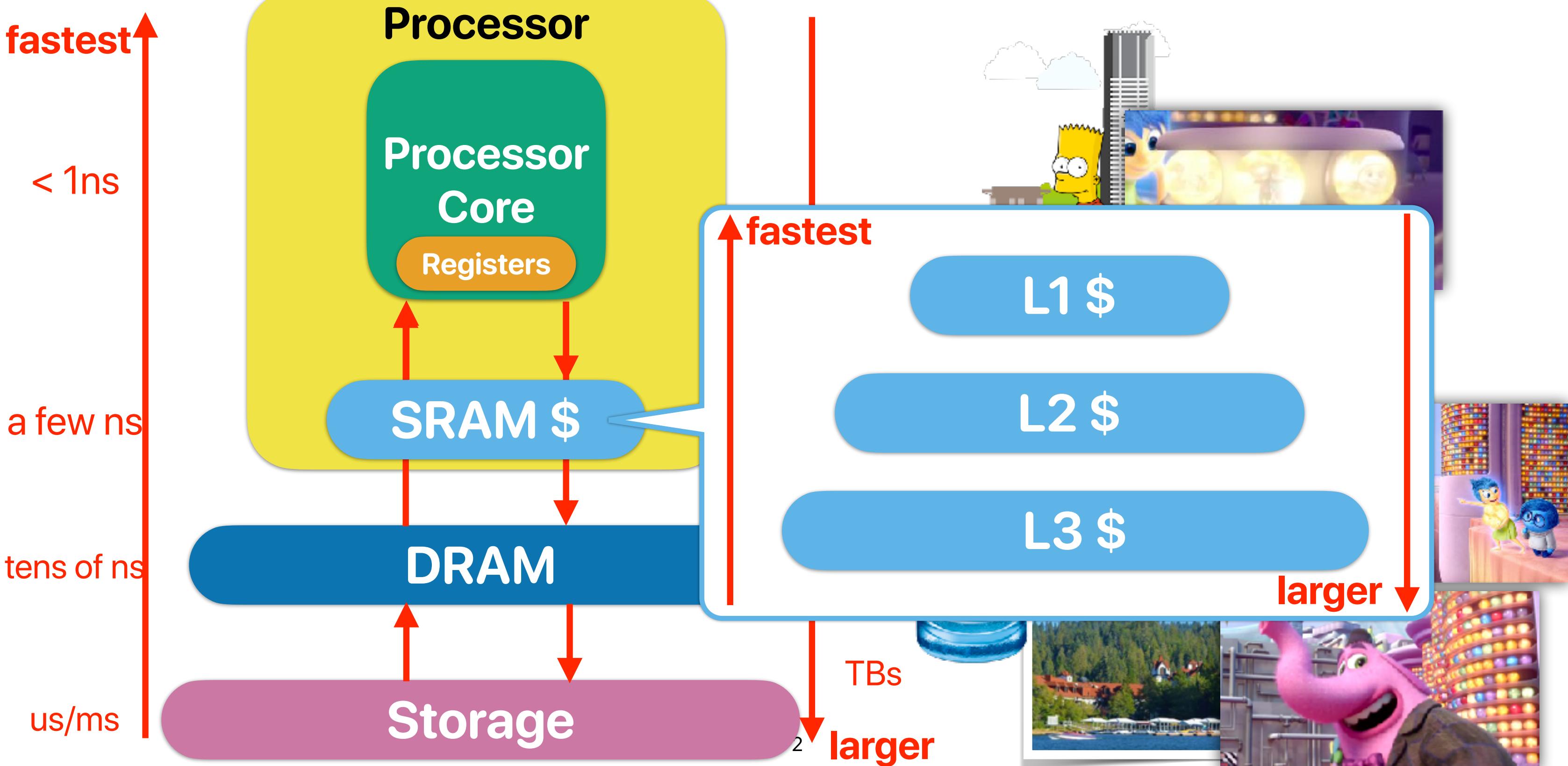


Memory Hierarchy (4): Cache Misses and How to Address Them (2)

Hung-Wei Tseng

Recap: Memory Hierarchy

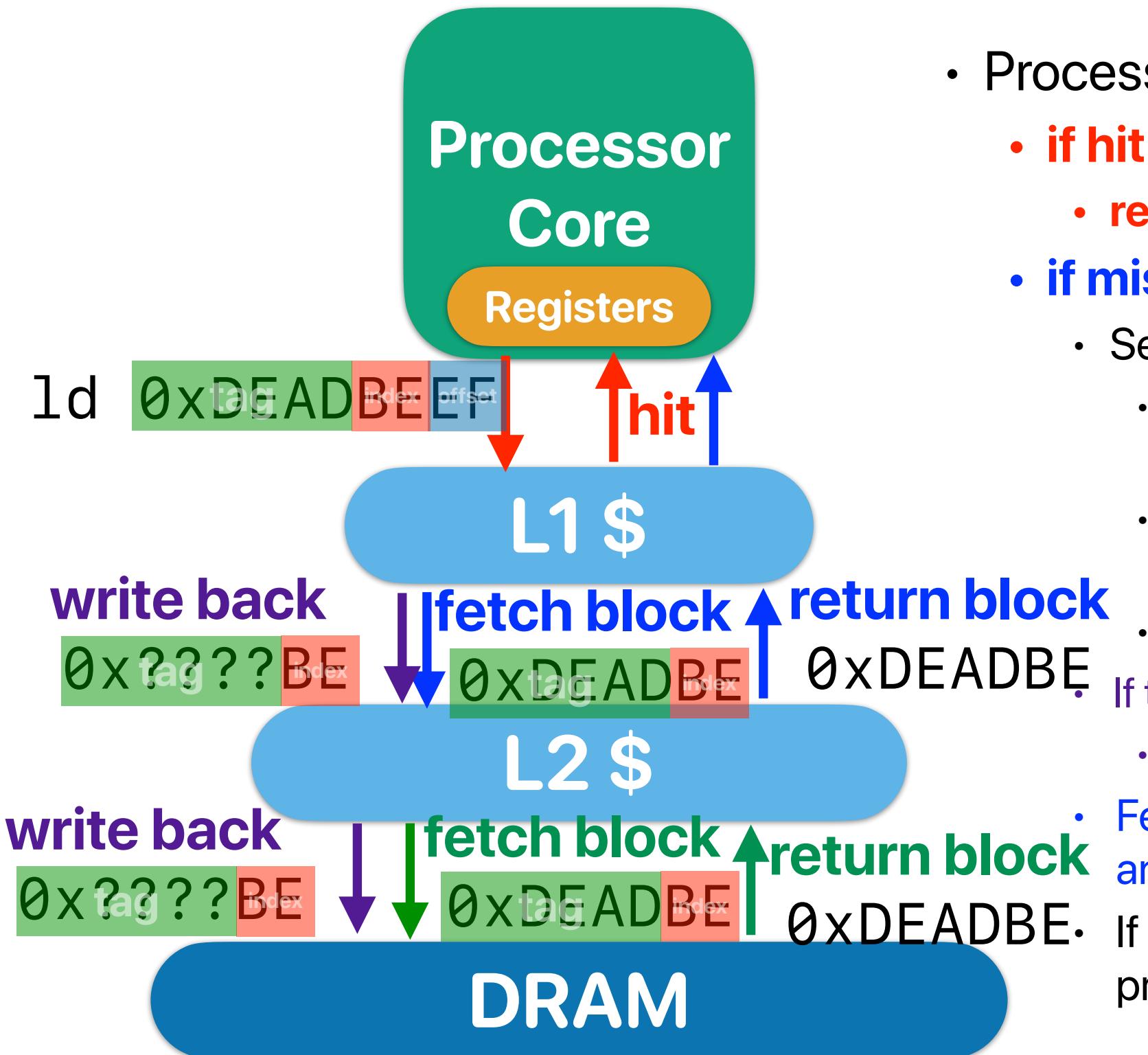


Recap: C = ABS

- **C:** Capacity in data arrays
- **A:** Way-Associativity — how many blocks within a set
 - N-way: N blocks in a set, A = N
 - 1 for direct-mapped cache
- **B:** Block Size (Cacheline)
 - How many bytes in a block
- **S:** Number of Sets:
 - A set contains blocks sharing the same index
 - 1 for fully associate cache
- number of bits in **block offset** — $\lg(B)$
- number of bits in **set index**: $\lg(S)$
- tag bits: **address_length** - $\lg(S)$ - $\lg(B)$
 - **address_length** is N bits for N-bit machines (e.g., 64-bit for 64-bit machines)
- $(\text{address} / \text{block_size}) \% S = \text{set index}$

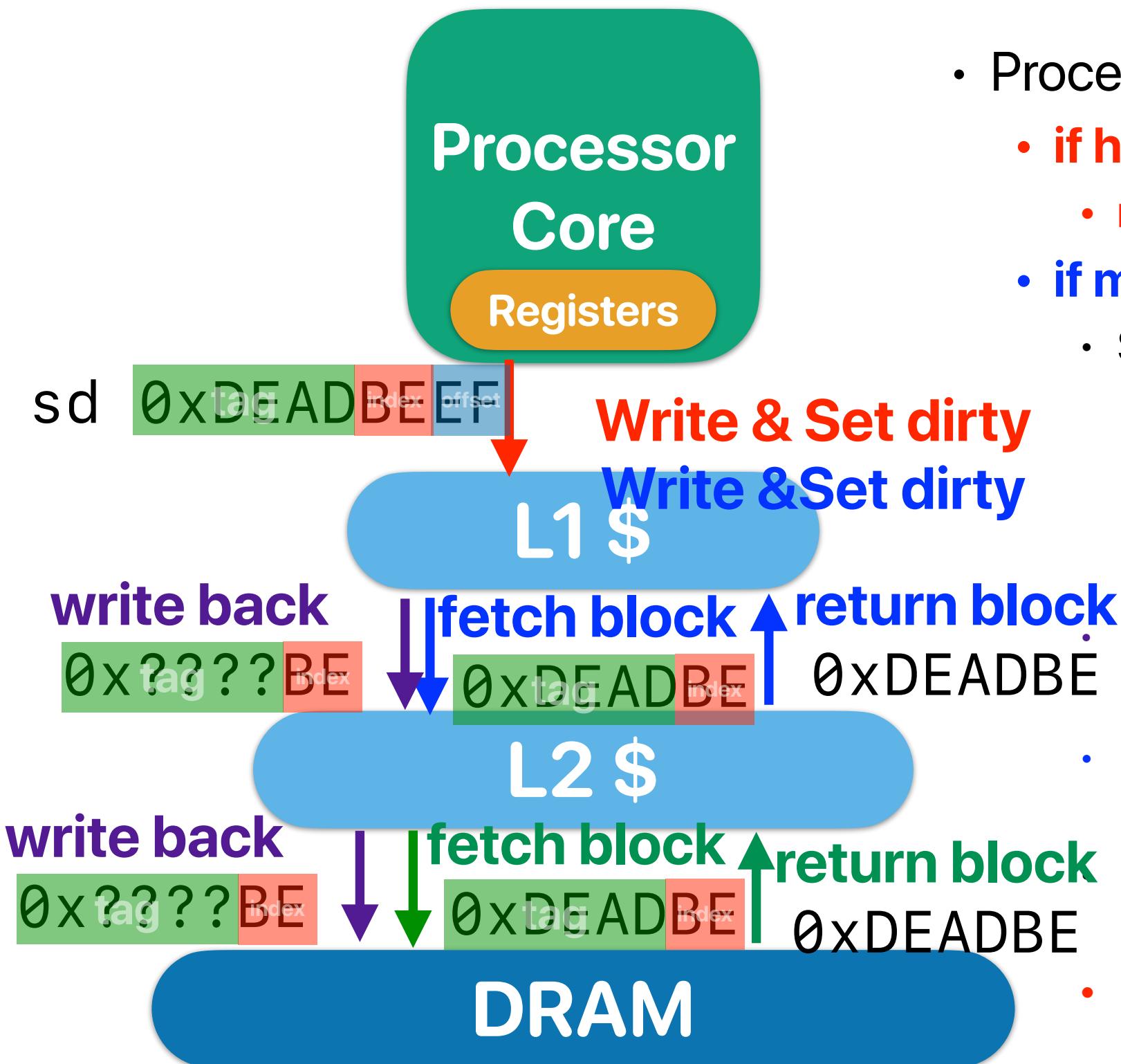


What happens when we read data



- Processor sends load request to L1-\$
 - if hit**
 - return data**
 - if miss**
 - Select a victim block
 - If the target “set” is not full — select an empty/invalidated block as the victim block
 - If the target “set” is full — select a victim block using some policy
 - LRU is preferred — to exploit temporal locality!
 - If the victim block is “dirty” & “valid”
 - Write back** the block to lower-level memory hierarchy
 - Fetch the requesting block from lower-level memory hierarchy and place in the victim block
- If write-back or fetching causes any miss, repeat the same process

What happens when we write data



- Processor sends load request to L1-\$
 - if hit**
 - return data — set DIRTY
 - if miss**
 - Select a victim block
 - If the target “set” is not full — select an empty/invalidated block as the victim block
 - If the target “set” is full — select a victim block using some policy
 - LRU is preferred — to exploit temporal locality!
 - If the victim block is “dirty” & “valid”
 - Write back** the block to lower-level memory hierarchy
 - Fetch the requesting block from lower-level memory hierarchy and place in the victim block
- If write-back or fetching causes any miss, repeat the same process
- Present the write “ONLY” in L1 and set DIRTY**

Recap: 3Cs of misses

- Compulsory miss
 - Cold start miss. First-time access to a block
- Capacity miss
 - The working set size of an application is bigger than cache size
- Conflict miss
 - Required data replaced by block(s) mapping to the same set
 - Similar collision in hash

NVIDIA Tegra X1

100% miss rate!

- Size 32KB, 4-way set associativity, 64B block, LRU policy, write-allocate, write-back, and assuming 64-bit address.

```
double a[8192], b[8192], c[8192], d[8192], e[8192];
/* a = 0x10000, b = 0x20000, c = 0x30000, d = 0x40000, e = 0x50000 */
for(i = 0; i < 512; i++) {
    e[i] = (a[i] * b[i] + c[i])/d[i];
    //load a[i], b[i], c[i], d[i] and then store to e[i]
```

C = ABS
 $32\text{KB} = 4 * 64 * S$
 $S = 128$
 $\text{offset} = \lg(64) = 6 \text{ bits}$
 $\text{index} = \lg(128) = 7 \text{ bits}$
 $\text{tag} = \text{the rest bits}$

	Address (Hex)	Address in binary	Tag	Index	Hit? Miss?	Replace?
a[0]	0x10000	0b000100000000000000000000000000	0x8	0x0	Compulsory Miss	
b[0]	0x20000	0b001000000000000000000000000000	0x10	0x0	Compulsory Miss	
c[0]	0x30000	0b001100000000000000000000000000	0x18	0x0	Compulsory Miss	
d[0]	0x40000	0b010000000000000000000000000000	0x20	0x0	Compulsory Miss	
e[0]	0x50000	0b010100000000000000000000000000	0x28	0x0	Compulsory Miss	a[0-7]
a[1]	0x10008	0b000100000000000000100000000000	0x8	0x0	Conflict Miss	b[0-7]
b[1]	0x20008	0b001000000000000000000000001000	0x10	0x0	Conflict Miss	c[0-7]
c[1]	0x30008	0b001100000000000000000000001000	0x18	0x0	Conflict Miss	d[0-7]
d[1]	0x40008	0b010000000000000000000000001000	0x20	0x0	Conflict Miss	e[0-7]
e[1]	0x50008	0b010100000000000000000000001000	0x28	0x0	Conflict Miss	a[0-7]

Loop optimizations

A

```
for(i = 0; i < ARRAY_SIZE; i++)  
{  
    for(j = 0; j < ARRAY_SIZE; j++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```

Loop interchange

B

```
for(j = 0; j < ARRAY_SIZE; j++)  
{  
    for(i = 0; i < ARRAY_SIZE; i++)  
    {  
        c[i][j] = a[i][j]+b[i][j];  
    }  
}
```

B

```
double a[8192], b[8192], c[8192], \  
       d[8192], e[8192];  
for(i = 0; i < 512; i++) {  
    e[i] = (a[i] * b[i] + c[i])/d[i];  
}
```

Loop fission

A

```
double a[8192], b[8192], c[8192], \  
       d[8192], e[8192];  
for(i = 0; i < 512; i++)  
    e[i] = a[i] * b[i] + c[i];  
for(i = 0; i < 512; i++)  
    e[i] /= d[i];
```

A

```
double a[8192], b[8192], c[8192], \  
       d[8192], e[8192];  
for(i = 0; i < 512; i++)  
    e[i] = a[i] * b[i] + c[i];  
for(i = 0; i < 512; i++)  
    e[i] /= d[i];
```

Loop fusion

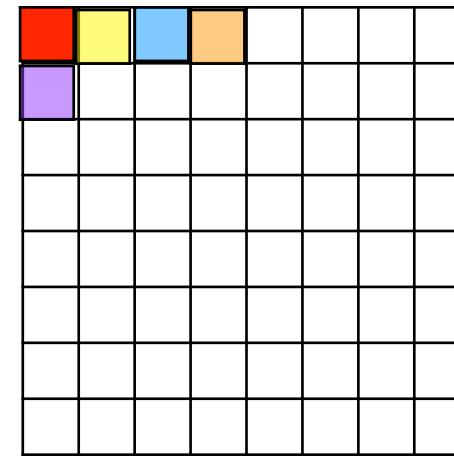
B

```
double a[8192], b[8192], c[8192], \  
       d[8192], e[8192];  
for(i = 0; i < 512; i++) {  
    e[i] = (a[i] * b[i] + c[i])/d[i];  
}
```

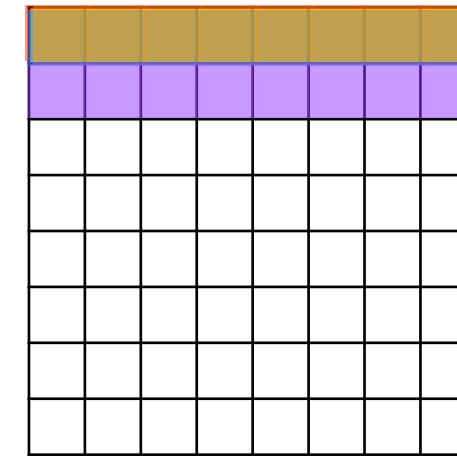
Matrix Multiplication

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

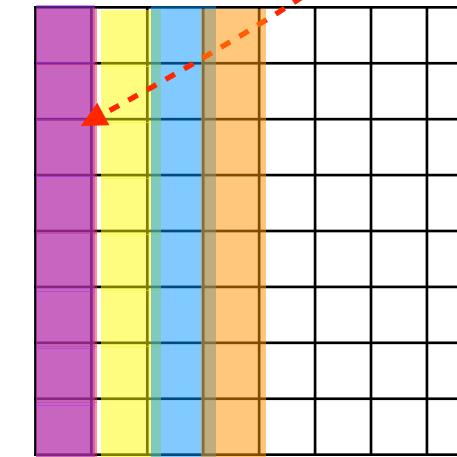
Very likely a miss if
array is large



c



a

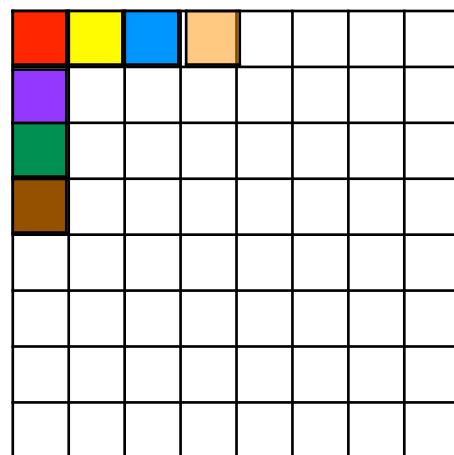
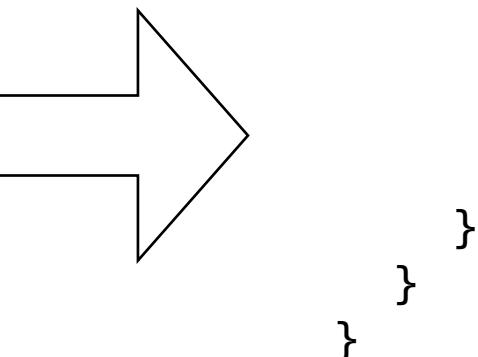


b

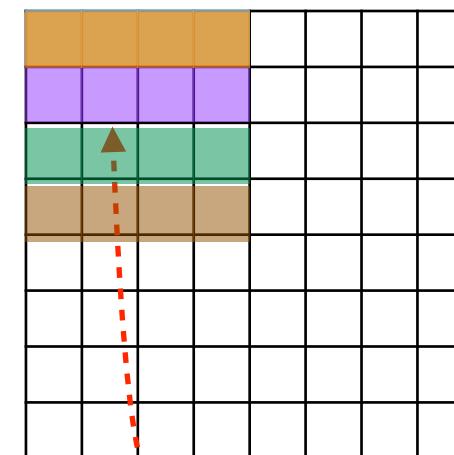
- If each dimension of your matrix is 2048
 - Each row takes 2048×8 bytes = 16KB
 - The L1 \$ of intel Core i7 is 48KB, 12-way, 64-byte blocked
 - You can only hold at most 3 rows/columns of each matrix!
 - You need the same row when j increase!

Blocking/tiling algorithm for matrix multiplications

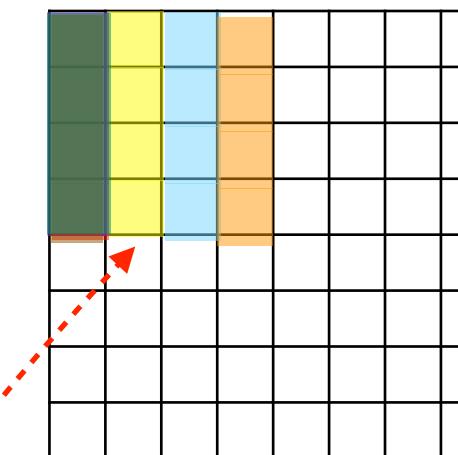
```
for(i = 0; i < ARRAY_SIZE; i++) {  
    for(j = 0; j < ARRAY_SIZE; j++) {  
        for(k = 0; k < ARRAY_SIZE; k++) {  
            c[i][j] += a[i][k]*b[k][j];  
        }  
    }  
}
```



c



a



b

You only need to hold these
sub-matrices in your cache

What kind(s) of misses can block algorithm remove?

- Comparing the naive algorithm and block algorithm on matrix multiplication, what kind of misses does block algorithm help to remove? (assuming an intel Core i7)

Naive

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

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Outline

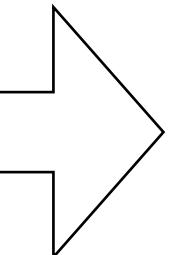
- Taxonomy/reasons of cache misses
- How to address cache misses: the software perspective

Matrix Transpose

