **Write-back** cache
  - only write to cache
  - Write to main memory only when cache block is replaced (evicted).
  - More complicated to implement, wasteful if write back every evicted cache block.
    - └ Add an additional bit ( Dirty bit) to each cache block —— indicate if data is changed by write operation.
      - Write operation will change dirty bit to 1
        - Only cache block is updated, no write to memory.
      - When a cache block is replaced:
        - Only write back to memory if dirty bit is 1.

---

Flowchart:

- Write address (**WA**) and **value** sent from processor → cache
- In cache?
  - **Write Hit** [Tags Match] AND [Valid]
  - **Write Miss** [Tags Mismatch] OR [!Valid]
- **Depends on Write Miss Policy** (Write Allocate) or (Write Around)
  - If Write Allocate
- **Depends on Write Policy** (Write Back) or (Write Through)

Diagram: Processor ↔ Cache ↔ DRAM, with Write Buffer between Cache and DRAM.