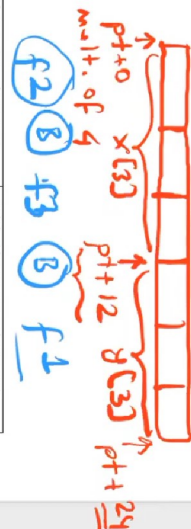
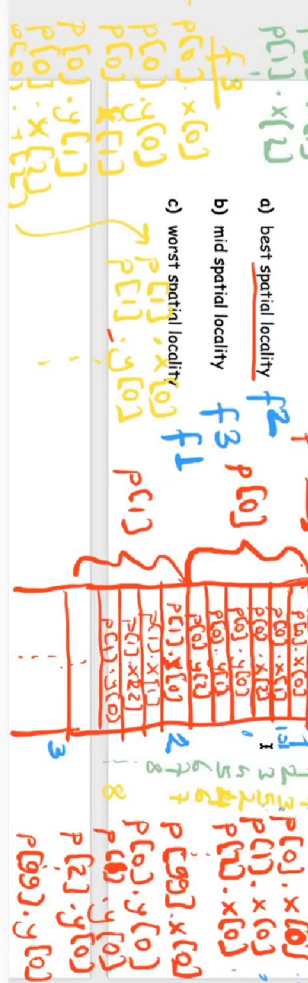


5. (12 pts) The three functions below perform the same operation with varying degrees of spatial locality. Rank-order the functions with respect to the spatial locality employed by each. Also, write an explanation for your ranking.

```
#define N 100
typedef struct {
    int x[3];
    int y[3];
} vec;
vec p[10];
```

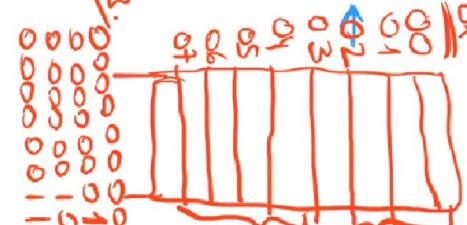
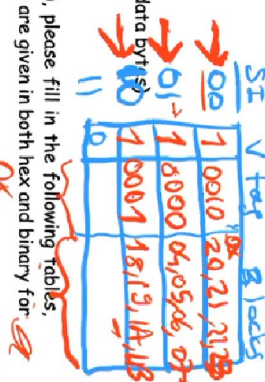
```
void f1(vec *p, int n) {
    int i;
    for (i=0; i<n; i++) {
        for (j=0; j<3; j++) {
            p[i].x[j] = 0;
        }
        for (j=0; j<3; j++) {
            p[i].y[j] = 0;
        }
    }
}
```

- a) best spatial locality  
b) mid spatial locality  
c) worst spatial locality



- The block size is 4 bytes ( $B = 4$ )
  - The cache has 4 sets ( $S = 4$ )
  - The cache is direct mapped ( $E = 1$ )
- a) What is the total capacity of the cache? (in number of data bytes)
- b) How long is a tag? (in number of bits)
- c) Assuming that the cache starts clean (all lines invalid), please fill in the following tables, describing what happens with each operation. Addresses are given in both hex and binary for your convenience.

Operation	Set index?	Hit or Miss?	Eviction?
load 0x00 (0000 0000) <sub>2</sub>	0	Miss	no
load 0x04 (0000 0100) <sub>2</sub>	1	Miss	no
load 0x08 (0000 1000) <sub>2</sub>	2	Miss	no
store 0x12 (0001 0010) <sub>2</sub>	0	Miss	yes
load 0x16 (0001 0100) <sub>2</sub>	1	Miss	yes
store 0x06 (0000 0110) <sub>2</sub>	1	Miss	yes
load 0x18 (0001 1000) <sub>2</sub>	2	Miss	yes
load 0x20 (0010 0000) <sub>2</sub>	0	Miss	yes
store 0x14 (0001 1010) <sub>2</sub>	2	Hit	no