```
// Transpose matrix b into b_t
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {
        b_t[i][j] += b[j][i];
    }
}

for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++) {
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++) {
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++) {
                        c[ii][jj] += a[ii][kk]*b[kk][jj];
                    }
                }
            }
        }
    }
}

// Compute on b_t
c[ii][jj] += a[ii][kk]*b_t[jj][kk];
```





What kind(s) of misses can matrix transpose remove?

- By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

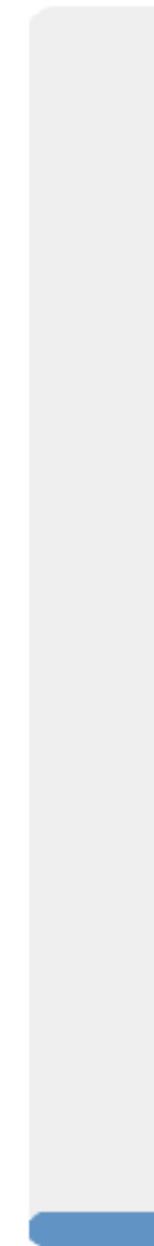
- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Block + Transpose

```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```

 0

0



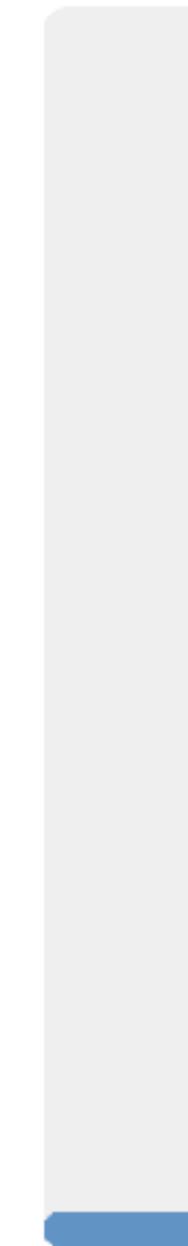
A

0



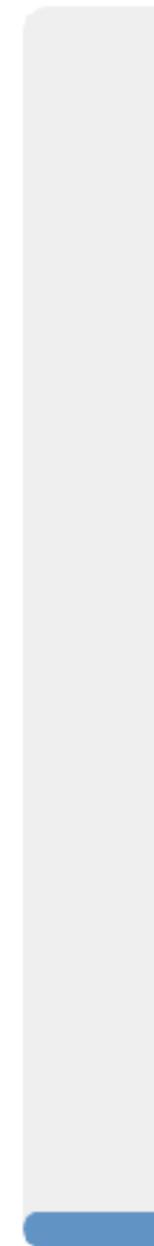
B

0



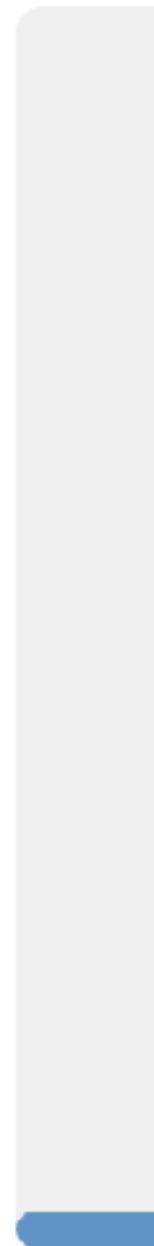
C

0



D

0



E



What kind(s) of misses can matrix transpose remove?

- By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

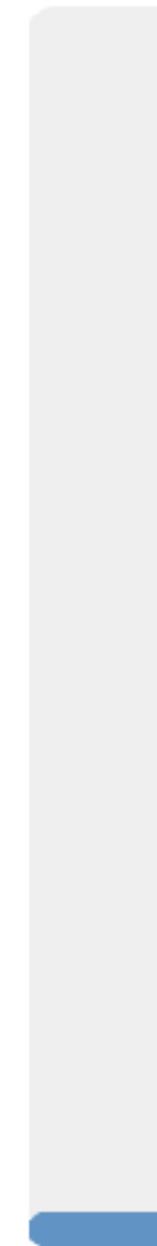
- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Block + Transpose

```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```

 0

0



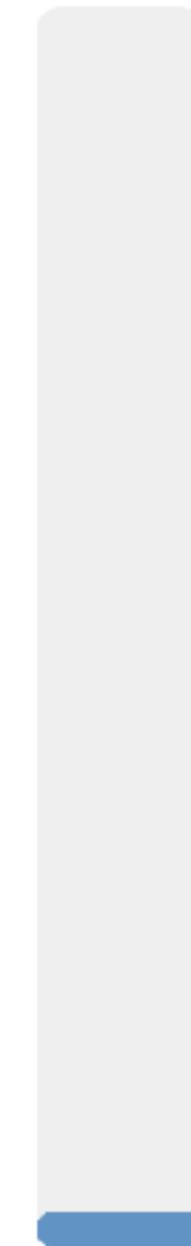
A

0



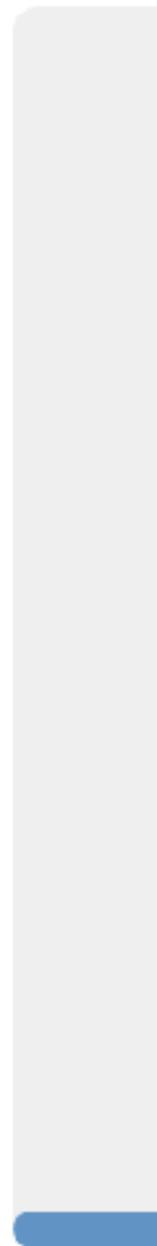
B

0



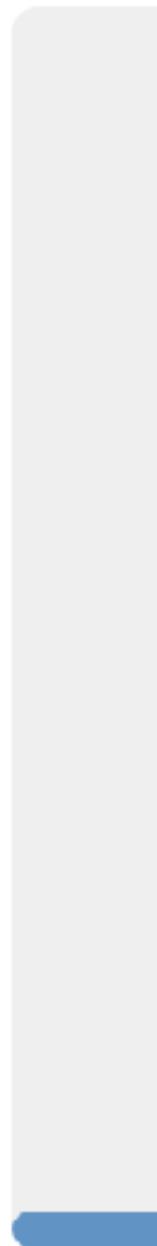
C

0



D

0



E

What kind(s) of misses can matrix transpose remove?

- By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Block + Transpose

```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```



Column-store or row-store

- Considering your the most frequently used queries in your database system are similar to

SELECT AVG(assignment_1) FROM table

Which of the following would be a data structure that better implements the table supporting this type of queries?

Array of objects	object of arrays
<pre>struct grades { int id; double *homework; double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre>	<pre>struct grades { int *id; double **homework; double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre>

- A. Array of objects
- B. Object of arrays

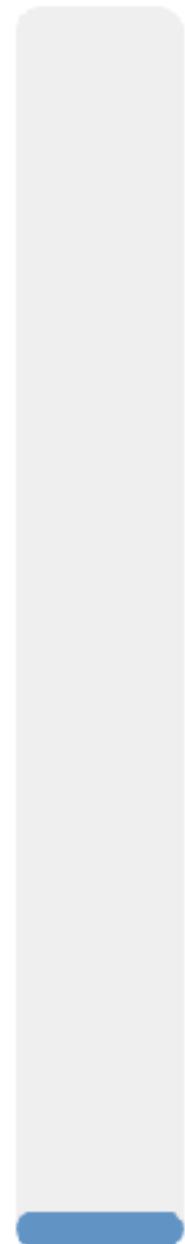
0

0



A

0



B



Column-store or row-store

- Considering your the most frequently used queries in your database system are similar to

SELECT AVG(assignment_1) FROM table

Which of the following would be a data structure that better implements the table supporting this type of queries?

Array of objects	object of arrays
<pre>struct grades { int id; double *homework; double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre>	<pre>struct grades { int *id; double **homework; double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre>

- A. Array of objects
- B. Object of arrays

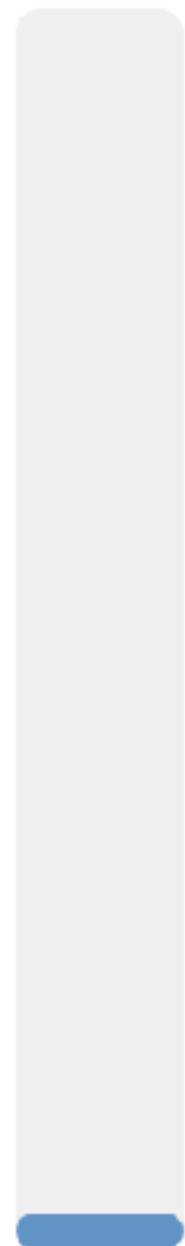
 0

0



A

0



B

Column-store or row-store

- Considering your the most frequently used queries in your database system are similar to

```
SELECT AVG(assignment_1) FROM table
```

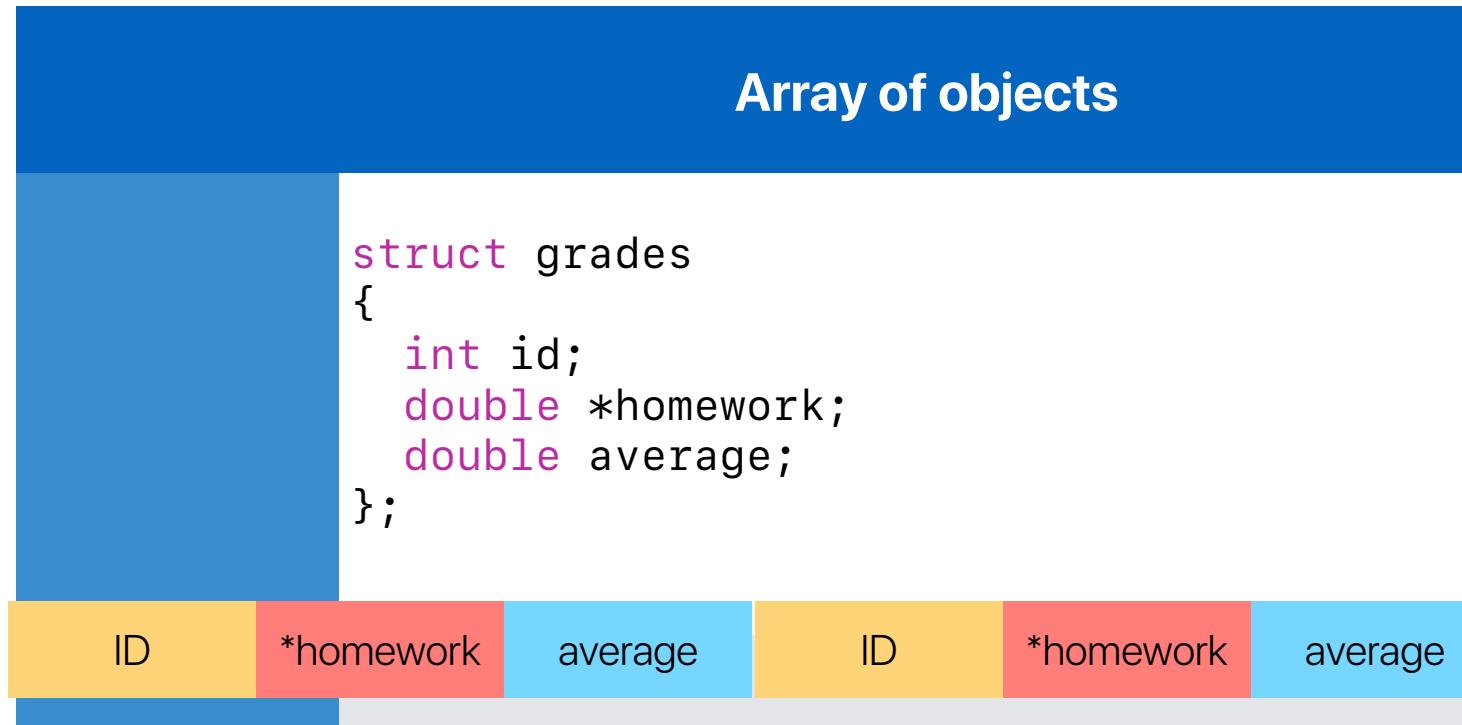
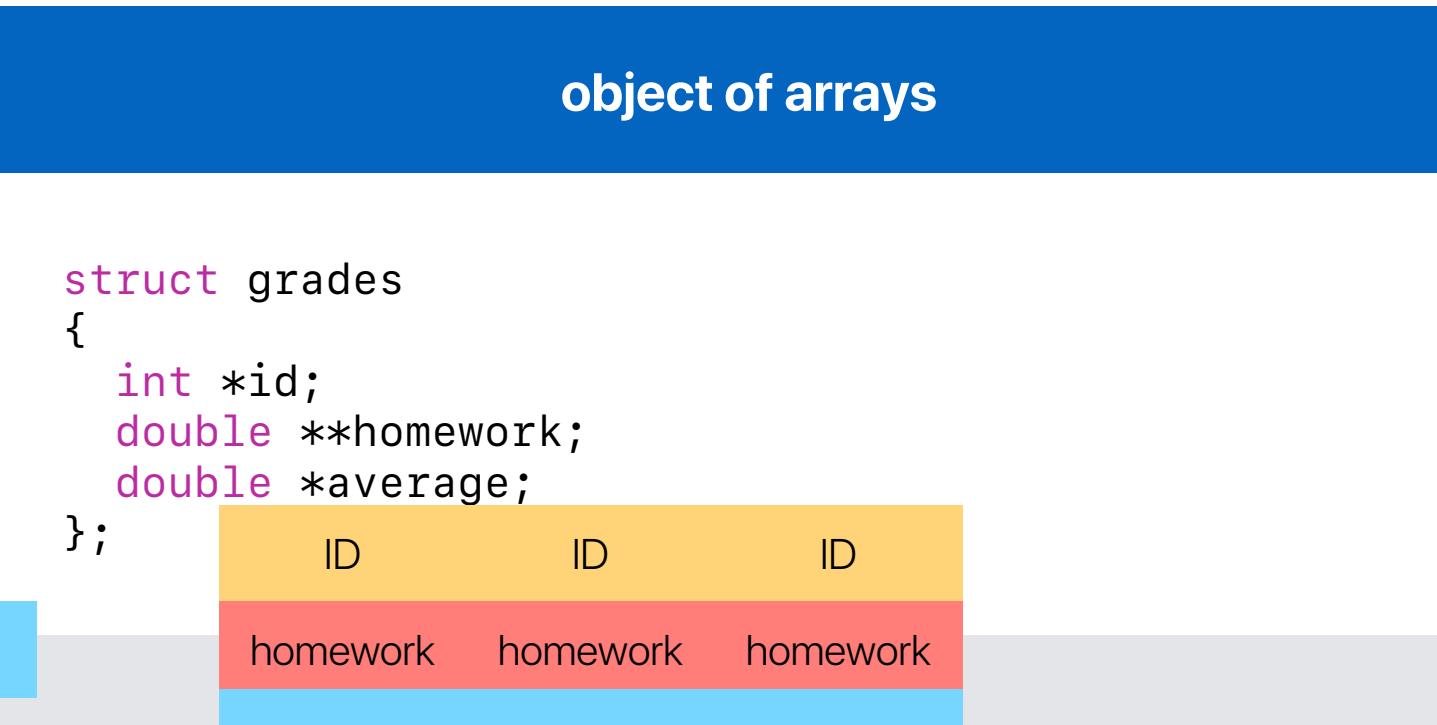
Which of the following would be a data structure that better implements the table supporting this type of queries?

Array of objects	object of arrays
<pre>struct grades { int id; double *homework; double average; }; table = (struct grades *) \ malloc(num_of_students*sizeof(struct</pre>	<pre>struct grades { int *id; double **homework; double *average; }; table =(struct grades *)malloc(sizeof(struct grades));</pre>

A. Array of objects **What if we want to calculate average scores for each student?**

B. Object of arrays

Array of structures or structure of arrays

	Array of objects	object of arrays									
	<pre>struct grades { int id; double *homework; double average; };</pre>  <pre>ID *homework average ID *homework average</pre>	<pre>struct grades { int *id; double **homework; double *average; };</pre>  <table border="1"><tr><td>ID</td><td>ID</td><td>ID</td></tr><tr><td>homework</td><td>homework</td><td>homework</td></tr><tr><td>average</td><td>average</td><td>average</td></tr></table>	ID	ID	ID	homework	homework	homework	average	average	average
ID	ID	ID									
homework	homework	homework									
average	average	average									
average of each homework	<pre>for(i=0;i<homework_items; i++) { gradesheet[total_number_students].homework[i] = 0.0; for(j=0;j<total_number_students;j++) gradesheet[total_number_students].homework[i] +=gradesheet[j].homework[i]; gradesheet[total_number_students].homework[i] /= (double)total_number_students; }</pre>	<pre>for(i = 0;i < homework_items; i++) { gradesheet.homework[i][total_number_students] = 0.0; for(j = 0; j <total_number_students;j++) { gradesheet.homework[i][total_number_students] += gradesheet.homework[i][j]; } gradesheet.homework[i][total_number_students] /= total_number_students; }</pre>									

Column-store or row-store

- If you're designing an in-memory database system, will you be using

RowId	Empld	Lastname	Firstname	Salary
1	10	Smith	Joe	40000
2	12	Jones	Mary	50000
3	11	Johnson	Cathy	44000
4	22	Jones	Bob	55000

- column-store — stores data tables column by column

10:001,12:002,11:003,22:004;

Smith:001,Jones:002,Johnson:003,Jones:004,

Joe:001,Mary:002,Cathy:003,Bob:004;

40000:001,50000:002,44000:003,55000:004;

if the most frequently used query looks like –
select Lastname, Firsntname from table

- row-store — stores data tables row by row

001:10,Smith,Joe,40000;

002:12,Jones,Mary,50000;

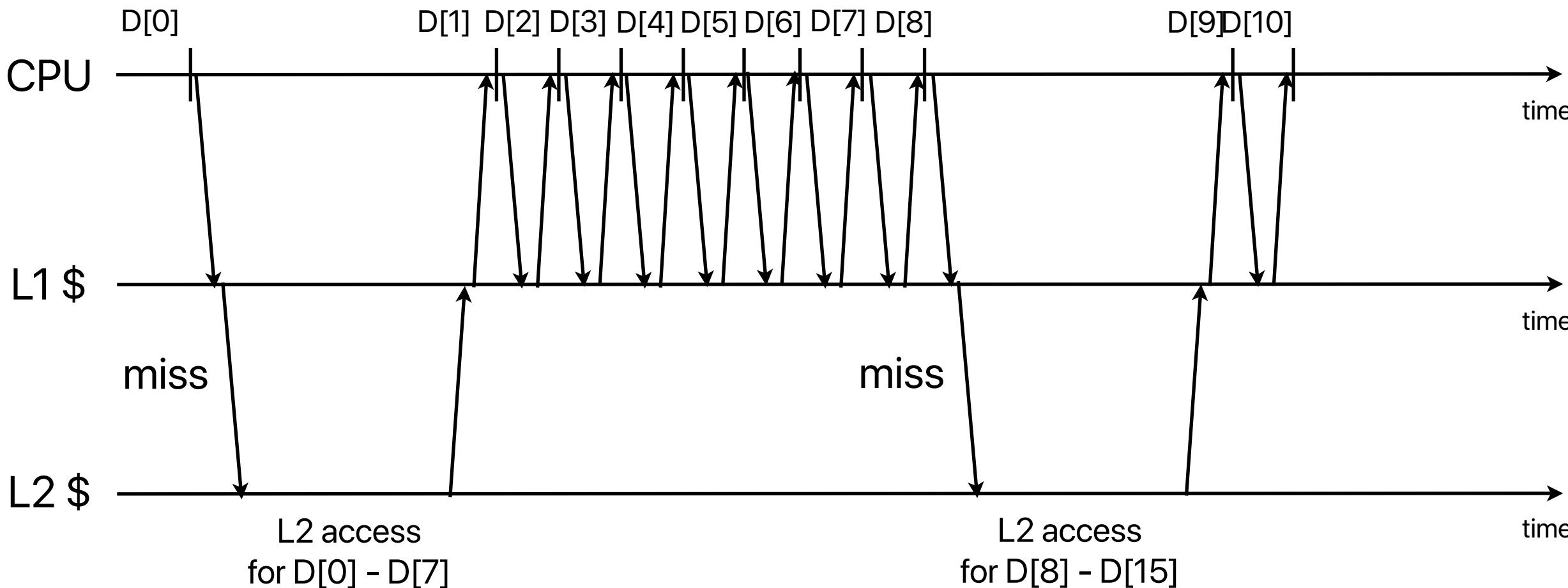
003:11,Johnson,Cathy,44000;

004:22,Jones,Bob,55000;

Prefetching

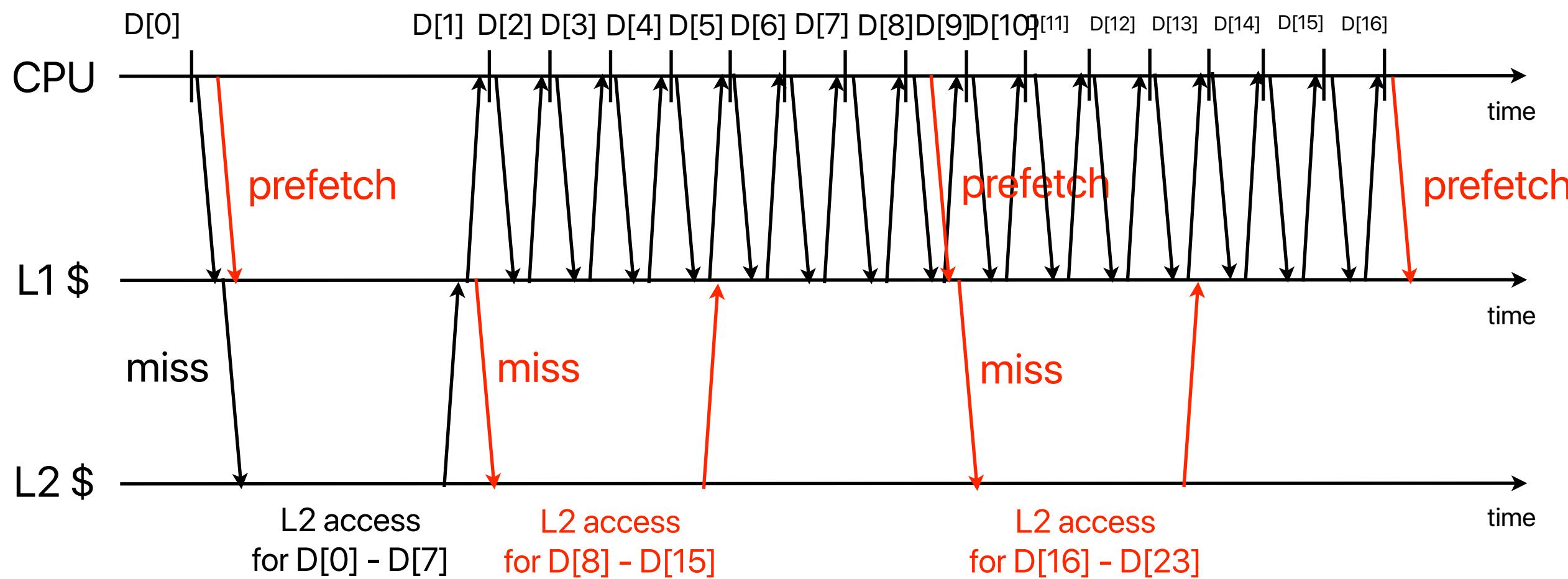
Characteristic of memory accesses

```
for(i = 0; i < 1000000; i++) {  
    D[i] = rand();  
}
```



Prefetching

```
for(i = 0; i < 1000000; i++) {  
    D[i] = rand();  
    // prefetch D[i+8] if i % 8 == 0  
}
```



Prefetching

- Identify the access pattern and proactively fetch data/instruction before the application asks for the data/instruction
 - Trigger the cache miss earlier to eliminate the miss when the application needs the data/instruction
- Hardware prefetch
 - The processor can keep track the distance between misses. If there is a pattern, fetch `miss_data_address+distance` for a miss
- Software prefetch
 - Load data into some register
 - Using prefetch instructions

Demo

- x86 provide prefetch instructions
- As a programmer, you may insert `_mm_prefetch` in x86 programs to perform software prefetch for your code
- gcc also has a flag “`-fprefetch-loop-arrays`” to automatically insert software prefetch instructions



Where can prefetch work effectively?

