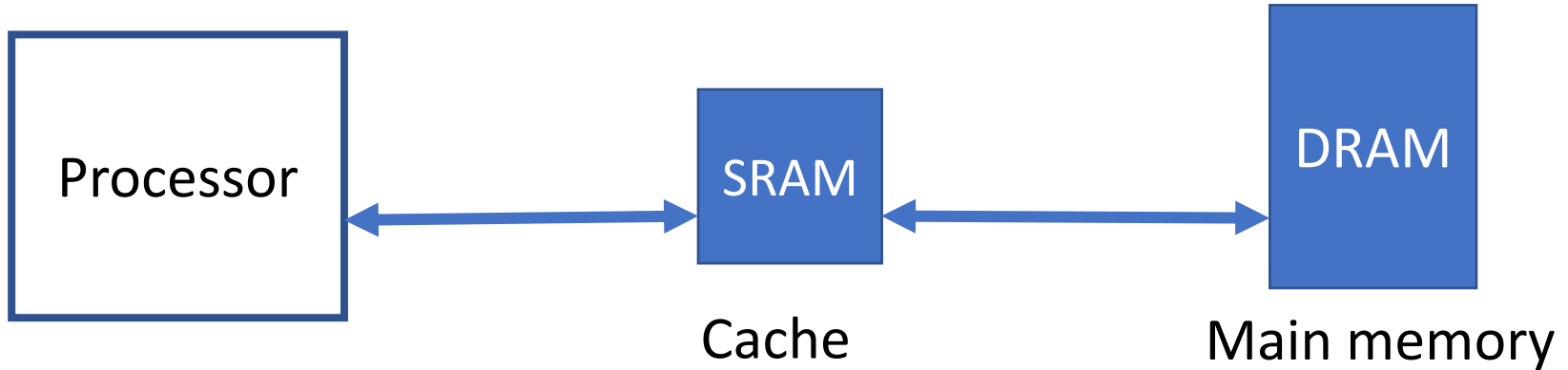


Application of memory management in C: Cache-efficient algorithms Matrix Multiplication

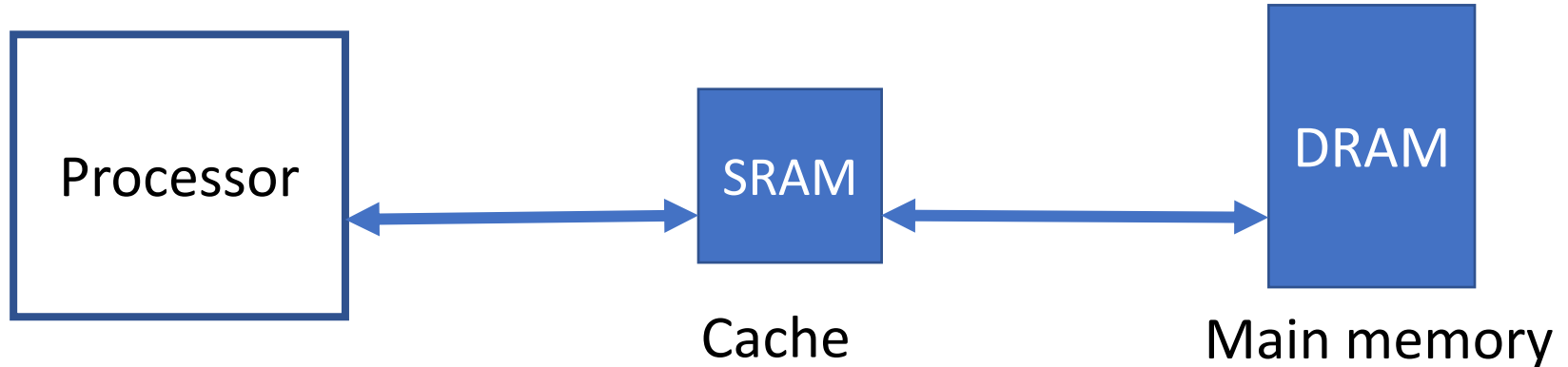
Sujoy Sinha Roy
School of Computer Science
University of Birmingham