## Array Example

### Practice Problem 3.38 (solution page 377)

Consider the following source code, where  $M$  and  $N$  are constants declared with #define:

```
long sum_element(long i, long j) {
    i in %rdi, j in %rsi
    1 sum_element:
    2 leaq 0(%rdi,8), %rdx
    3 subq %rdi, %rdx
    4 addq %rsi, %rdx
    5 leaq (%rsi,%rdi,4), %rax
    6 addq %rax, %rdi
    7 movq Q(%rdi,8), %rax
    8 addq P(%rdx,8), %rax
    9 ret
}
```

long P[M][N];  
long Q[N][M];  
long sum\_element(long i, long j) {  
return P[i][j] + Q[j][i];  
}

- 8 is scale value
- Offset of matrix P  $7i + j \rightarrow$
- P has 7 columns
- Offset of matrix Q  $5j + i \rightarrow$
- Q has 5 columns
- So,  $M = 5$  and  $N = 7$

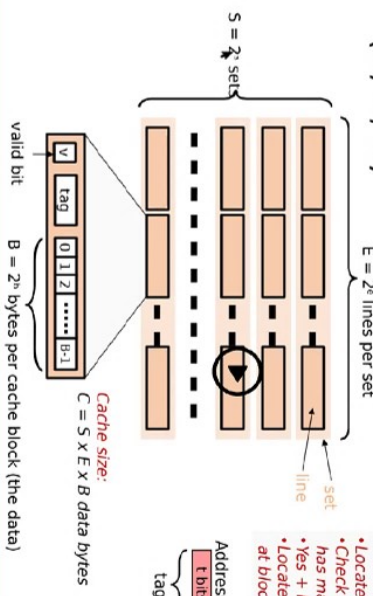
Compute  $8i \rightarrow rdx$   
Compute  $7i \rightarrow rdx$   
Compute  $7i + j \rightarrow rdx$   
Compute  $5j \rightarrow rax$   
Compute  $i + 5j \rightarrow rdi$   
Retrieve  $M[8i+8*5j + i] \rightarrow rax$   
Add  $M[8j + 8*(7i+j)] \rightarrow rax$

Handwritten notes and diagrams for the first problem. The notes include calculations for spatial locality, such as  $5+2+2=9$  and  $12 \text{ bytes} = 12 \text{ bytes}$ . There are also diagrams showing memory access patterns and a table of memory addresses.