- How many of the following code snippet can “prefetching” effectively help improving performance?

(1)
while(node){
 node = node->next;
}

(3)
while (root != NULL){
 if (key > root->data)
 root = root->right;

 else if (key < root->data)
 root = root->left;
 else
 return true;
}

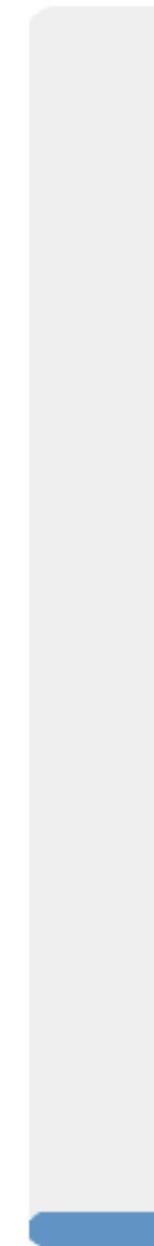
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

(2)
while(++i<100000)
 a[i]=rand();

(4)
for (i = 0; i < 65536; i++) {
 mix_i = ((i * 167) + 13) & 65536;
 results[mix_i]++;
}

 0

0



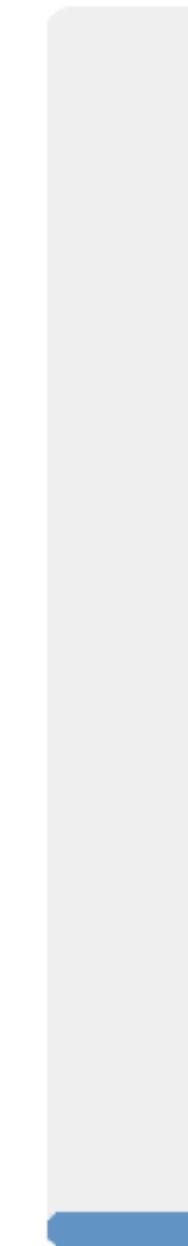
A

0



B

0



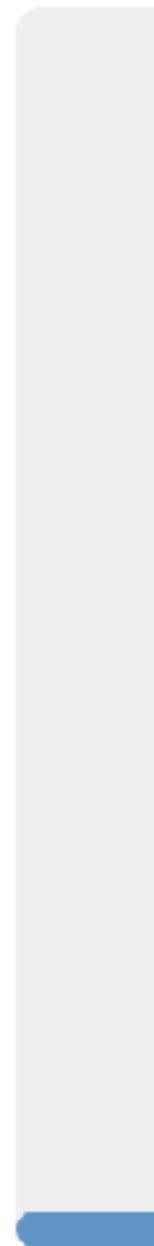
C

0



D

0



E



Where can prefetch work effectively?

- How many of the following code snippet can “prefetching” effectively help improving performance?

(1)
while(node){
 node = node->next;
}

(3)
while (root != NULL){
 if (key > root->data)
 root = root->right;

 else if (key < root->data)
 root = root->left;
 else
 return true;
}

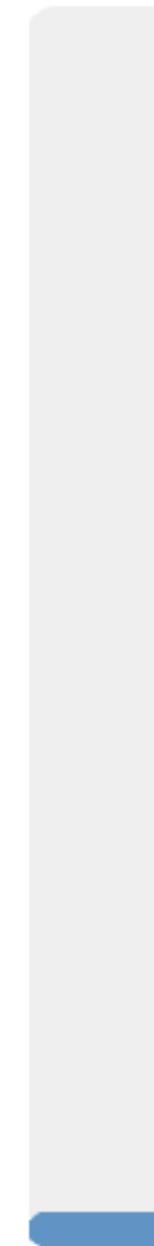
- A. 0
- B. 1
- C. 2
- D. 3
- E. 4

(2)
while(++i<100000)
 a[i]=rand();

(4)
for (i = 0; i < 65536; i++) {
 mix_i = ((i * 167) + 13) & 65536;
 results[mix_i]++;
}

 0

0



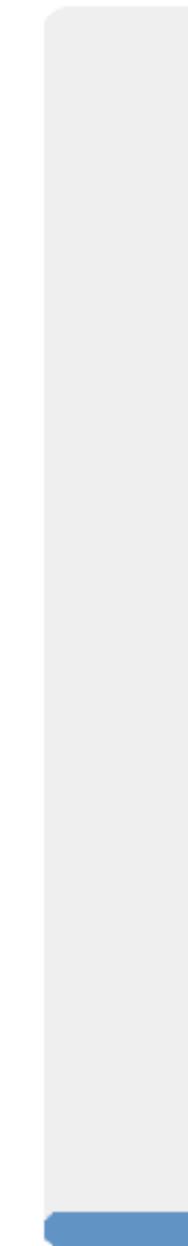
A

0



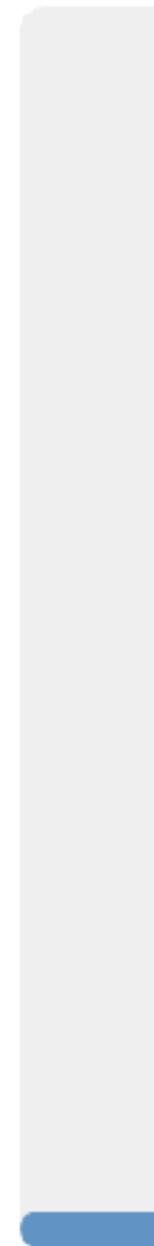
B

0



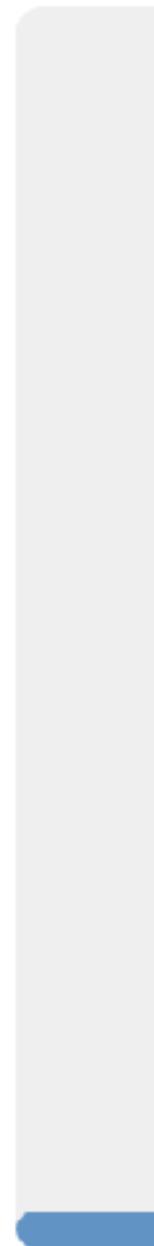
C

0



D

0



E

Where can prefetch work effectively?

- How many of the following code snippet can “prefetching” effectively help improving performance?

(1)
while(node){
 node = node->next;
}

— where the next pointing to is hard to predict

(3)
while (root != NULL){
 if (key > root->data)
 root = root->right;

 else if (key < root->data)
 root = root->left;
 else
 return true;
}

A. 0

B. 1

C. 2

D. 3

E. 4

(2) 
while(++i<100000)
 a[i]=rand();

(4)
for (i = 0; i < 65536; i++) {
 mix_i = ((i * 167) + 13) & 65536;
 results[mix_i]++;
}

— the stride to the next element is hard to predict...

— where the next node is also hard to predict

Data structures



The result of sizeof(struct student)

- Consider the following data structure:

```
struct student {  
    int id;  
    double *homework;  
    int participation;  
    double midterm;  
    double average;  
};
```

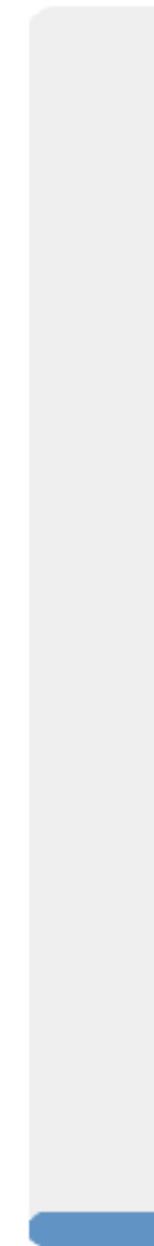
What's the output of

```
printf("%lu\n", sizeof(struct student));
```

- A. 20
- B. 28
- C. 32
- D. 36
- E. 40

 0

0



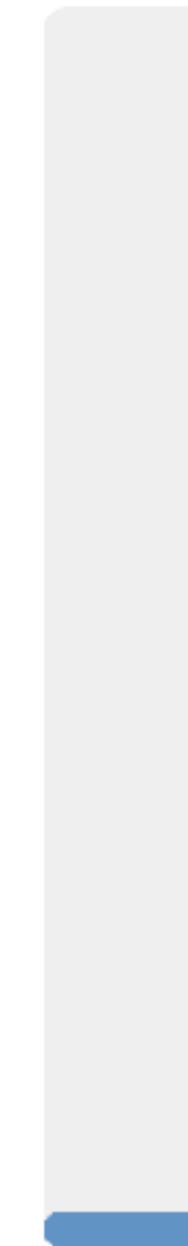
A

0



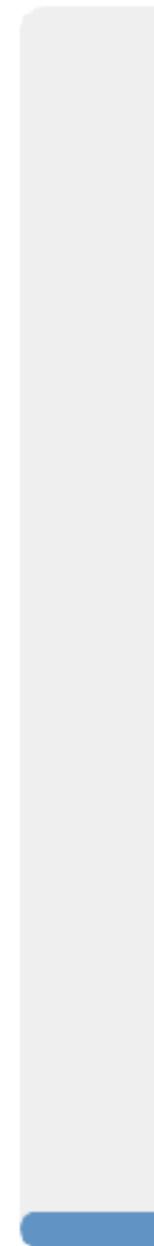
B

0



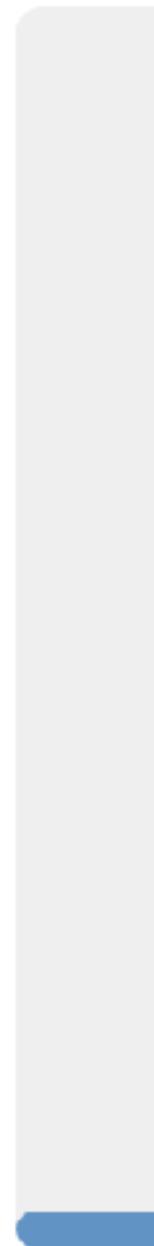
C

0



D

0



E



The result of sizeof(struct student)

- Consider the following data structure:

```
struct student {  
    int id;  
    double *homework;  
    int participation;  
    double midterm;  
    double average;  
};
```

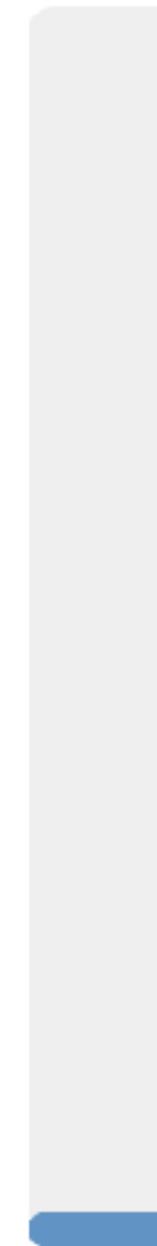
What's the output of

```
printf("%lu\n", sizeof(struct student));
```

- A. 20
- B. 28
- C. 32
- D. 36
- E. 40

 0

0



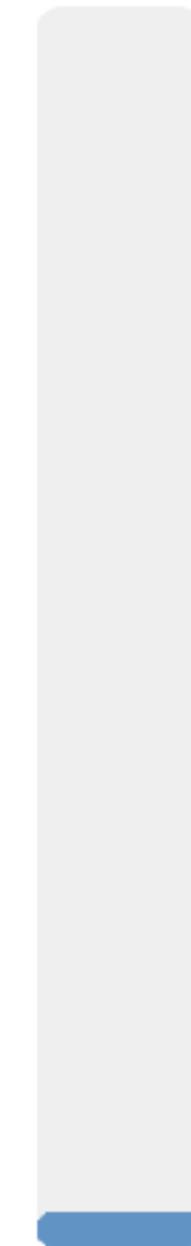
A

0



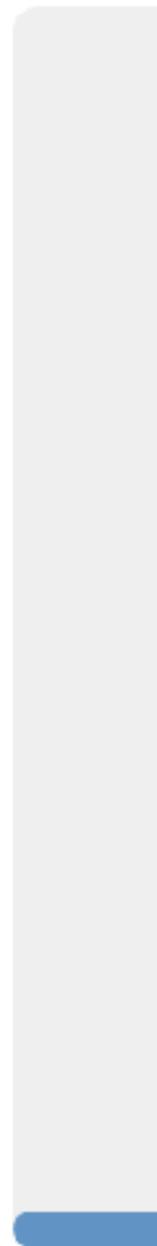
B

0



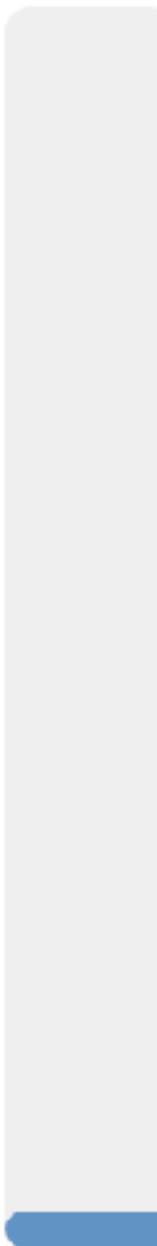
C

0



D

0



E

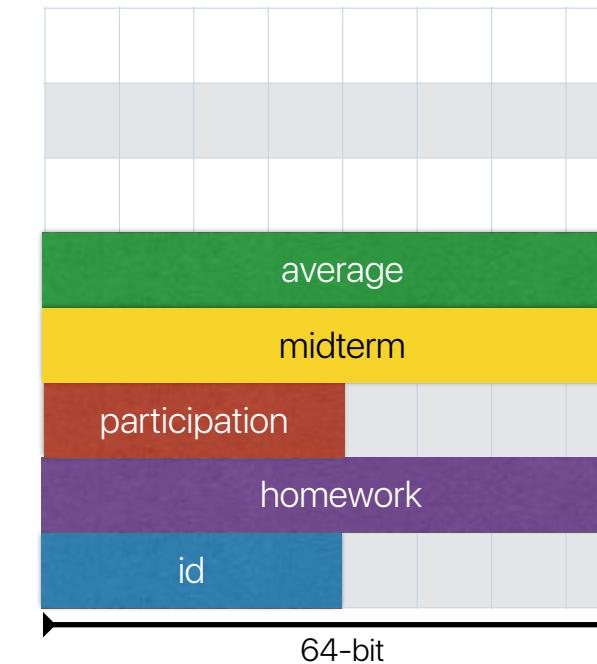
Memory addressing/alignment

- Almost every popular ISA architecture uses “byte-addressing” to access memory locations
- Instructions generally work faster when the given memory address is aligned
 - Aligned — if an instruction accesses an object of size n at address X , the access is **aligned** if $X \bmod n = 0$.
 - Some architecture/processor does not support aligned access at all
 - Therefore, compilers only allocate objects on “aligned” address

The result of sizeof(struct student)

- Consider the following data structure:

```
struct student {  
    int id;  
    double *homework;  
    int participation;  
    double midterm;  
    double average;  
};
```



What's the output of
`printf("%lu\n", sizeof(struct student))?`

- A. 20
- B. 28
- C. 32
- D. 36
- E. 40

Summary of Software Optimizations

- Data layout — capacity miss, conflict miss, compulsory miss
 - Matrix transpose
 - Column-store vs. row-store
- Blocking/tiling — capacity miss, conflict miss
 - Matrix multiplications
- Loop fission — conflict miss — when \$ has limited way associativity
- Loop fusion — capacity miss — when \$ has enough way associativity
- Loop interchange — conflict/capacity miss

**Let's take a look of another aspect
of memory systems**

Let's dig into this code

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```

What will happen?

- If we execute the code on the right-hand side code on a machine with only 32 GB of physical memory installed and the dim is "70000" (requires $70000 \times 70000 \times 8$ bytes ~ 37 GB memory at least), What will happen?
 - A. The program will crash in one of the malloc function call
 - B. The program will crash due to a "segmentation fault" that caused by accessing NULL pointer
 - C. The program will be killed automatically by the OS as it uses more than installed physical main memory
 - D. The program will finish without any issue

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```



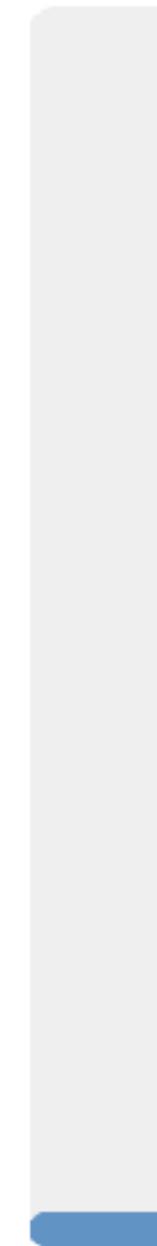
What will happen?

- If we execute the code on the right-hand side code on a machine with only 32 GB of physical memory installed and the dim is "70000" (requires $70000 \times 70000 \times 8$ bytes ~ 37 GB memory at least), What will happen?
 - A. The program will crash in one of the malloc function call
 - B. The program will crash due to a "segmentation fault" that caused by accessing NULL pointer
 - C. The program will be killed automatically by the OS as it uses more than installed physical main memory
 - D. The program will finish without any issue

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```

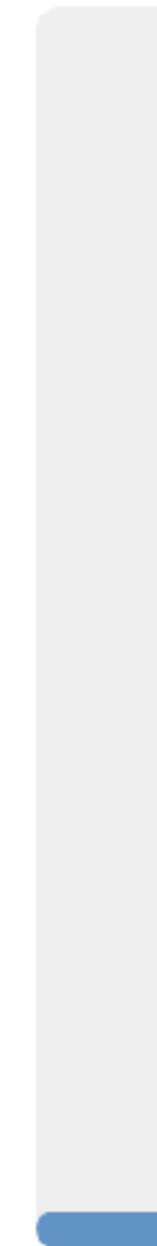
 0

0



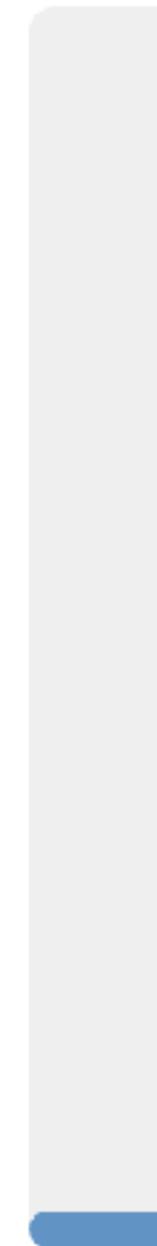
A

0



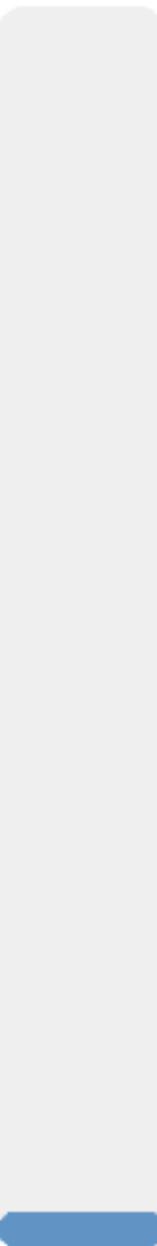
B

0



C

0



D



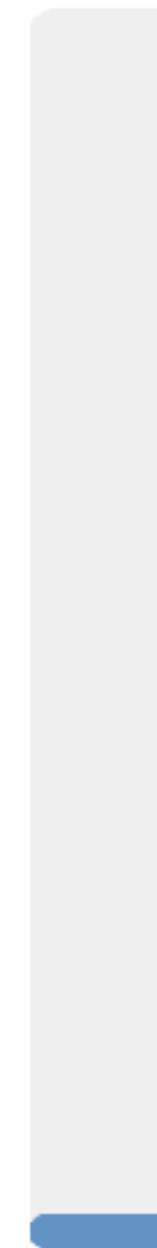
What will happen?