Hypothetical computer



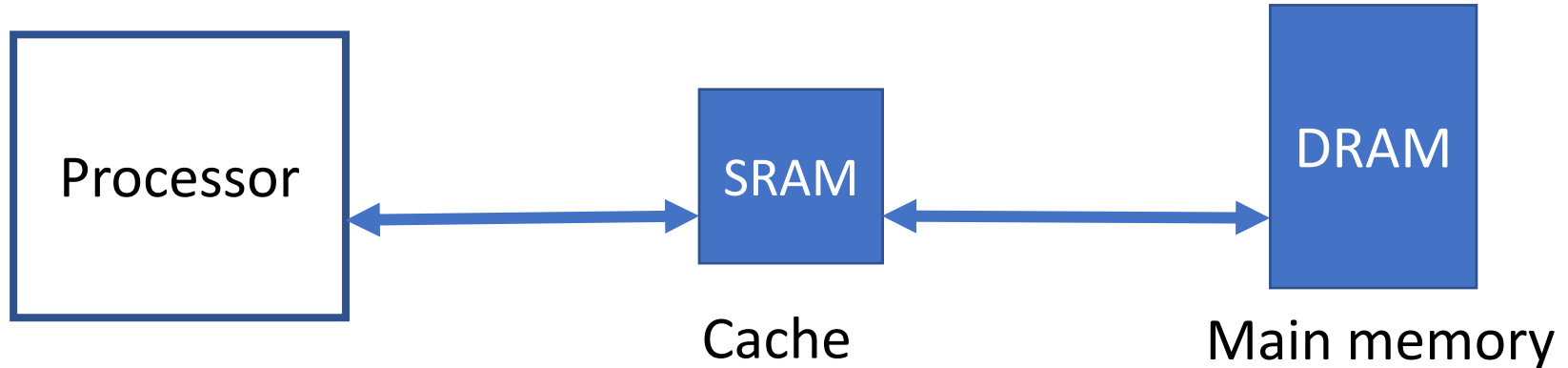
Assumption: processor can access fast cache in negligible time

Simplified model for computation-communication tradeoff



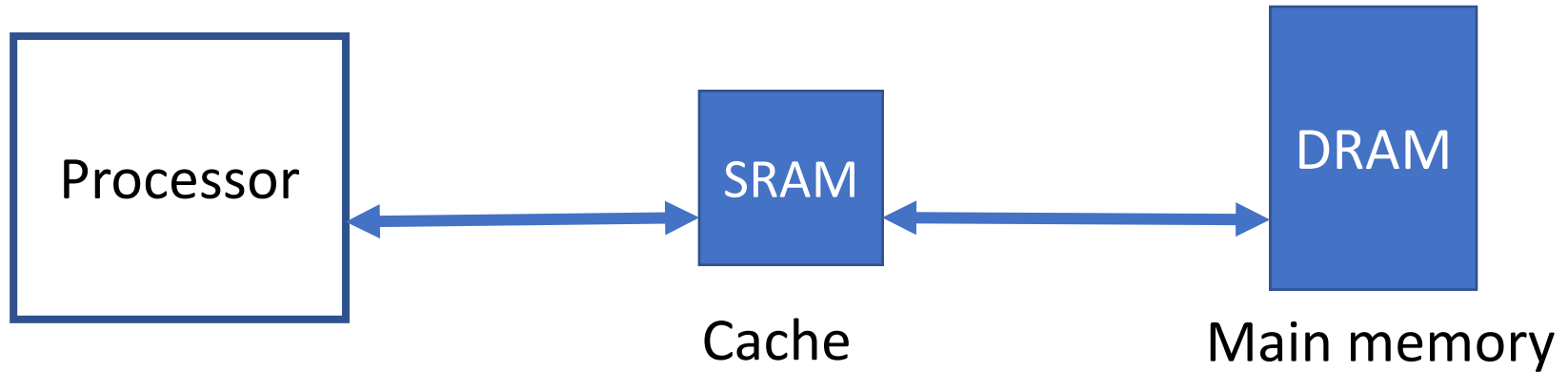
- Assume a generic program which performs
 - m = number of data movements between fast and slow memory
 - f = number of arithmetic operations