Address	Content	Header/Footer?	is pointed?
0x0000	00000000000000000000000000000000	Header	
0x0002	Free		
0x0004	Free		
0x0006	00000000000000000000000000000000	Header	
0x0008	00000000000000000000000000000000	Header	
0x000a	Allocated		
0x000c	Allocated		
0x000e	Allocated		
0x0010	00000000000000000000000000000000	Header	
0x0012	00000000000000000000000000000000	Header	
0x0014	00000000000000000000000000000000	Header	
0x0016	00000000000000000000000000000000	Header	
0x0018	00000000000000000000000000000000	Header	
0x001a	00000000000000000000000000000000	Header	
0x001c	00000000000000000000000000000000	Header	
0x001e	00000000000000000000000000000000	Header	
0x0020	00000000000000000000000000000000	Header	
0x0022	00000000000000000000000000000000	Header	
0x0024	00000000000000000000000000000000	Header	
0x0026	00000000000000000000000000000000	Header	
0x0028	00000000000000000000000000000000	Header	
0x002a	00000000000000000000000000000000	Header	
0x002c	00000000000000000000000000000000	Header	
0x002e	00000000000000000000000000000000	Header	
0x0030	00000000000000000000000000000000	Header	
0x0032	00000000000000000000000000000000	Header	
0x0034	00000000000000000000000000000000	Header	
0x0036	00000000000000000000000000000000	Header	
0x0038	00000000000000000000000000000000	Header	
0x003a	00000000000000000000000000000000	Header	
0x003c	00000000000000000000000000000000	Header	
0x003e	00000000000000000000000000000000	Header	
0x0040	00000000000000000000000000000000	Header	
0x0042	00000000000000000000000000000000	Header	
0x0044	00000000000000000000000000000000	Header	
0x0046	00000000000000000000000000000000	Header	
0x0048	00000000000000000000000000000000	Header	
0x004a	00000000000000000000000000000000	Header	
0x004c	00000000000000000000000000000000	Header	
0x004e	00000000000000000000000000000000	Header	
0x0050	00000000000000000000000000000000	Header	
0x0052	00000000000000000000000000000000	Header	
0x0054	00000000000000000000000000000000	Header	
0x0056	00000000000000000000000000000000	Header	
0x0058	00000000000000000000000000000000	Header	
0x005a	00000000000000000000000000000000	Header	
0x005c	00000000000000000000000000000000	Header	
0x005e	00000000000000000000000000000000	Header	
0x0060	00000000000000000000000000000000	Header	
0x0062	00000000000000000000000000000000	Header	
0x0064	00000000000000000000000000000000	Header	
0x0066	00000000000000000000000000000000	Header	
0x0068	00000000000000000000000000000000	Header	
0x006a	00000000000000000000000000000000	Header	
0x006c	00000000000000000000000000000000	Header	
0x006e	00000000000000000000000000000000	Header	
0x0070	00000000000000000000000000000000	Header	
0x0072	00000000000000000000000000000000	Header	
0x0074	00000000000000000000000000000000	Header	
0x0076	00000000000000000000000000000000	Header	
0x0078	00000000000000000000000000000000	Header	
0x007a	00000000000000000000000000000000	Header	
0x007c	00000000000000000000000000000000	Header	
0x007e	00000000000000000000000000000000	Header	
0x0080	00000000000000000000000000000000	Header	
0x0082	00000000000000000000000000000000	Header	
0x0084	00000000000000000000000000000000	Header	
0x0086	00000000000000000000000000000000	Header	
0x0088	00000000000000000000000000000000	Header	
0x008a	00000000000000000000000000000000	Header	
0x008c	00000000000000000000000000000000	Header	
0x008e	00000000000000000000000000000000	Header	
0x0090	00000000000000000000000000000000	Header	
0x0092	00000000000000000000000000000000	Header	
0x0094	00000000000000000000000000000000	Header	
0x0096	00000000000000000000000000000000	Header	
0x0098	00000000000000000000000000000000	Header	
0x009a	00000000000000000000000000000000	Header	
0x009c	00000000000000000000000000000000	Header	
0x009e	00000000000000000000000000000000	Header	
0x00a0	00000000000000000000000000000000	Header	
0x00a2	00000000000000000000000000000000	Header	
0x00a4	00000000000000000000000000000000	Header	
0x00a6	00000000000000000000000000000000	Header	
0x00a8	00000000000000000000000000000000	Header	
0x00aa	00000000000000000000000000000000	Header	
0x00ac	00000000000000000000000000000000	Header	
0x00ae	00000000000000000000000000000000	Header	
0x00b0	00000000000000000000000000000000	Header	
0x00b2	00000000000000000000000000000000	Header	
0x00b4	00000000000000000000000000000000	Header	
0x00b6	00000000000000000000000000000000	Header	
0x00b8	00000000000000000000000000000000	Header	
0x00ba	00000000000000000000000000000000	Header	
0x00bc	00000000000000000000000000000000	Header	
0x00be	00000000000000000000000000000000	Header	
0x00c0	00000000000000000000000000000000	Header	
0x00c2	00000000000000000000000000000000	Header	
0x00c4	00000000000000000000000000000000	Header	
0x00c6	00000000000000000000000000000000	Header	
0x00c8	00000000000000000000000000000000	Header	
0x00ca	00000000000000000000000000000000	Header	
0x00cc	00000000000000000000000000000000	Header	
0x00ce	00000000000000000000000000000000	Header	
0x00d0	00000000000000000000000000000000	Header	
0x00d2	00000000000000000000000000000000	Header	
0x00d4	00000000000000000000000000000000	Header	
0x00d6	00000000000000000000000000000000	Header	
0x00d8	00000000000000000000000000000000	Header	
0x00da	00000000000000000000000000000000	Header	
0x00dc	00000000000000000000000000000000	Header	
0x00de	00000000000000000000000000000000	Header	
0x00e0	00000000000000000000000000000000	Header	
0x00e2	00000000000000000000000000000000	Header	
0x00e4	00000000000000000000000000000000	Header	
0x00e6	00000000000000000000000000000000	Header	
0x00e8	00000000000000000000000000000000	Header	
0x00ea	00000000000000000000000000000000	Header	
0x00ec	00000000000000000000000000000000	Header	
0x00ee	00000000000000000000000000000000	Header	
0x00f0	00000000000000000000000000000000	Header	
0x00f2	00000000000000000000000000000000	Header	
0x00f4	00000000000000000000000000000000	Header	
0x00f6	00000000000000000000000000000000	Header	
0x00f8	00000000000000000000000000000000	Header	
0x00fa	00000000000000000000000000000000	Header	
0x00fc	00000000000000000000000000000000	Header	
0x00fe	00000000000000000000000000000000	Header	
0x0100	00000000000000000000000000000000	Header	
0x0102	00000000000000000000000000000000	Header	
0x0104	00000000000000000000000000000000	Header	
0x0106	00000000000000000000000000000000	Header	
0x0108	00000000000000000000000000000000	Header	
0x010a	00000000000000000000000000000000	Header	
0x010c	00000000000000000000000000000000	Header	
0x010e	00000000000000000000000000000000	Header	
0x0110	00000000000000000000000000000000	Header	
0x0112	00000000000000000000000000000000	Header	
0x0114	00000000000000000000000000000000	Header	
0x0116	00000000000000000000000000000000	Header	
0x0118	00000000000000000000000000000000	Header	
0x011a	00000000000000000000000000000000	Header	
0x011c	00000000000000000000000000000000	Header	
0x011e	00000000000000000000000000000000	Header	
0x0120	00000000000000000000000000000000	Header	
0x0122	00000000000000000000000000000000	Header	
0x0124	00000000000000000000000000000000	Header	
0x0126	00000000000000		