- If we execute the code on the right-hand side code on a machine with only 32 GB of physical memory installed and the dim is "70000" (requires $70000 \times 70000 \times 8$ bytes ~ 37 GB memory at least), What will happen?
 - A. The program will crash in one of the malloc function call
 - B. The program will crash due to a "segmentation fault" that caused by accessing NULL pointer
 - C. The program will be killed automatically by the OS as it uses more than installed physical main memory
 - D. The program will finish without any issue

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```

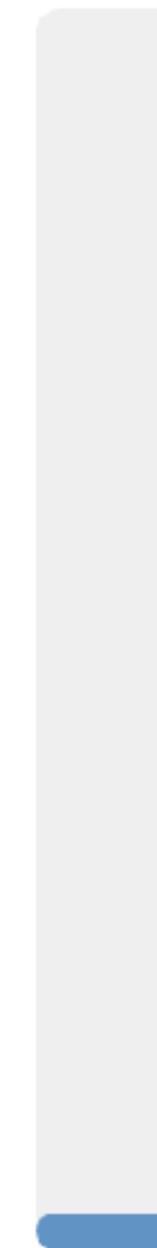
 0

0



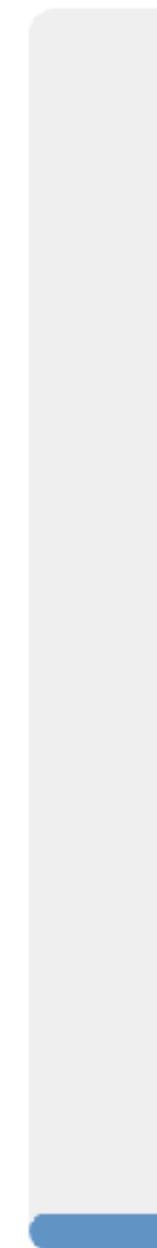
A

0



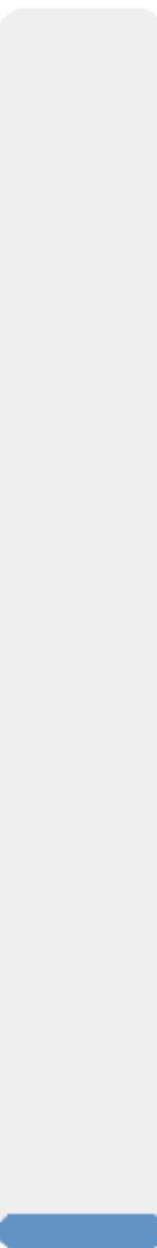
B

0



C

0



D

```
終端機 — -tosh — 102x25
top ... -tosh + E

#include <stdlib.h>
#include <cassert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

//#define dim 32768
//#define dim 49152
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
File memory_allocation.c not changed so no update needed
BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory] -bunny- ./memory_allocation

BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory] -bunny- □
```

What will happen?

- If we execute the code on the right-hand side code on a machine with only 32 GB of physical memory installed and the dim is "70000" (requires $70000 \times 70000 \times 8$ bytes ~ 37 GB memory at least), What will happen?
 - A. The program will crash in one of the malloc function call
 - B. The program will crash due to a "segmentation fault" that caused by accessing NULL pointer
 - C. The program will be killed automatically by the OS as it uses more than installed physical main memory
 - D. The program will finish without any issue

```
int main(int argc, char *argv[])
{
    int i,j;
    double **a;
    double sum=0, average;
    int dim=32768;
    if(argc < 2)
    {
        fprintf(stderr, "Usage: %s dimension\n", argv[0]);
        exit(1);
    }
    dim = atoi(argv[1]);
    a = (double **)malloc(sizeof(double *)*dim);
    for(i = 0 ; i < dim; i++)
        a[i] = (double *)malloc(sizeof(double)*dim);
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            a[i][j] = rand();
    for(i = 0 ; i < dim; i++)
        for(j = 0 ; j < dim; j++)
            sum+=a[i][j];
    average = sum/(dim*dim);
    fprintf(stderr,"average: %lf\n",average);
    for(i = 0 ; i < dim; i++)
        free(a[i]);
    free(a);
    return 0;
}
```

Let's dig into this code

```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    // Create processes
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    // Generate rand seed
    srand((int)time(NULL)+(int)getpid());
    a = rand();
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n", getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n", getpid(), a, &a);
    return 0;
}
```



Consider the following code ...

- Consider the case when we run multiple instances of the given program at the same time on modern machines, which pair of statements is correct?
 - ① The printed “address of a” is the same for every running instances
 - ② The printed “address of a” is different for each instance
 - ③ All running instances will print the same value of a
 - ④ Some instances will print the same value of a
 - ⑤ Each instance will print a different value of a

A. (1) & (3)
B. (1) & (4)
C. (1) & (5)
D. (2) & (3)
E. (2) & (4)

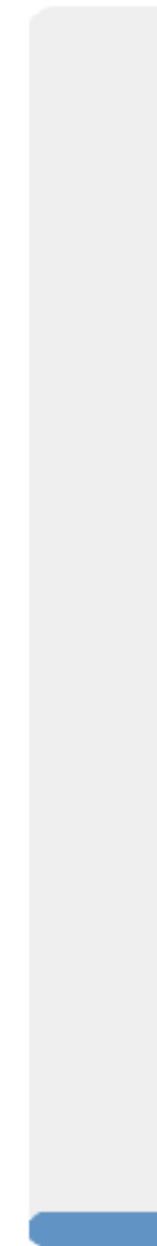
```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

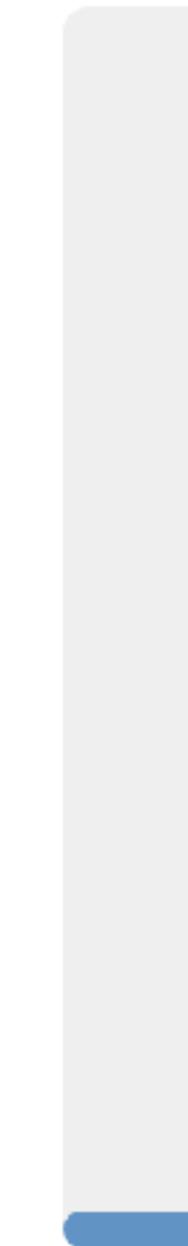
int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    srand((int)time(NULL)+(int)getpid());
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    return 0;
}
```

 0

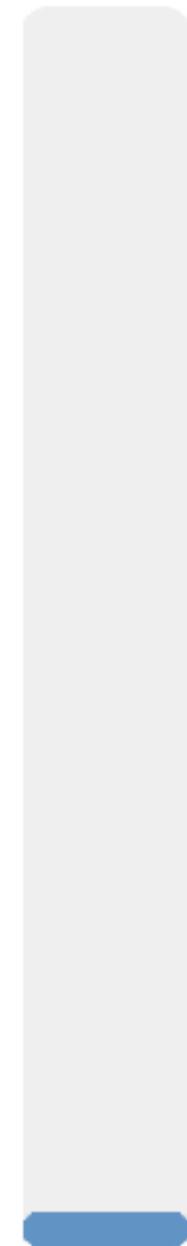
0



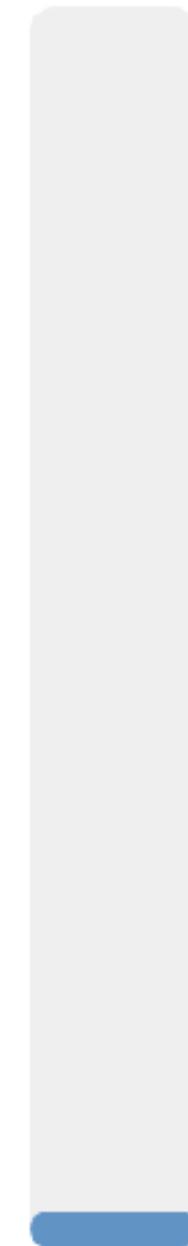
0



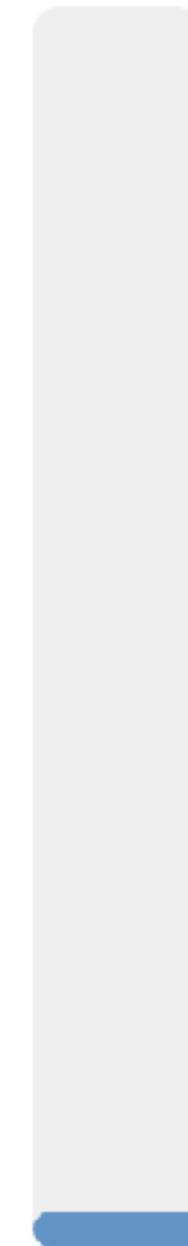
0



0



0



A

B

C

D

E



Consider the following code ...

- Consider the case when we run multiple instances of the given program at the same time on modern machines, which pair of statements is correct?
 - ① The printed “address of a” is the same for every running instances
 - ② The printed “address of a” is different for each instance
 - ③ All running instances will print the same value of a
 - ④ Some instances will print the same value of a
 - ⑤ Each instance will print a different value of a

A. (1) & (3)
B. (1) & (4)
C. (1) & (5)
D. (2) & (3)
E. (2) & (4)

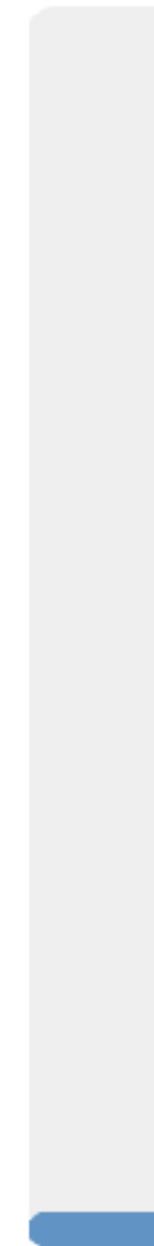
```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    srand((int)time(NULL)+(int)getpid());
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    return 0;
}
```

 0

0



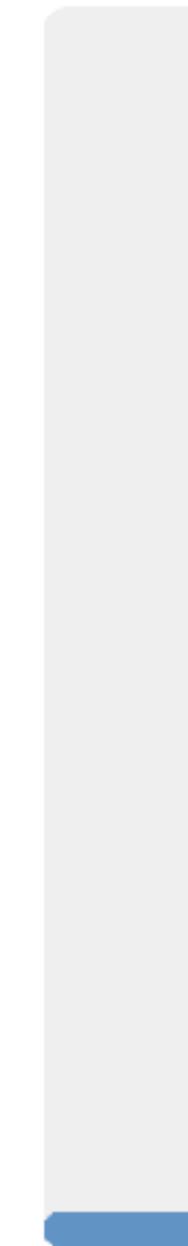
A

0



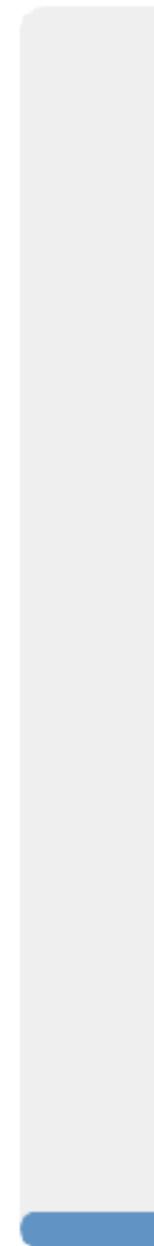
B

0



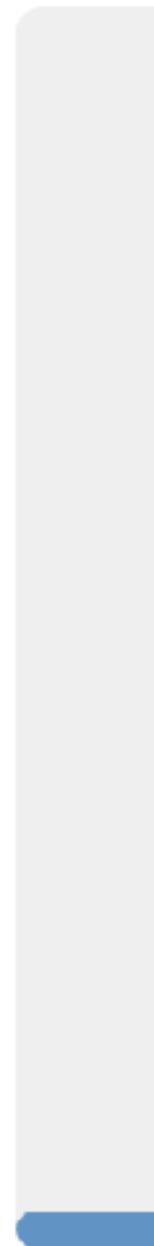
C

0



D

0



E

終端機 — tcsh — 102x25
[BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory]] -bunny- make clean
rm -f virtualization
BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory]] -bunny- []

Demo revisited

```
sleep(10);
fprintf(stderr, "\nProcess %d: Value of a is %lf and address of a is %p\n", (int)getpid(), a, &a);
return 0;
}
```

File virtualization.c not changed so no update needed

```
BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory] -bunny- make
gcc -O3 virtualization.c -o virtualization
```

```
BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory] -bunny- ./virtualization 4
```

```
Process 19719: Value of a is 1671139616.000000 and address of a is 0x104967050
```

```
Process 19720: Value of a is 1671156423.000000 and address of a is 0x104967050
```

```
Process 19718: Value of a is 1671122809.000000 and address of a is 0x104967050
```

```
Process 19721: Value of a is 1671173230.000000 and address of a is 0x104967050
    Different values
```

```
Process 19719: Value of a is 1671139616.000000 and address of a is 0x104967050
```

```
Process 19721: Value of a is 1671173230.000000 and address of a is 0x104967050
```

```
Process 19720: Value of a is 1671156423.000000 and address of a is 0x104967050
```

```
Process 19718: Value of a is 1671122809.000000 and address of a is 0x104967050
```

```
BunnyMACProRetina [/Users/bunny/Dropbox/CSE203/GitHub/demo/virtual_memory] -bunny-
```

Different values are
preserved

The same memory
address!

Consider the following code ...

- Consider the case when we run multiple instances of the given program at the same time on modern machines, which pair of statements is correct?
 - The printed “address of a” is the same for every running instances
 - The printed “address of a” is different for each instance
 - All running instances will print the same value of a
 - Some instances will print the same value of a
 - Each instance will print a different value of a

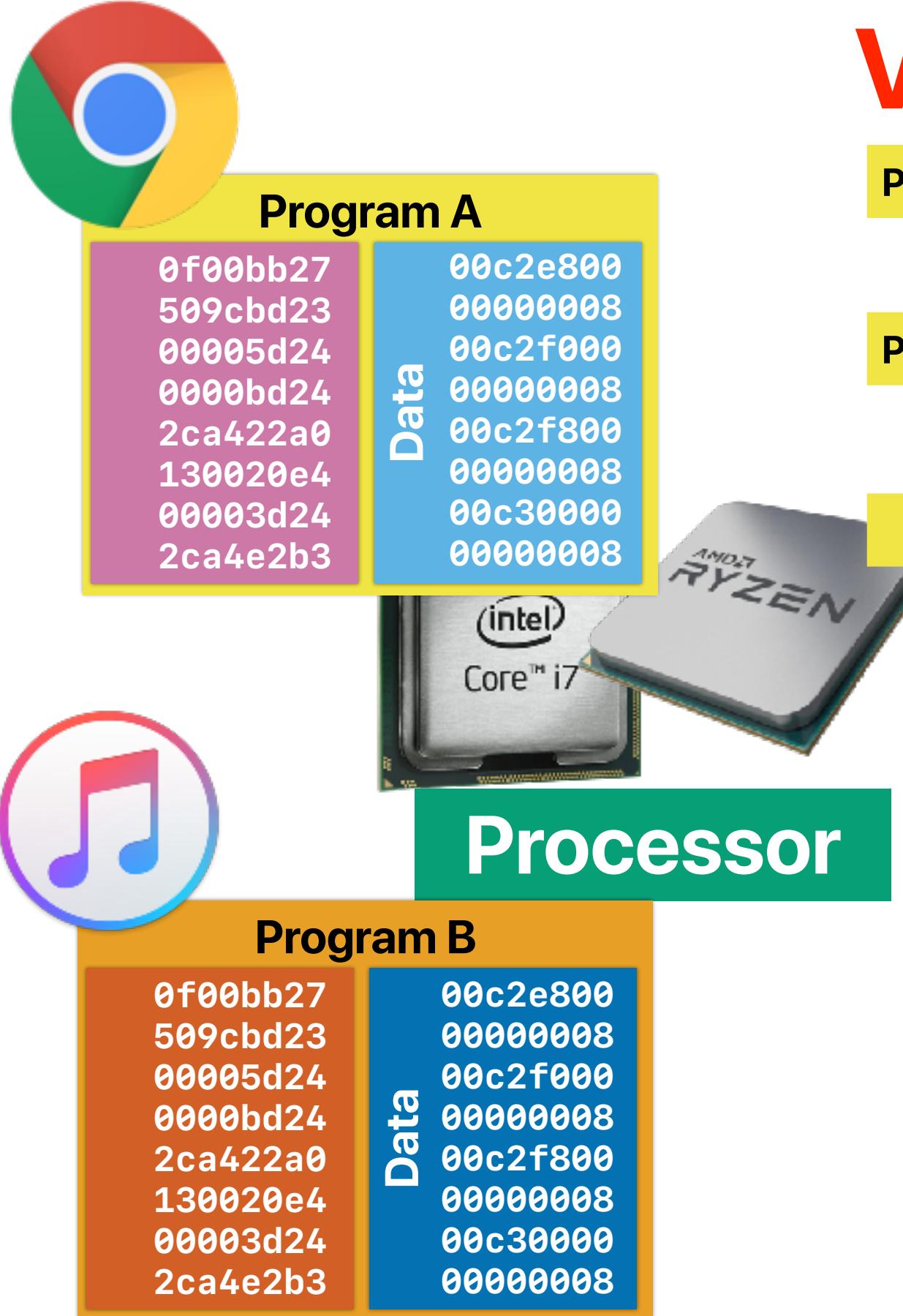
A. (1) & (3)
B. (1) & (4)
C. (1) & (5)
D. (2) & (3)
E. (2) & (4)

```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

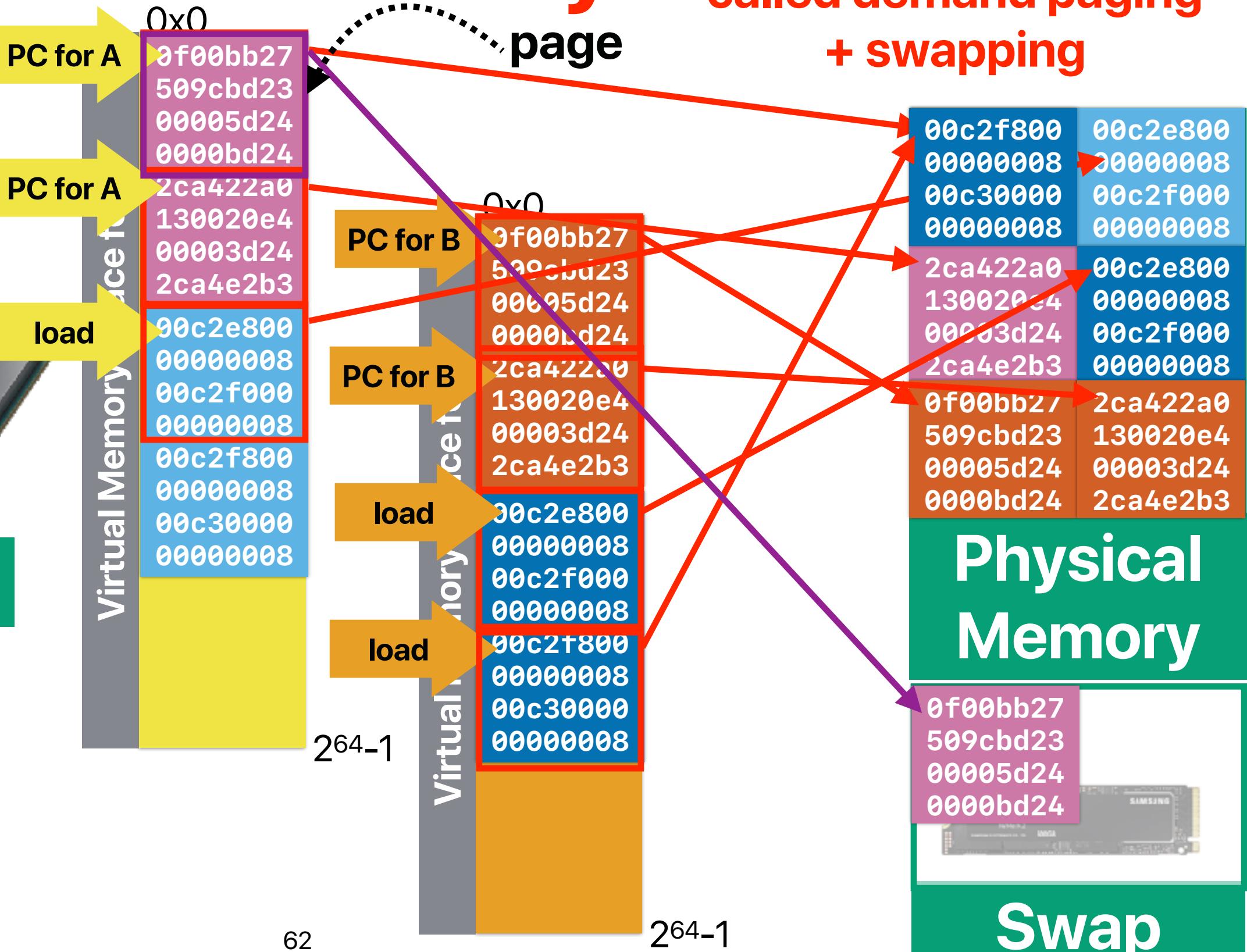
int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    srand((int)time(NULL)+(int)getpid());
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address
of a is %p\n",getpid(), a, &a);
    return 0;
}
```