Simplified model for computation-communication tradeoff



- Assume a generic program which performs
 - m = number of data movements between fast and slow memory
 - f = number of arithmetic operations
 - t_m = time per data movement
 - t_f = time per arithmetic operation
- Total time for program execution = $f * t_f + m * t_m$

Simplified model for computation-communication tradeoff



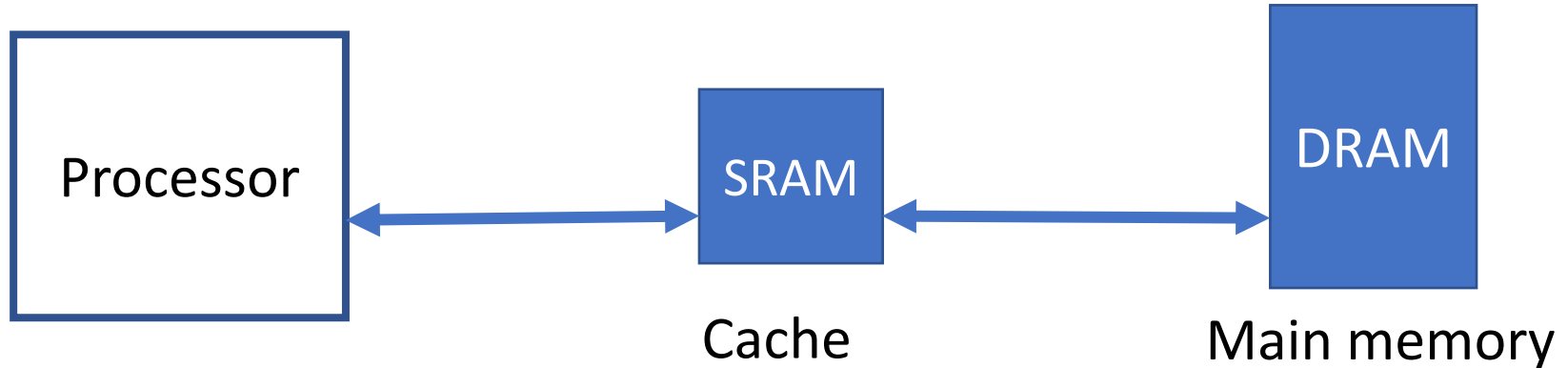
- Assume a generic program which performs
 - m = number of data movements between fast and slow memory
 - f = number of arithmetic operations
 - t_m = time per data movement
 - t_f = time per arithmetic operation
 - $q = f/m$ average number of arithmetic operations per slow memory access

Higher q indicates that the algorithm is ***computation-centric*** and performs less data movement.

So, a higher q results in a better system-performance.



Simplified model for computation-communication tradeoff



- Assume a generic program which performs
 - m = number of data movements between fast and slow memory
 - f = number of arithmetic operations
 - t_m = time per data movement
 - t_f = time per arithmetic operation
 - $q = f/m$ average number of arithmetic operations per slow memory access
- Total time for program execution = $f * t_f + m * t_m$
 $= f * t_f * (1 + t_m/t_f * \mathbf{1/q})$

$$\begin{aligned}\text{Total time for program execution} &= f * t_f + m * t_m \\ &= f * t_f * (1 + t_m/t_f * \mathbf{1/q})\end{aligned}$$

where

- t_m = time per data movement (speed of memory)
- t_f = time per arithmetic operation (speed of processor)

Hardware constants
(programmer cannot control)

Goal: Cache-efficient algorithms increase \mathbf{q} by reducing the number of slow main-memory access \rightarrow improves system performance.