# General Cache Organization (S, E, B)



## Satisfying Alignment with Structures

- Within structure:
- Must satisfy each element's alignment requirement
- Overall structure placement
- Each structure has alignment requirement  $K$
- $K$  = Largest alignment of any element
- Initial address & structure length must be multiples of  $K$

Example:

- $K = 8$ , due to double element

Diagram showing memory alignment for a structure with elements of size 4 bytes. The structure starts at address  $P+0$  and ends at  $P+24$ , with a total size of 24 bytes (multiple of 8).

## Locality

- Principle of Locality:** Programs tend to use data and instructions with addresses near or equal to those they have used recently
- Temporal locality:**
  - Recently referenced items are likely to be referenced again in the near future
- Spatial locality:**
  - Items with nearby addresses tend to be referenced close together in time

## Keeping Track of Free Blocks

- Method 1: Implicit list using length—links all blocks**
- Method 2: Explicit list among the free blocks using pointers**
- Method 3: Segregated free list**
  - Different free lists for different size classes
- Method 4: Blocks sorted by size**
  - Can use a balanced tree (e.g. Red-Black tree) with pointers within each free block, and the length used as a key

Parameter	Description
Fundamental parameters	
$S = 2^s$	Number of sets
$E$	Number of lines per set
$B = 2^b$	Block size (bytes)
$m = \log_2(M)$	Number of physical (main memory) address bits
Derived quantities	
$M = 2^m$	Maximum number of unique memory addresses
$s = \log_2(S)$	Number of set index bits
$b = \log_2(B)$	Number of block offset bits
$t = m - (s + b)$	Number of tag bits
$C = B \times E \times S$	Cache size (bytes), not including overhead such as the valid and tag bits

## Union Allocation

- Allocate according to largest element
- Can only use one field at a time

Diagram showing union allocation for a structure with elements of size 4 bytes. The structure starts at address  $P+0$  and ends at  $P+24$ , with a total size of 24 bytes (multiple of 8).