Virtual Memory



Virtual memory

**This approach is
called demand paging
+ swapping**



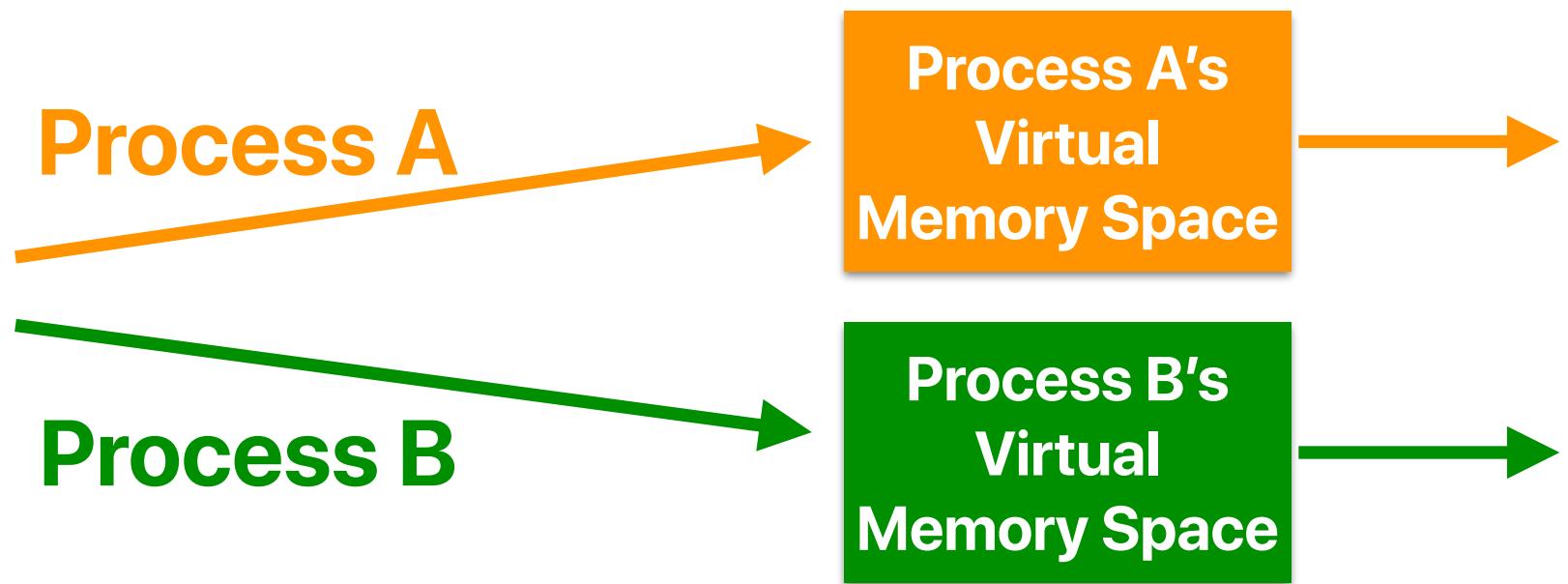
Demo revisited

```
#define _GNU_SOURCE
#include <unistd.h>
#include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <sched.h>
#include <sys/syscall.h>
#include <time.h>

double a;

int main(int argc, char *argv[])
{
    int i, number_of_total_processes=4;
    number_of_total_processes = atoi(argv[1]);
    for(i = 0; i< number_of_total_processes-1 && fork(); i++);
    srand((int)time(NULL)+(int)getpid());
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
    sleep(10);
    fprintf(stderr, "\nProcess %d. Value of a is %lf and address of a is %p\n",getpid(), a, &a);
    return 0;
}
```

&a = 0x601090



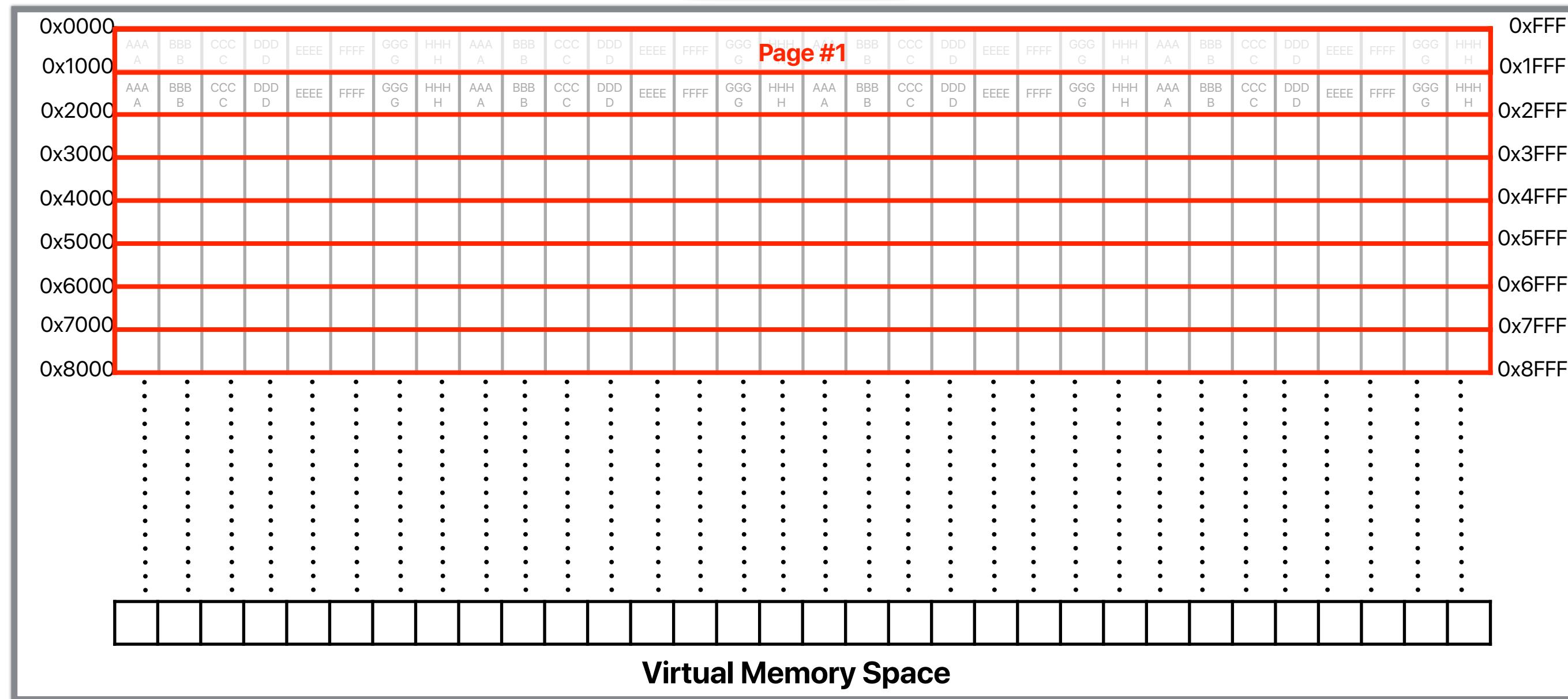
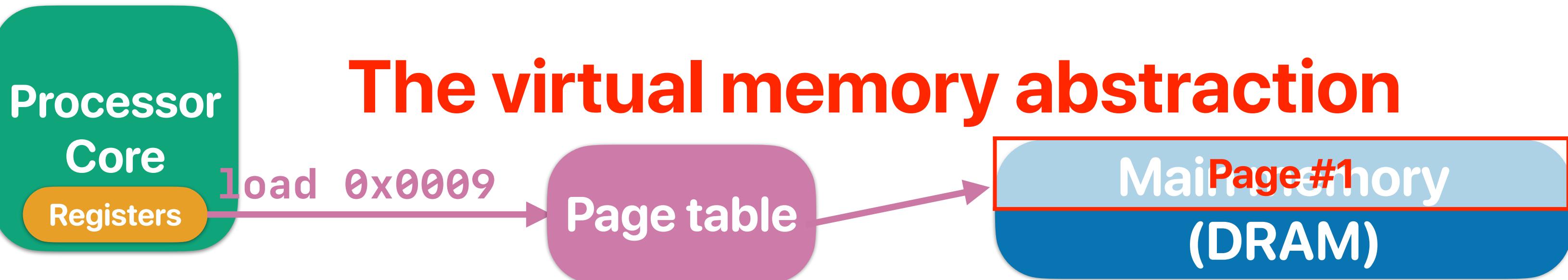
Virtual memory

- An **abstraction** of memory space available for programs/software/programmer
- Programs execute using virtual memory address
- The operating system and hardware work together to handle the mapping between virtual memory addresses and real/physical memory addresses
- Virtual memory organizes memory locations into “**pages**”

Why Virtual memory?

- Allowing multiple applications to share physical main memory
 - Memory protection/isolation among programs/processes is automatically achieved
- Allowing applications to work even though the installed physical memory or available physical memory is smaller than the working set of the application
 - Programmer does not need to worry about the physical memory capacity of different machines — make compiled program compatible
 - Multiple programs can work concurrently even though their total memory demand is larger than the installed physical memory

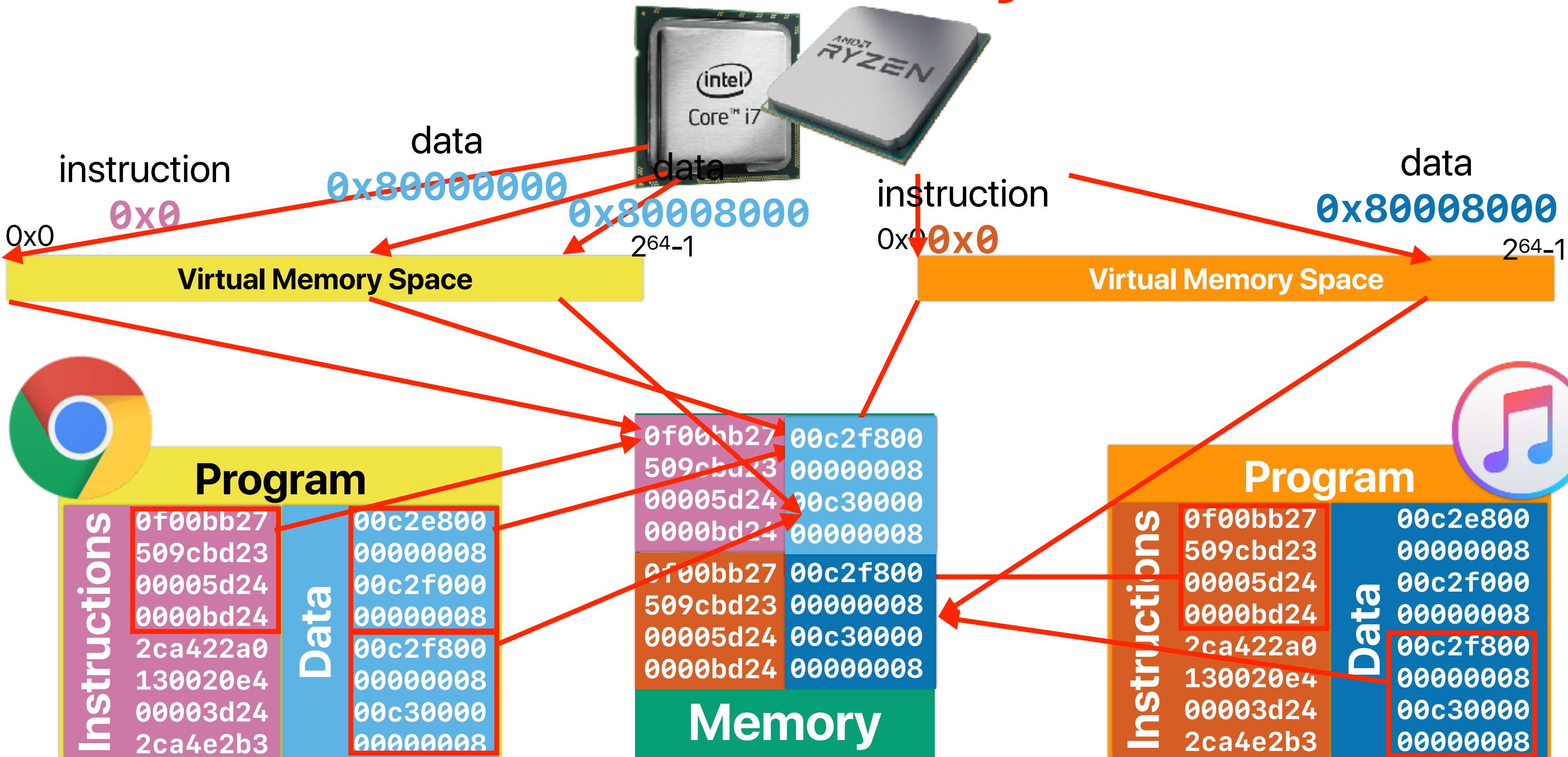
The virtual memory abstraction



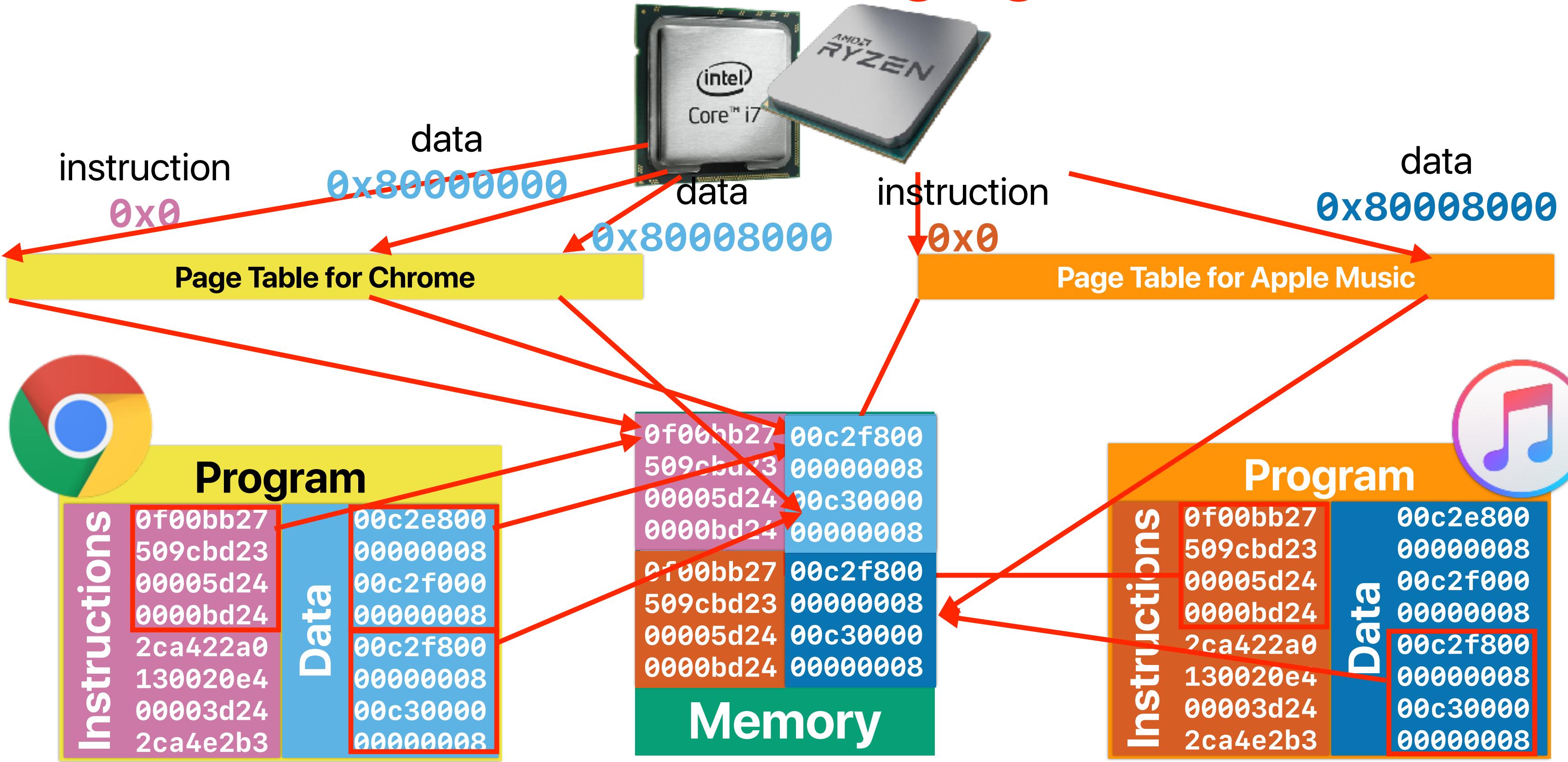
Demand paging

- Treating physical main memory as a “cache” of virtual memory
- The block size is the “page size”
- The page table is the “tag array”
- It’s a “fully-associate” cache — a virtual page can go anywhere in the physical main memory

Virtual memory

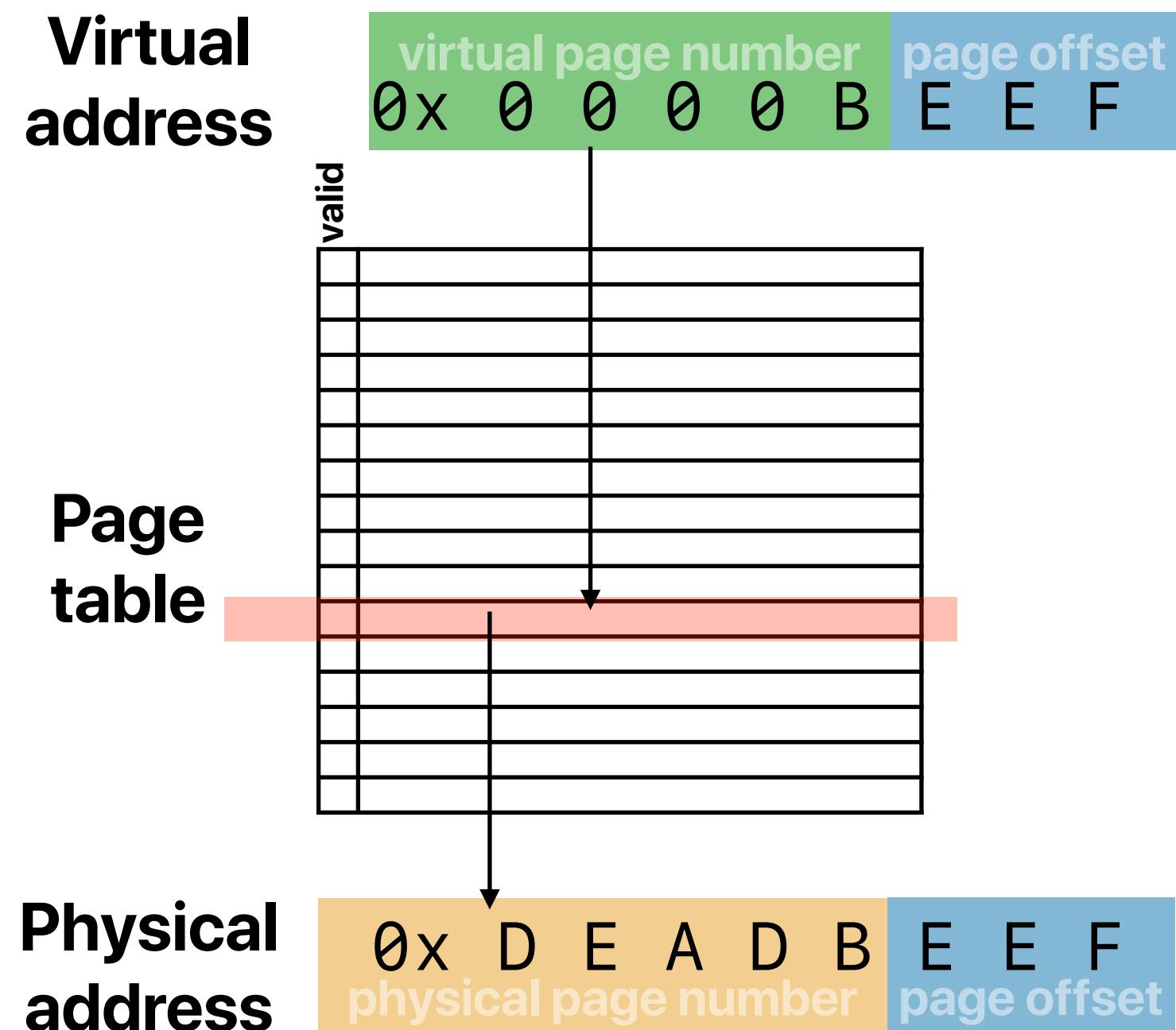


Demand paging



Address translation

- Processor receives virtual addresses from the running code, main memory uses physical memory addresses
- Virtual address space is organized into “pages”
- The system references the **page table** to translate addresses
 - Each process has its own page table
 - The page table content is maintained by OS





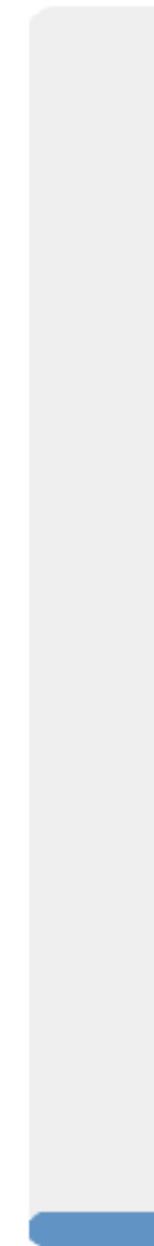
Size of page table

- Assume that we have **64-bit** virtual address space, each page is 4KB, each page table entry is 8 Bytes, what magnitude in size is the page table for a process?
 - A. MB — 2^{20} Bytes
 - B. GB — 2^{30} Bytes
 - C. TB — 2^{40} Bytes
 - D. PB — 2^{50} Bytes
 - E. EB — 2^{60} Bytes

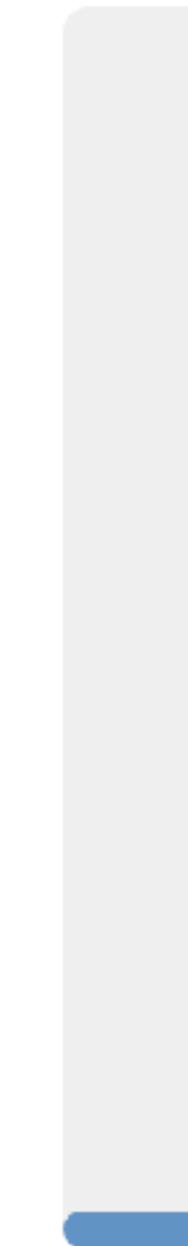


 0

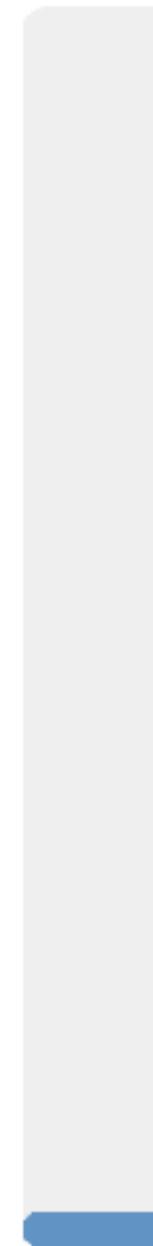
0



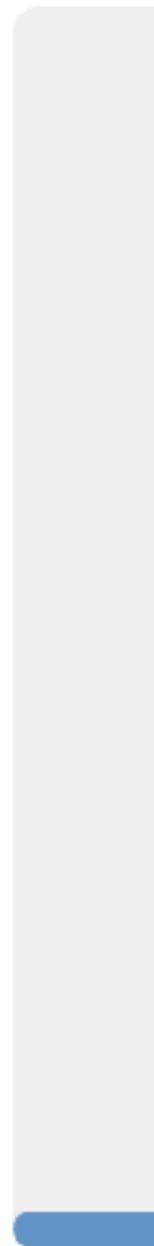
0



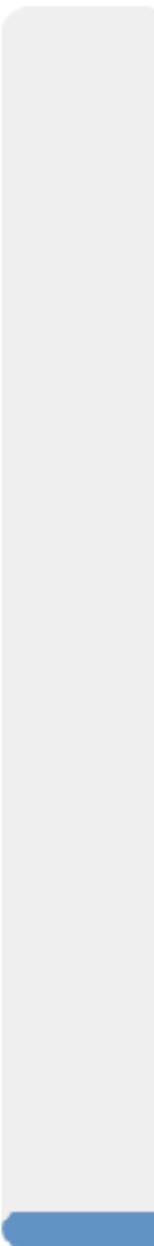
0



0



0



A

B

C

D

E



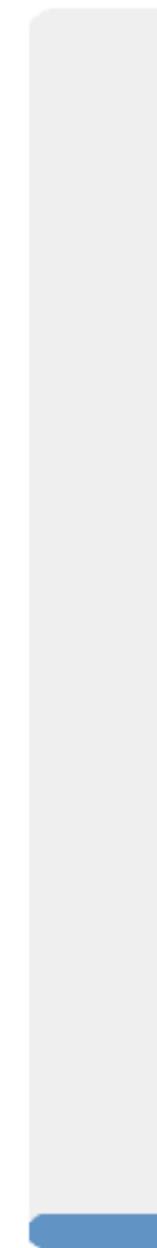
Size of page table

- Assume that we have **64-bit** virtual address space, each page is 4KB, each page table entry is 8 Bytes, what magnitude in size is the page table for a process?
 - A. MB — 2^{20} Bytes
 - B. GB — 2^{30} Bytes
 - C. TB — 2^{40} Bytes
 - D. PB — 2^{50} Bytes
 - E. EB — 2^{60} Bytes