## General Caching Concepts: Types of Cache Misses

- Cold (compulsory) miss**
  - Cold misses occur because the cache is empty.
- Conflict miss**
  - Most caches limit blocks at level  $k+1$  to a small subset (sometimes a singleton) of the block positions at level  $k$ .
  - E.g. Block  $i$  at level  $k+1$  must be placed in block  $(i \bmod 4)$  at level  $k$ .
- Capacity miss**
  - Occurs when the set of active cache blocks (working set) is larger than the cache.

## Method 1: Implicit List

- For each block we need both size and allocation status
- Could store this information in two words: wasteful!
- Standard trick**
  - If blocks are aligned, some low-order address bits are always 0
  - Instead of storing an always-0 bit, use it as a allocated/free flag
  - When reading size word, must mask out this bit

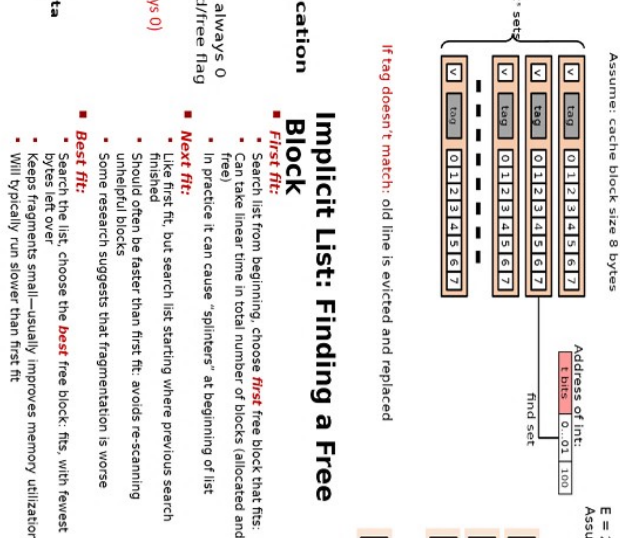


## FP Basics

- Arguments passed in  $\%xmm0, \%xmm1, \dots$
- Result returned in  $\%xmm0$
- All XMM registers caller-saved

Diagram showing floating-point register usage. A function call `float fadd(float x, float y)` returns the result in  $\%xmm0$ . The function body uses  $\%xmm0$  and  $\%xmm1$  for calculations.

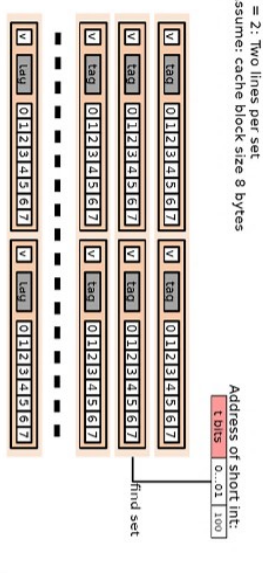
## Example: Direct Mapped Cache (E = 1)



## What about writes?

- Multiple copies of data exist:**
  - L1, L2, L3, Main Memory, Disk
- What to do on a write-hit?**
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
  - Need a dirty bit (line different from memory or not)
- What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
  - Good if more writes to the location follow
  - No-write-allocate (writes straight to memory, does not load into cache)
- Typical**
  - Write-through + No-write-allocate
  - Write-back + Write-allocate

## E-way Set Associative Cache (Here: E = 2)



## Nested Array Element Access

- Array Elements
- $A[i][j]$  is element of type  $T$ , which requires  $K$  bytes
- Address  $A + i * (C * K) + j * K$
- $= A + (i * C + j) * K$

Diagram showing nested array access. A 2D array  $A$  of type  $T$  (size  $K$  bytes) is accessed at row  $i$  and column  $j$ . The address is calculated as  $A + (i * C + j) * K$ .