 0

0



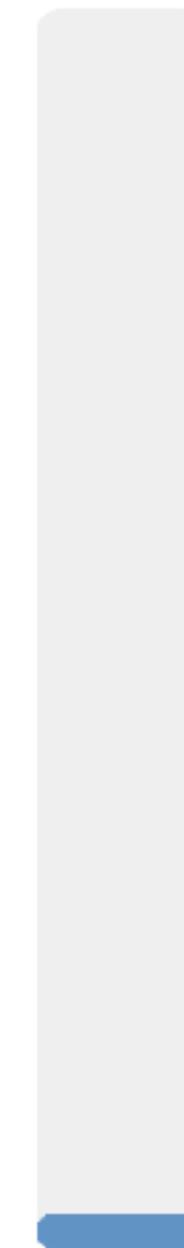
A

0



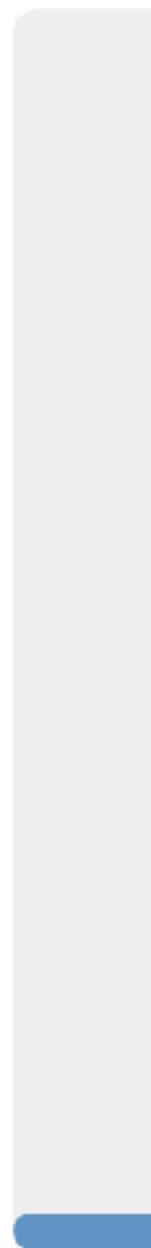
B

0



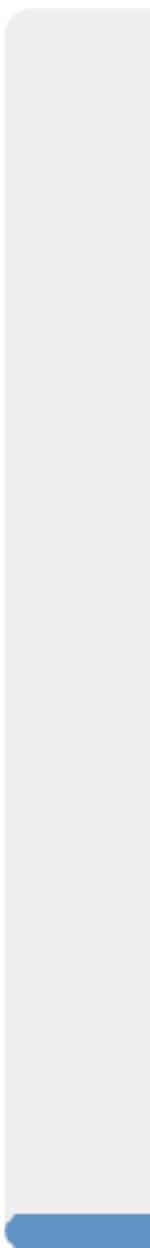
C

0



D

0



E

Size of page table

- Assume that we have **64-bit** virtual address space, each page is 4KB, each page table entry is 8 Bytes, what magnitude in size is the page table for a process?
 - MB — 2^{20} Bytes
 - GB — 2^{30} Bytes
 - TB — 2^{40} Bytes
 - PB — 2^{50} Bytes
 - EB — 2^{60} Bytes

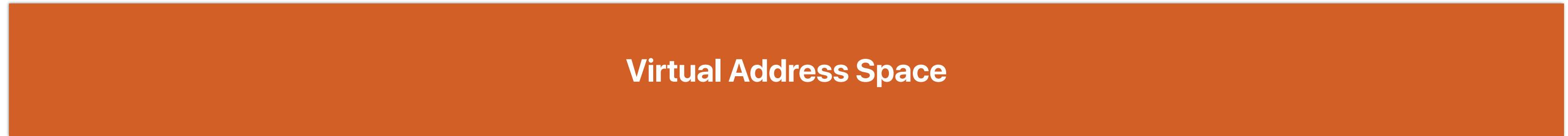
$$\frac{2^{64} \text{ Bytes}}{4 \text{ KB}} \times 8 \text{ Bytes} = 2^{55} \text{ Bytes} = 32 \text{ PB}$$

If you still don't know why — you need to take CS202

Conventional page table

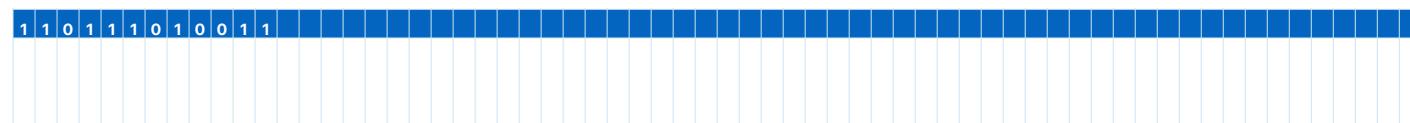
0x0

0xFFFFFFFFFFFFFFFFF

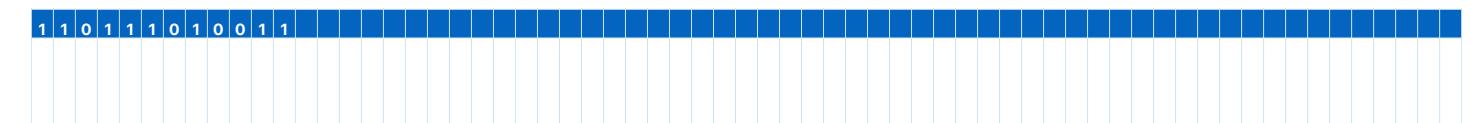


- must be consecutive in the physical memory
- need a big segment! — difficult to find a spot
- simply too big to fit in memory if address space is large!

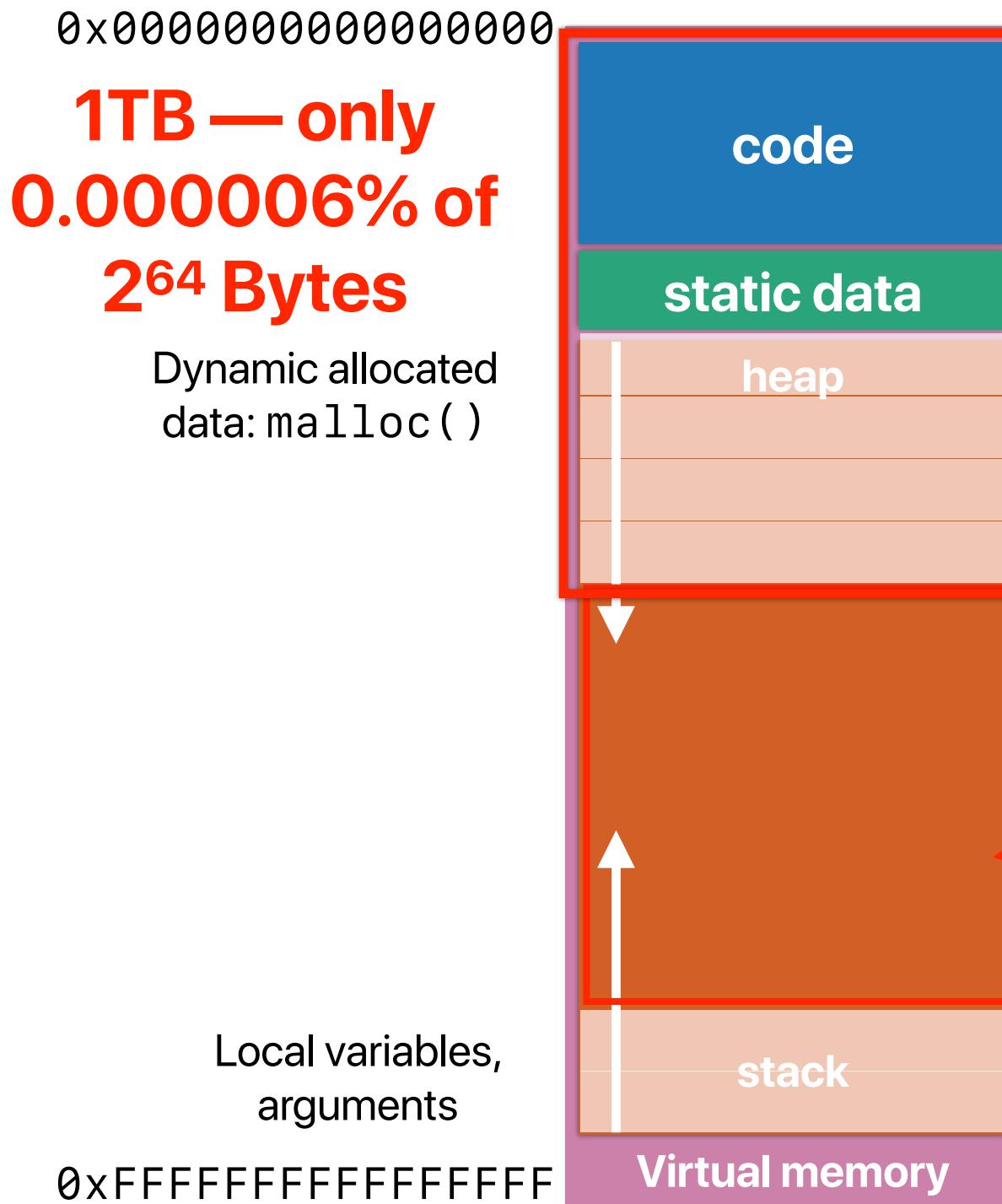
$\frac{2^{64} \text{ } B}{2^{12} \text{ } B}$ page table entries/leaf nodes



.....



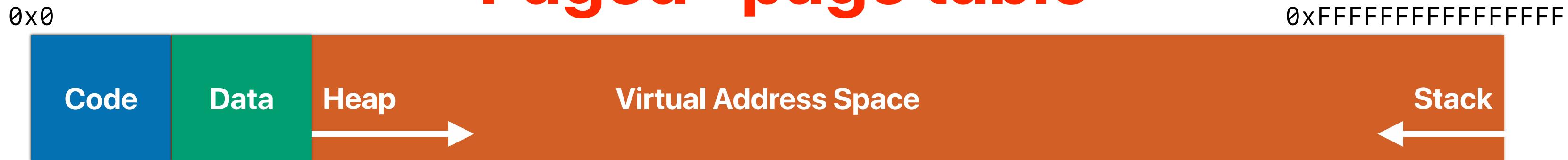
Do we really need a large table?



**Your program probably
never uses this huge area!**

If you still don't know why — you need to take CS202

"Paged" page table

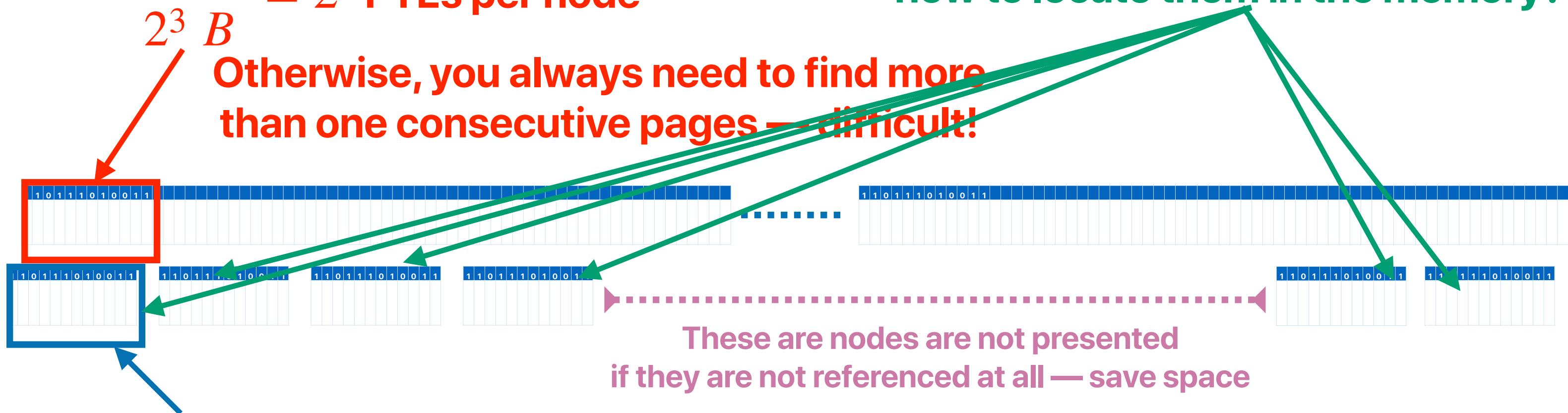


Break up entries into pages!

Each of these occupies exactly a page

$$-\frac{2^{12} \text{ B}}{2^3 \text{ B}} = 2^9 \text{ PTEs per node}$$

Otherwise, you always need to find more than one consecutive pages — difficult!

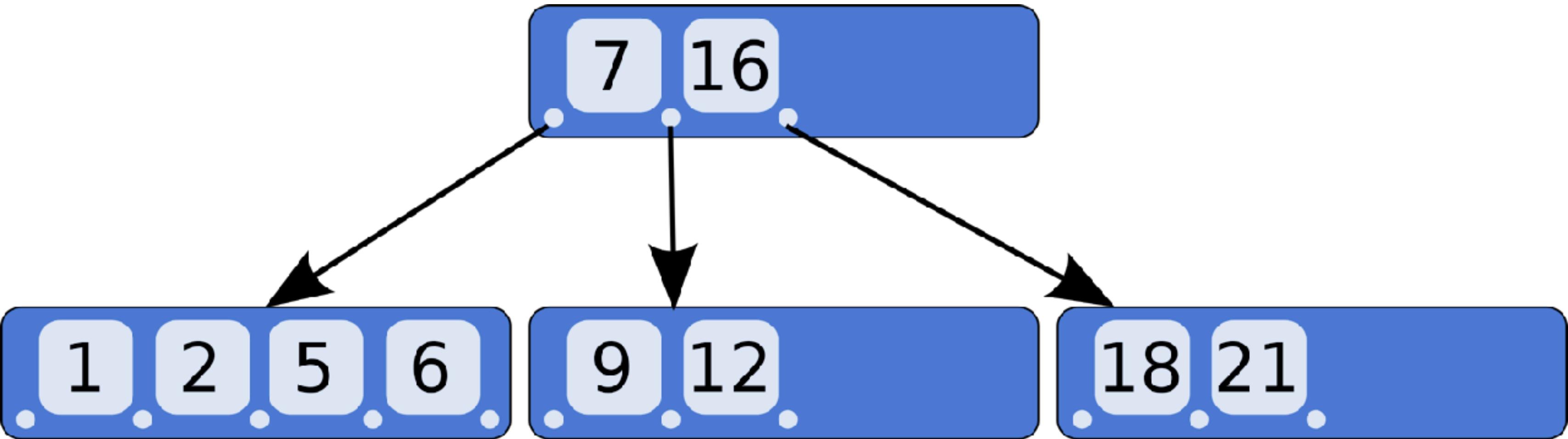


Allocate page table entry nodes “on demand”⁷⁸

Question:

These nodes are spread out, how to locate them in the memory?

B-tree

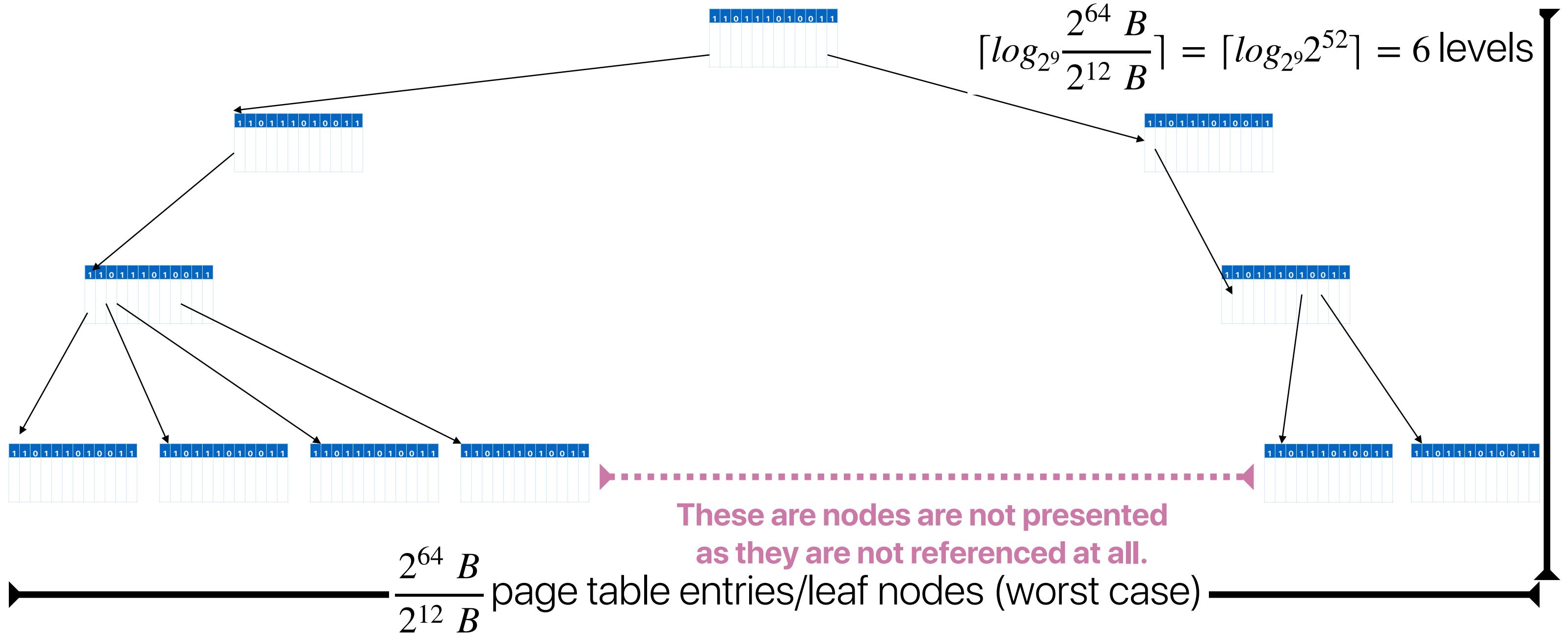


<https://en.wikipedia.org/wiki/B-tree#/media/File:B-tree.svg>

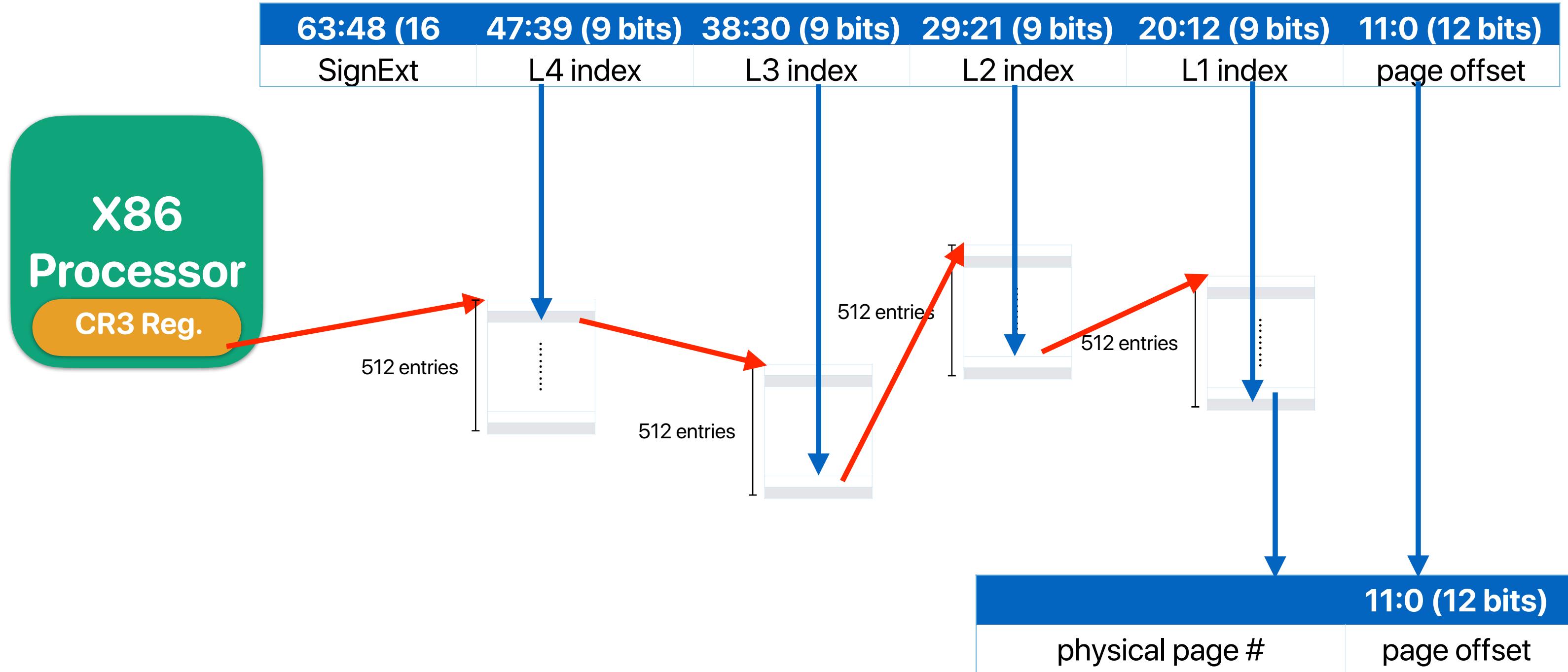
Hierarchical Page Table

0x0

0xFFFFFFFFFFFFFFFFF



Address translation in x86-64





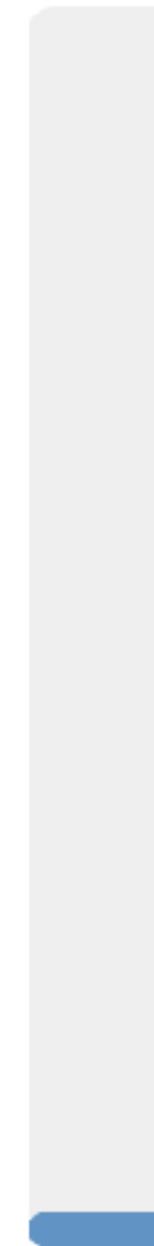
When we have virtual memory...

- If an x86 processor supports virtual memory through the basic format of the page table as shown in the previous slide, how many memory accesses can a **mov** instruction that access data memory once incur?
 - 2
 - 4
 - 6
 - 8
 - 10



 0

0



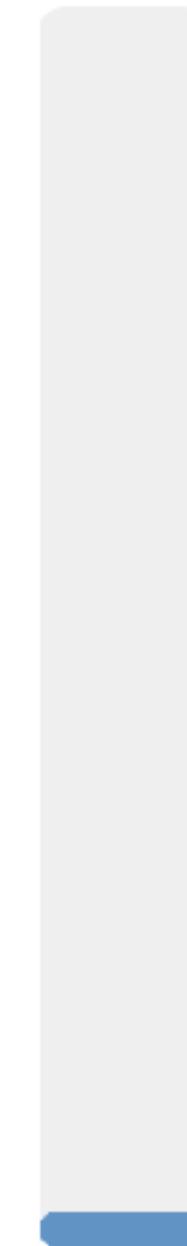
A

0



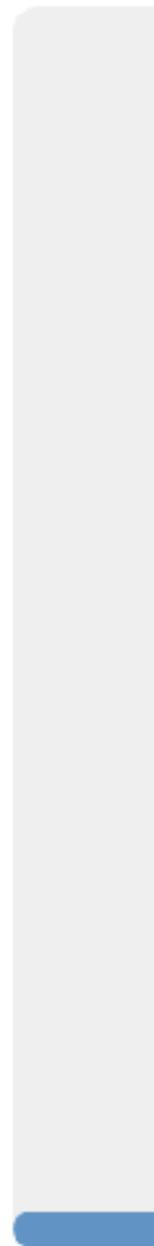
B

0



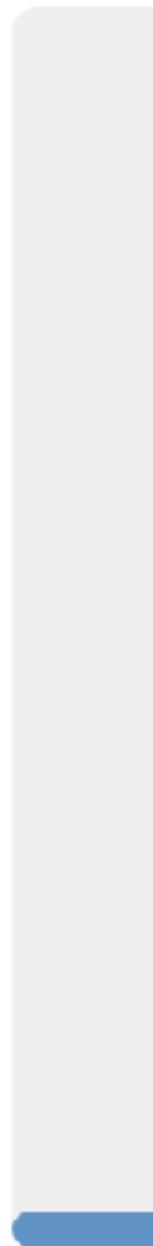
C

0



D

0



E



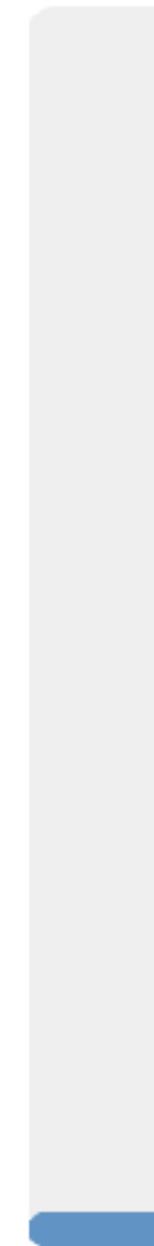
When we have virtual memory...

- If an x86 processor supports virtual memory through the basic format of the page table as shown in the previous slide, how many memory accesses can a **mov** instruction that access data memory once incur?
 - 2
 - 4
 - 6
 - 8
 - 10



 0

0



A

0



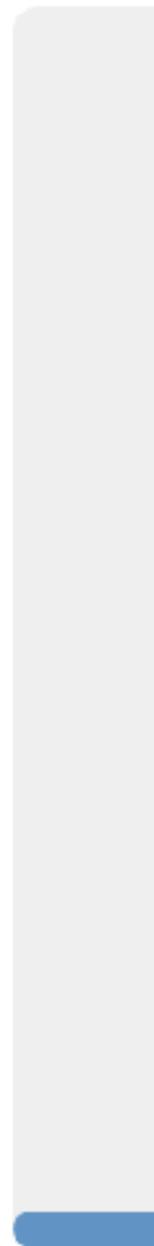
B

0



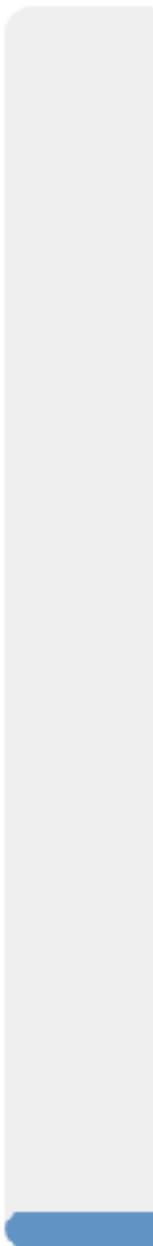
C

0



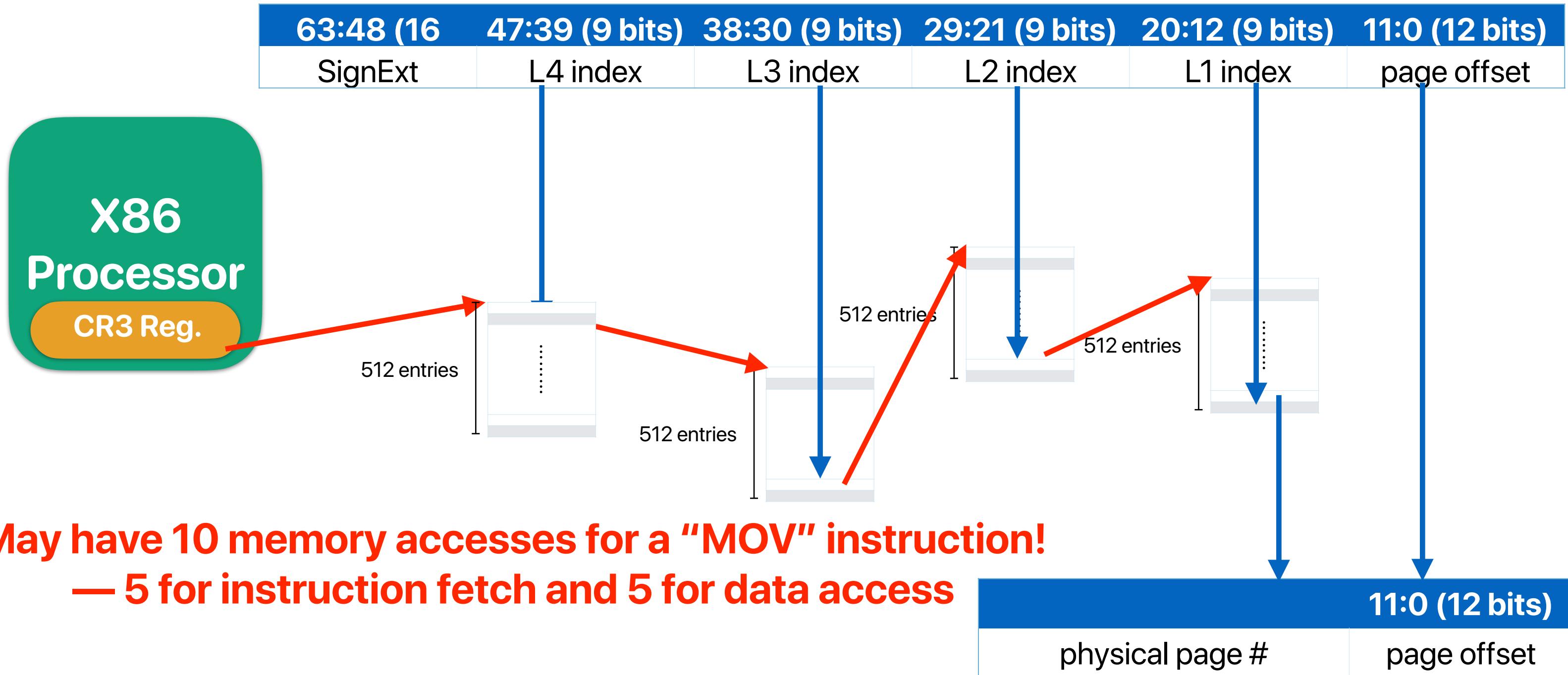
D

0



E

Address translation in x86-64

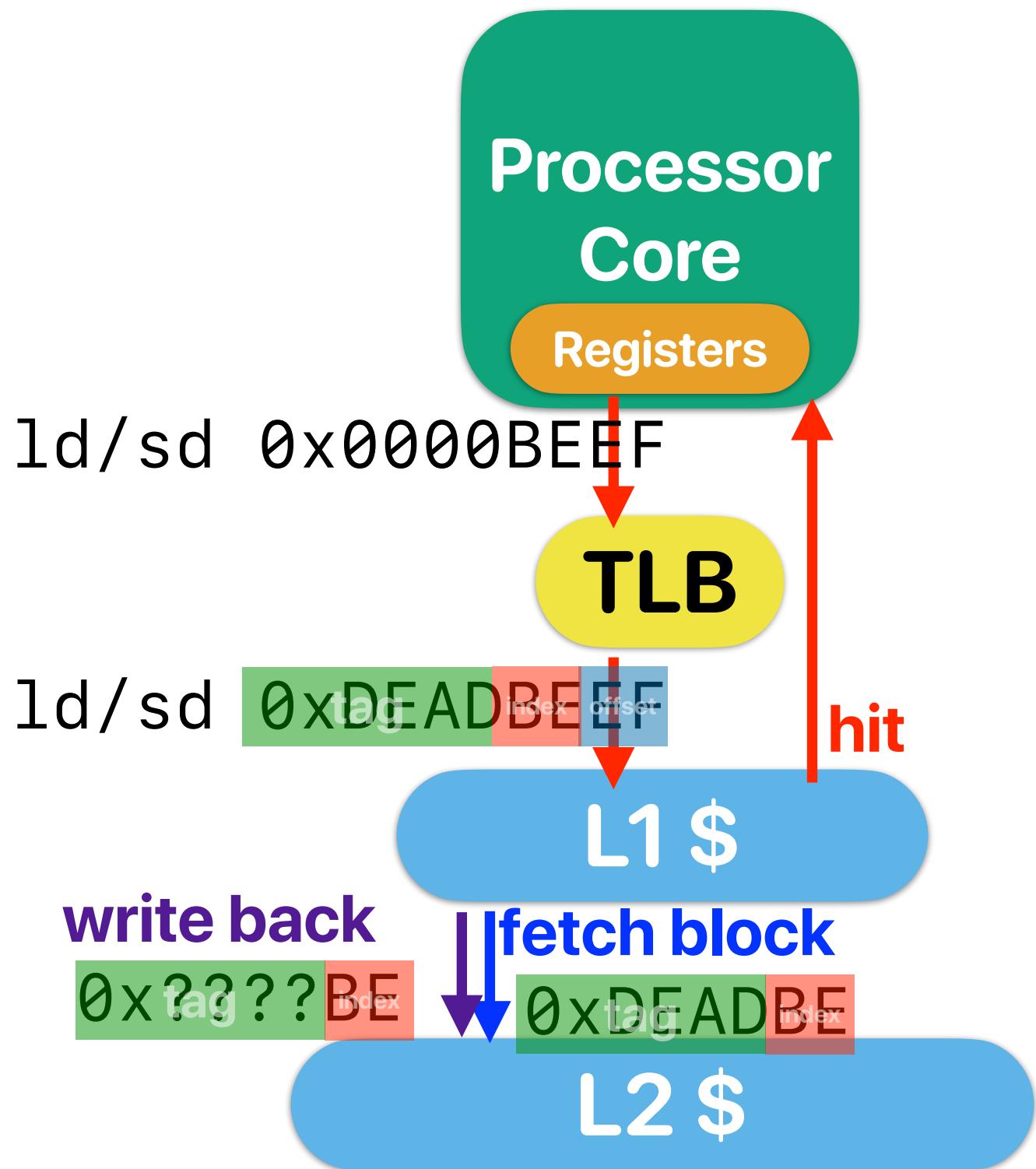


When we have virtual memory...

- If an x86 processor supports virtual memory through the basic format of the page table as shown in the previous slide, how many memory accesses can a **mov** instruction that access data memory incur?
 - 2
 - 4
 - 6
 - 8
 - 10

Avoiding the address translation overhead

TLB: Translation Look-aside Buffer



- TLB — a small SRAM stores frequently used page table entries
- Good — A lot faster than having everything going to the DRAM
- Bad — Still on the critical path

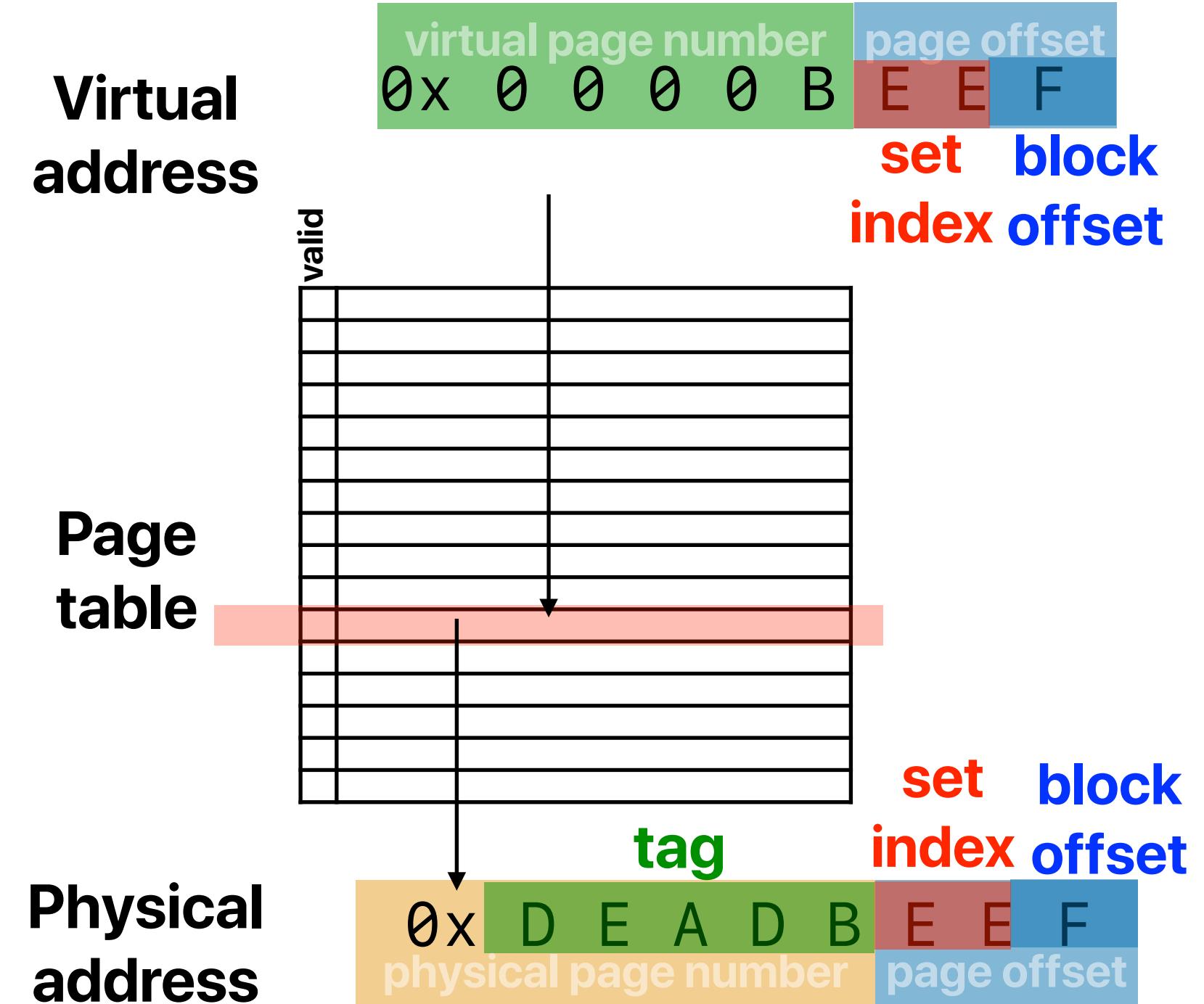
TLB + Virtual cache

- L1 \$ accepts virtual address — you don't need to translate
- Good — you can access both TLB and L1-\$ at the same time and physical address is only needed if L1-\$ misses
- Bad — it doesn't work in practice
 - Many applications have the same virtual address but should be pointing different **physical addresses**
 - An application can have “aliasing virtual addresses” pointing to the same **physical address**



Virtually indexed, physically tagged cache

- Can we find physical address directly in the virtual address
 - Not everything — but the page offset isn't changing!
- Can we indexing the cache using the "partial physical address"?
 - Yes — Just make set index + block set to be exactly the page offset

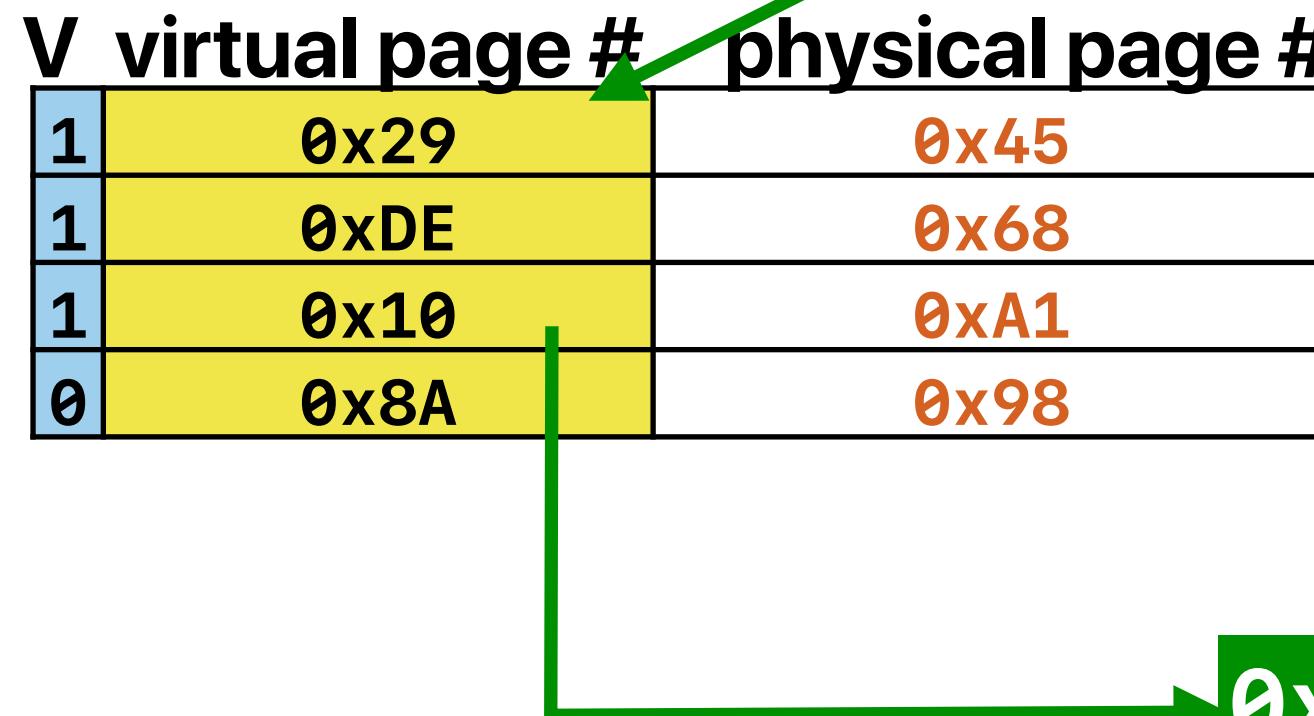


Virtually indexed, physically tagged cache

memory address:

0x0 8 2 4 set block

memory address:



V	D	tag		data
1	1	0x00	AABBCCDDEEGGFFHH	
1	1	0x10	IIJJKKLLMMNNOOPP	
1	0	0xA1	QQRRSSTTUUUVVWWXX	
0	1	0x10	YYZZAABBCCDDEEFF	
1	1	0x31	AABBCCDDEEGGFFHH	
1	1	0x45	IIJJKKLLMMNNOOPP	
0	1	0x41	QQRRSSTTUUUVVWWXX	
0	1	0x68	YYZZAABBCCDDEEFF	

hit'

Virtually indexed, physically tagged cache

- If page size is 4KB —

$$\lg(B) + \lg(S) = \lg(4096) = 12$$

$$C = ABS$$

$$C = A \times 2^{12}$$

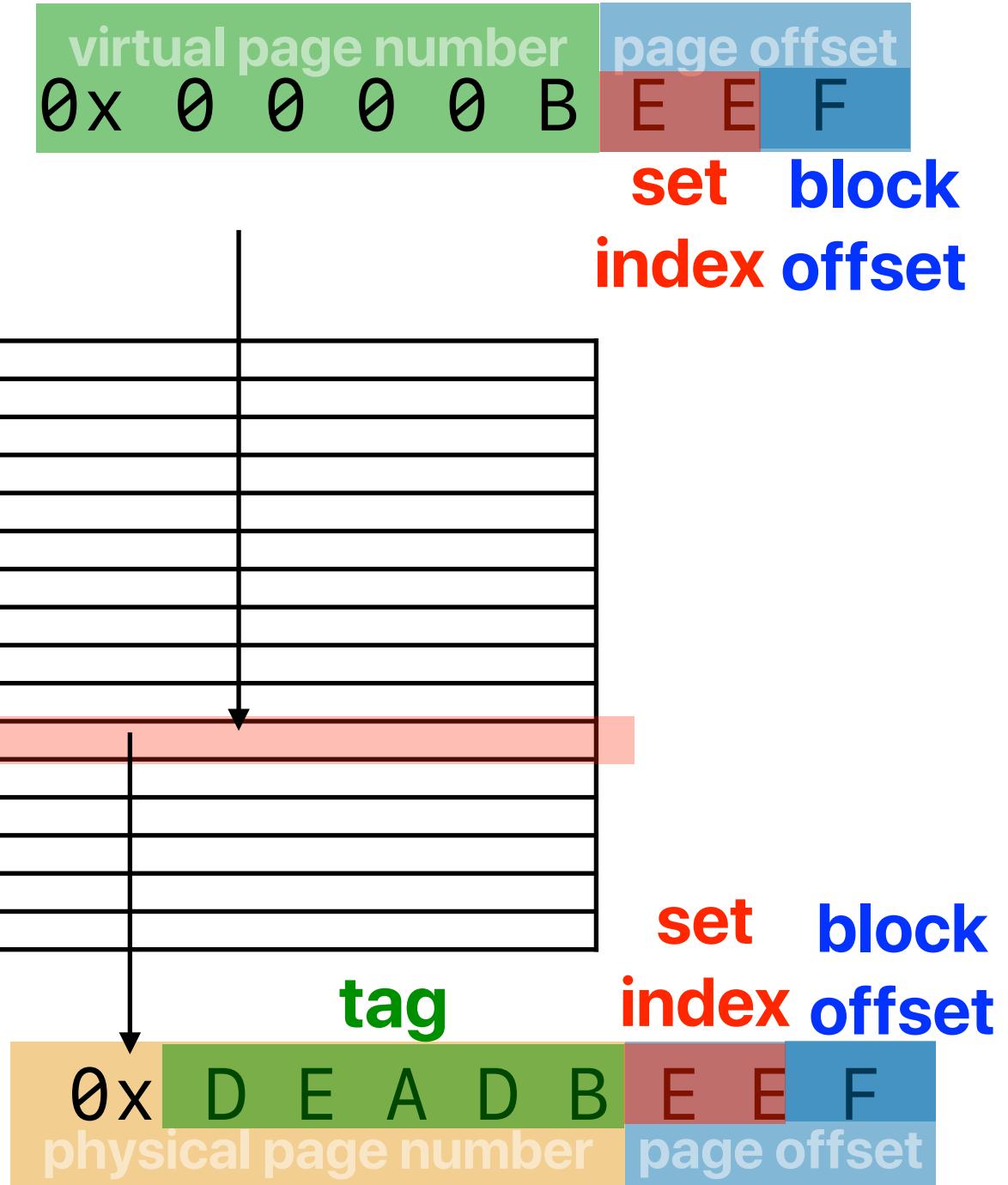
$$\text{if } A = 1$$

$$C = 4KB$$

Virtual address

Page table

Physical address



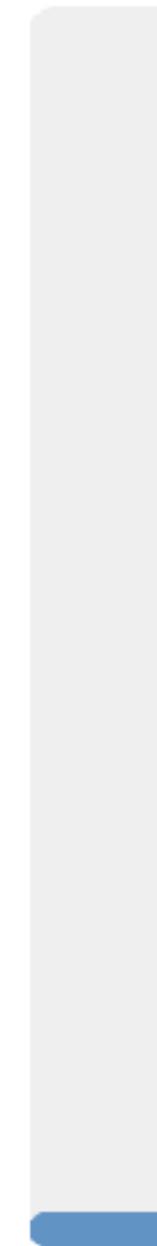


Virtual indexed, physical tagged cache limits the cache size

- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the operating system use 4K pages.
 - A. 32B blocks, 2-way
 - B. 32B blocks, 4-way
 - C. 64B blocks, 4-way
 - D. 64B blocks, 8-way

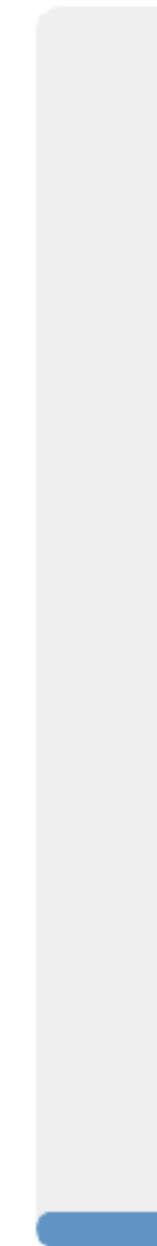
 0

0



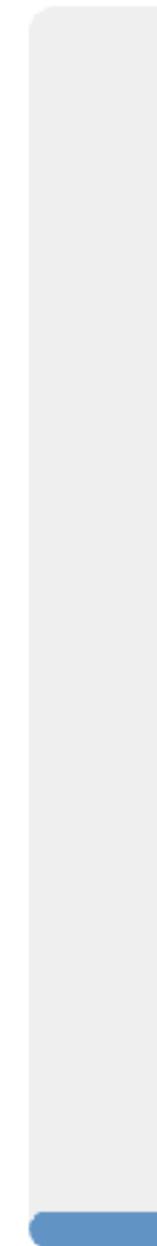
A

0



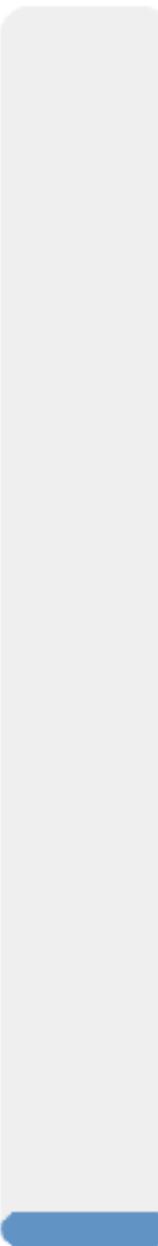
B

0



C

0



D

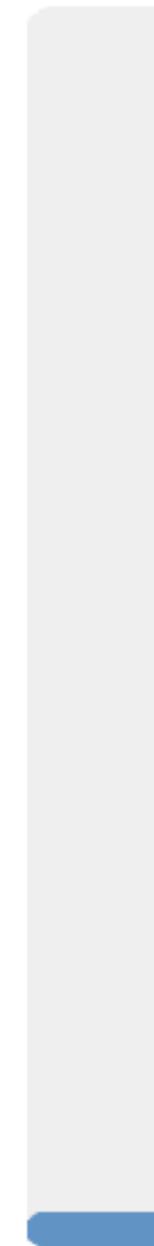


Virtual indexed, physical tagged cache limits the cache size

- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the operating system use 4K pages.
 - A. 32B blocks, 2-way
 - B. 32B blocks, 4-way
 - C. 64B blocks, 4-way
 - D. 64B blocks, 8-way

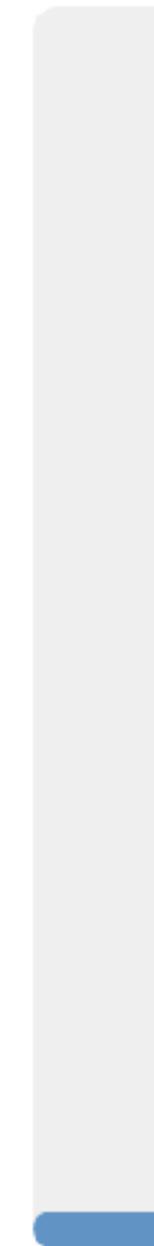
 0

0



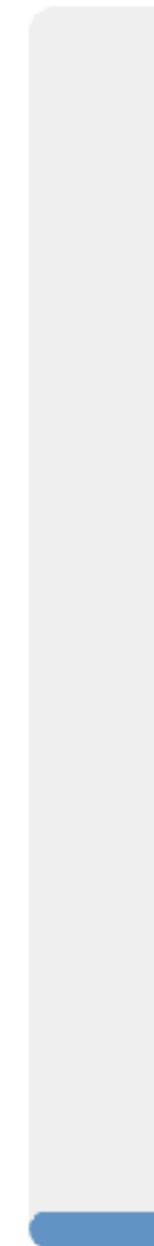
A

0



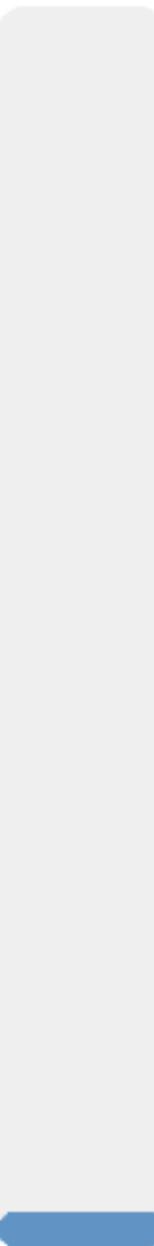
B

0



C

0



D

Virtual indexed, physical tagged cache limits the cache size

- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the operating system use 4K pages.

- A. 32B blocks, 2-way
- B. 32B blocks, 4-way
- C. 64B blocks, 4-way
- D. 64B blocks, 8-way

$$\lg(B) + \lg(S) = \lg(4096) = 12$$

$$C = ABS$$

$$32KB = A \times 2^{12}$$

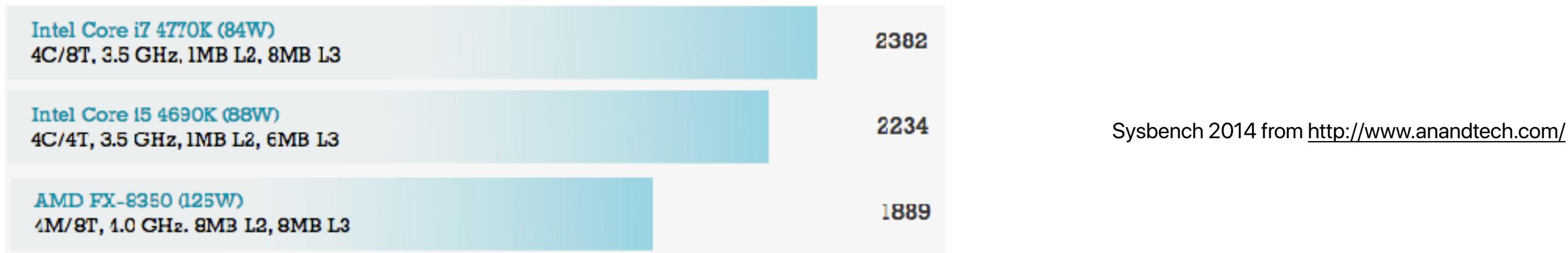
$$A = 8$$

Exactly how Core i7 configures
its own cache

Sample Midterm

Identify the performance bottleneck

- Why does an Intel Core i7 @ 3.5 GHz usually perform better than an Intel Core i5 @ 3.5 GHz or AMD FX-8350@4GHz?



- A. Because the instruction count of the program are different
- B. Because the clock rate of AMD FX is higher
- C. Because the CPI of Core i7 is better
- D. Because the clock rate of AMD FX is higher and CPI of Core i7 is better
- E. None of the above

Amdahl's Law on Multicore Architectures

- Regarding Amdahl's Law on multicore architectures, how many of the following statements is/are correct?
 - ① If we have unlimited parallelism, the performance of each parallel piece does not matter as long as the performance slowdown in each piece is bounded
 - ② With unlimited amount of parallel hardware units, single-core performance does not matter anymore
 - ③ With unlimited amount of parallel hardware units, the maximum speedup will be bounded by the fraction of parallel parts
 - ④ With unlimited amount of parallel hardware units, the effect of scheduling and data exchange overhead is minor

A. 0
B. 1
C. 2
D. 3
E. 4

How programmer affects performance?

- Performance equation consists of the following three factors
 - ① IC
 - ② CPI
 - ③ CT

How many can a **programmer** affect?

- A. 0
- B. 1
- C. 2
- D. 3

Demo — programmer & performance

A

```
for(i = 0; i < ARRAY_SIZE; i++)
{
    for(j = 0; j < ARRAY_SIZE; j++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

B

```
for(j = 0; j < ARRAY_SIZE; j++)
{
    for(i = 0; i < ARRAY_SIZE; i++)
    {
        c[i][j] = a[i][j]+b[i][j];
    }
}
```

- How many of the following make(s) the performance of A better than B?

① IC

② CPI

③ CT

A. 0

B. 1

C. 2

D. 3

Fair comparison

- How many of the following comparisons are fair?
 - ① Comparing the frame rates of Halo 5 on AMD RyZen 1600X and civilization on Intel Core i7 7700X
 - ② Using bit torrent to compare the network throughput on two machines
 - ③ Comparing the frame rates of Halo 5 using medium settings on AMD RyZen 1600X and low settings on Intel Core i7 7700X
 - ④ Using the peak floating point performance to judge the gaming performance of machines using AMD RyZen 1600X and Intel Core i7 7700X
- A. 0
B. 1
C. 2
D. 3
E. 4

Locality

- Which description about locality of arrays sum and A in the following code is the most accurate?

```
for(i = 0; i < 100000; i++)
{
    sum[i%10] += A[i];
}
```

- A. Access of A has temporal locality, sum has spatial locality
- B. Both A and sum have temporal locality, and sum also has spatial locality
- C. Access of A has spatial locality, sum has temporal locality
- D. Both A and sum have spatial locality
- E. Both A and sum have spatial locality, and sum also has temporal locality

intel Core i7

- L1 data (D-L1) cache configuration of Core i7
 - Size 32KB, 8-way set associativity, 64B block
 - Assume 64-bit memory address
 - Which of the following is NOT correct?
 - A. Tag is 52 bits
 - B. Index is 6 bits
 - C. Offset is 6 bits
 - D. The cache has 128 sets

Virtual indexed, physical tagged cache limits the cache size

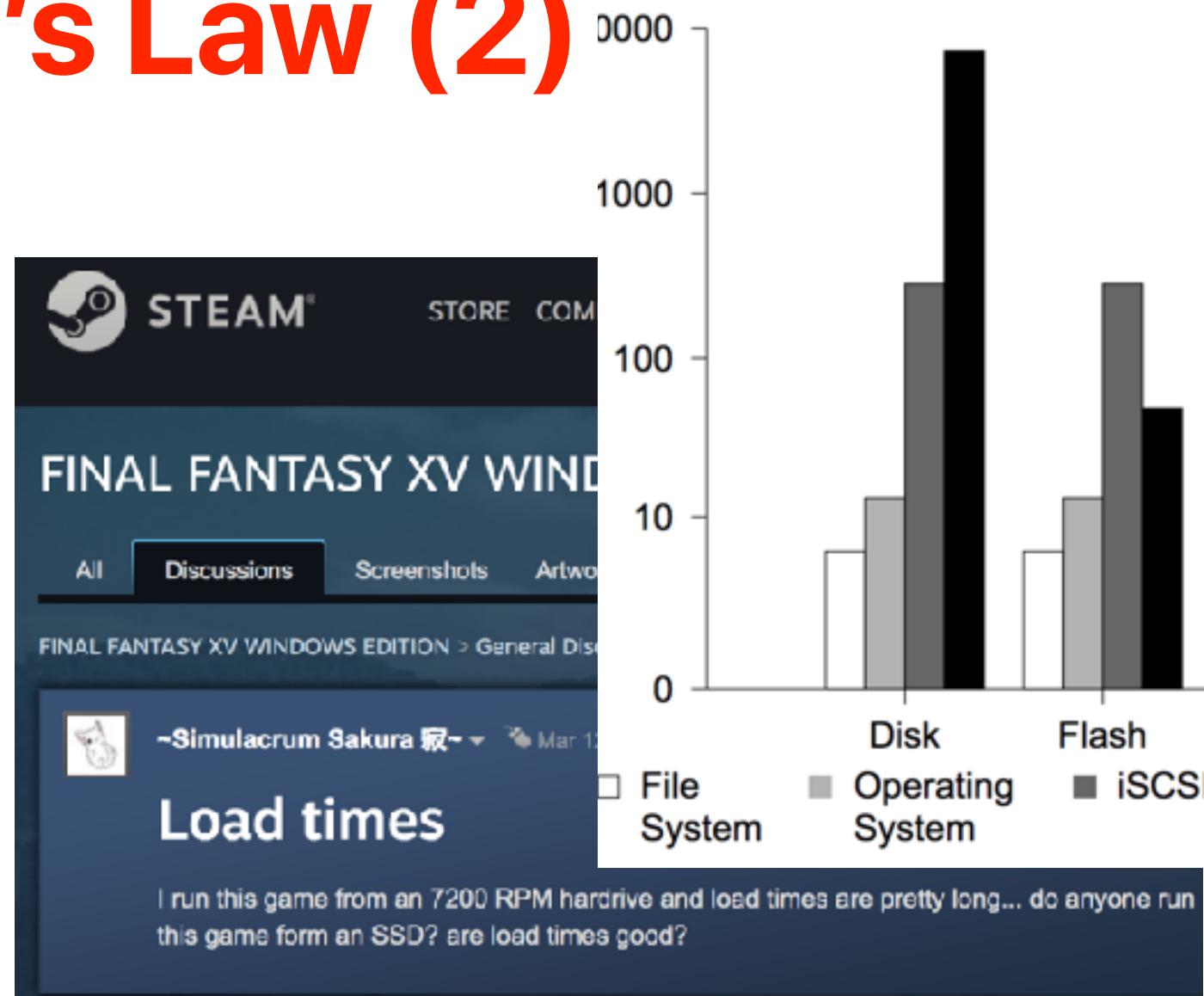
- If you want to build a virtual indexed, physical tagged cache with 32KB capacity, which of the following configuration is possible? Assume the system use 4K pages.
 - A. 32B blocks, 2-way
 - B. 32B blocks, 4-way
 - C. 64B blocks, 4-way
 - D. 64B blocks, 8-way

When we have virtual memory...

- In a modern x86-64 processor supports virtual memory through, how many memory accesses can an instruction incur?
 - A. 2
 - B. 4
 - C. 6
 - D. 8
 - E. More than 10

Practicing Amdahl's Law (2)

- Final Fantasy XV spends lots of time loading a map — within which period that 95% of the time on the accessing the H.D.D., the rest in the operating system, file system and the I/O protocol. If we replace the H.D.D. with a flash drive, which provides 100x faster access time and a better processor to accelerate the software overhead by 2x. By how much can we speed up the map loading process?
 - ~7x
 - ~10x
 - ~17x
 - ~29x
 - ~100x



What kind(s) of misses can matrix transpose remove?

- By transposing a matrix, the performance of matrix multiplication can be further improved. What kind(s) of cache misses does matrix transpose help to remove?

Block

```
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        c[ii][jj] += a[ii][kk]*b[kk][jj];  
        }  
    }  
}
```

- A. Compulsory miss
- B. Capacity miss
- C. Conflict miss
- D. Capacity & conflict miss
- E. Compulsory & conflict miss

Block + Transpose

```
// Transpose matrix b into b_t  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        b_t[i][j] += b[j][i];  
    }  
}  
  
for(i = 0; i < ARRAY_SIZE; i+=(ARRAY_SIZE/n)) {  
    for(j = 0; j < ARRAY_SIZE; j+=(ARRAY_SIZE/n)) {  
        for(k = 0; k < ARRAY_SIZE; k+=(ARRAY_SIZE/n)) {  
            for(ii = i; ii < i+(ARRAY_SIZE/n); ii++)  
                for(jj = j; jj < j+(ARRAY_SIZE/n); jj++)  
                    for(kk = k; kk < k+(ARRAY_SIZE/n); kk++)  
                        // Compute on b_t  
                        c[ii][jj] += a[ii][kk]*b_t[jj][kk];  
        }  
    }  
}
```

Sample short answer questions (< 30 words)

- What are the limitations of compiler optimizations? Can you list two?
- Please define Amdahl's Law and explain each term in it
- Please define the CPU performance equation and explain each term.
- Can you list two things affecting each term in the performance equation?
- What's the difference between latency and throughput? When should you use latency or throughput to judge performance?
- What's "benchmark" suite? Why is it important?
- Why TFLOPS or inferences per second is not a good metrics?

Amdahl's Law for multiple optimizations

- Assume that memory access takes 30% of execution time.
 - Cache can speedup 80% of memory operation by a factor of 4
 - L2 cache can speedup 50% of the remaining 20% by a factor of 2
- What's the total speedup?

Speedup of Y over X

- Consider the same program on the following two machines, X and Y. By how much Y is faster than X?

	Clock Rate	Instructions	Percentage of Type-A	CPI of Type-A	Percentage of Type-B	CPI of Type-B	Percentage of Type-C	CPI of Type-C
Machine X	3 GHz	5000000000	20%	8	20%	4	60%	1
Machine Y	5 GHz	5000000000	20%	13	20%	4	60%	1

Performance evaluation with cache

- Consider the following cache configuration on RISC-V processor:

	I-L1	D-L1	L2	DRAM
size	32K	32K	256K	Big enough
block size	64 Bytes	64 Bytes	64 Bytes	4KB pages
associativity	2-way	2-way	8-way	
access time	1 cycle (no penalty if it's a hit)	1 cycle (no penalty if it's a hit)	10 cycles	100 cycles
local miss rate	2%	10%, 20% dirty	15% (i.e., 15% of L1 misses, also miss in the L2), 30% dirty	
Write policy	N/A	Write-back, write allocate		
Replacement	LRU replacement policy			

The application has 20% branches, 10% loads/stores, 70% integer instructions.

Assume that TLB miss rate is 2% and it requires 100 cycles to handle a TLB miss. Also assume that the branch predictor has a hit rate of 87.5%, what's the CPI of branch, L/S, and integer instructions? What is the average CPI?

Cache simulation

- The processor has a 8KB, 256B blocked, 2-way L1 cache. Consider the following code:

```
for(i=0;i<256;i++) {  
    a[i] = b[i] + c[i];  
    // load a[i] and load b[i], store to c[i]  
    // &a[0] = 0x10000, &b[0] = 0x20000, &c[0] = 0x30000  
}
```

- What's the total miss rate? How many of the misses are compulsory misses? How many of the misses are conflict misses?
- How can you improve the cache performance of the above code through changing hardware?
- How can you improve the performance **without** changing hardware?



Start the presentation to see live content. Still no live content? Install the app or get help at PollEv.com/app

Q & A



Announcement

- Assignment #2 due the upcoming Sunday
 - Assignments SHOULD BE done/submitted individually — if discussed with others, make sure their names on your submission
 - We will drop your least performing assignment as well
 - Check the website tonight for the template and questions
- Midterm next Monday — will release a midterm review online only lecture by Friday
 - It will be held during the lecture time — both in-person and online sessions — 8/21 2p-3:20p
 - Both formats contain multiple choices, short answers, free answers
 - Online examine will have 1.5x questions compared to in-person since typing is faster and online is open-book.

Computer Science & Engineering

142

つづく