Case study: Arithmetic on dynamically allocated matrices

Assume we have three matrices:

$$A = \begin{vmatrix} A00 & A01 & A02 & A03 \\ A10 & A11 & A12 & A13 \\ A20 & A21 & A22 & A23 \\ A30 & A31 & A32 & A33 \end{vmatrix}$$

$$B = \begin{vmatrix} B00 & B01 & B02 & B03 \\ B10 & B11 & B12 & B13 \\ B20 & B21 & B22 & B23 \\ B30 & B31 & B32 & B33 \end{vmatrix}$$

$$C = \begin{vmatrix} C00 & C01 & C02 & C03 \\ C10 & C11 & C12 & C13 \\ C20 & C21 & C22 & C23 \\ C30 & C31 & C32 & C33 \end{vmatrix}$$

We want to compute $[C] = [C] + [A] * [B]$

Creating a dynamically-allocated matrix

We can create a matrix with n rows and n columns as

```
T *A;  
  
A = (T *) malloc(n*n*sizeof(T));
```

If we want to store data in row-major order then

```
for (int i = 0; i < n; ++i) {  
    for (int j = 0; j < n; ++j)  
        scanf("%f", p+i*n+j);  
}
```

If we want to store data in column-major order then

```
for (int i = 0; i < n; ++i) {  
    for (int j = 0; j < n; ++j)  
        scanf("%f", p+i+j*n);  
}
```

Case study: Arithmetic on dynamically allocated matrices

Assume that matrices are stored in column-major order.

$$A = \begin{vmatrix} A00 & A01 & A02 & A03 \\ A10 & A11 & A12 & A13 \\ A20 & A21 & A22 & A23 \\ A30 & A31 & A32 & A33 \end{vmatrix} \quad B = \begin{vmatrix} B00 & B01 & B02 & B03 \\ B10 & B11 & B12 & B13 \\ B20 & B21 & B22 & B23 \\ B30 & B31 & B32 & B33 \end{vmatrix} \quad \text{Logical view}$$

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33
-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

Memory layout

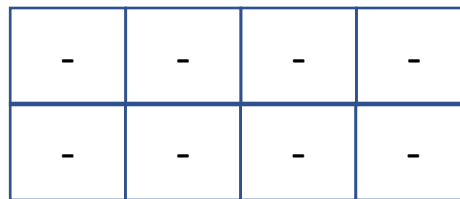
Note: All data is initially in slow Main memory

Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view Memory layout
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

A and B are stored
in slow memory in
column-major order



Cache (fast but small)

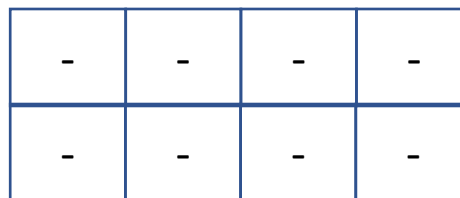
Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view Memory layout
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00 * B00 + A01 * B10 + A02 * B20 + A03 * B30$

A and B are stored
in slow memory in
column-major order



Cache (fast but small)

Compute $C = C + A * B$

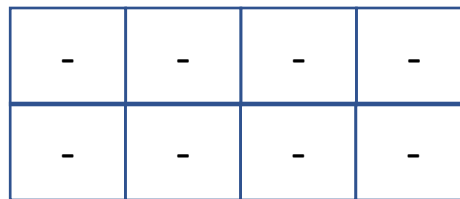
A00	A01	A02	A03	B00	B01	B02	B03	Logical view Memory layout
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00 * B00 + A01 * B10 + A02 * B20 + A03 * B30$

1. Read A00 into fast memory
2. Read B00 into fast memory

A and B are stored
in slow memory in
column-major order



Cache (fast but small)

Compute $C = C + A * B$

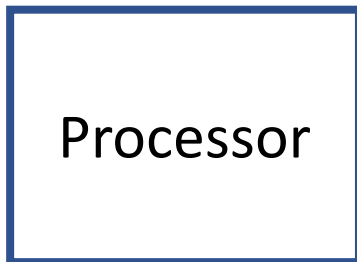
A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
Memory layout								

Spatial locality

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00 * B00 + A01 * B10 + A02 * B20 + A03 * B30$

1. Read A00 into fast memory
2. Read B00 into fast memory



B00	B10	B20	B30
A00	A10	A20	A30

Cache (fast but small)

A and B are stored in slow memory in column-major order

Due to spatial locality, not just A00 and B00, but also nearby elements get stored in the cache.

Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00 * B00 + A01 * B10 + A02 * B20 + A03 * B30$

A and B are stored
in slow memory in
column-major order



B00	B10	B20	B30
A00	A10	A20	A30

Cache (fast but small)

Processor can read only
{A00, B00} from fast mem.
So, $C00 = A00 * B00$

Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = (A00 * B00) + A01 * B10 + A02 * B20 + A03 * B30$

A and B are stored
in slow memory in
column-major order



B00	B10	B20	B30
A00	A10	A20	A30

Cache (fast but small)

Processor can get only B10
from Cache.
A01 is needed from Main
memory.

Compute $C = C + A * B$

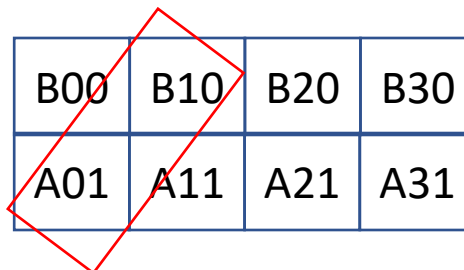
A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = (A00 * B00) + A01 * B10 + A02 * B20 + A03 * B30$

3. Read A01 into fast memory (cache miss)

A and B are stored
in slow memory in
column-major order



Cache (fast but small)

Processor can read only
{A01, B10} from fast mem
 $C00 = C00 + A01 * B10$

Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = (A00 * B00 + A01 * B10) + A02 * B20 + A03 * B30$

4. Read A02 into fast memory (again cache miss!)

A and B are stored
in slow memory in
column-major order



B00	B10	B20	B30
A02	A12	A22	A32

Cache (fast but small)

Processor can read only
{A02, B10} from fast mem
 $C00 = C00 + A02 * B20$

Compute $C = C + A * B$

A00	A01	A02	A03	B00	B01	B02	B03	Logical view Memory layout
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00 * B00 + A01 * B10 + A02 * B20 + A03 * B30$
 5. ... and so on

A and B are stored
 in slow memory in
 column-major order



B00	B10	B20	B30
A03	A13	A23	A33

Cache (fast but small)

Always cache miss for $A[]$
Observation:
 Processor cannot exploit
 spatial locality

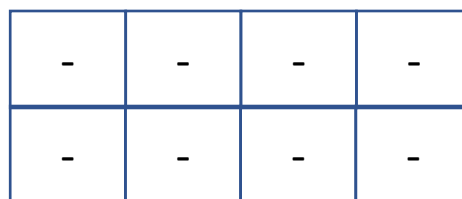
What if A and B are in row-major order?

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	Memory layout

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B01	B02	B03	B10	B11	B12	B13	B20	B21	B22	B30	B30	B31	B32	B33

Compute $C00 = A00*B00 + A01*B10 + A02*B20 + A03*B30$

A and B are stored
in slow memory in
row-major order



Cache (fast but small)

Can processor exploit
spatial locality?

What if A and B are in row-major order? (Answer)

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B01	B02	B03	B10	B11	B12	B13	B20	B21	B22	B30	B30	B31	B32	B33

Compute $C00 = A00*B00 + A01*B10 + A02*B20 + A03*B30$

A and B are stored
in slow memory in
row-major order



Cache (fast but small)

Can processor exploit
spatial locality?
Always cache miss for B[]

Column-major order: #comp, #comm

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	
								Memory layout

A00	A10	A20	A30	A01	A11	A21	A31	A02	A12	A22	A32	A03	A13	A23	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

To compute $C00 = A00*B00+A01*B10+A02*B20+A03*B30$

$m_{C00} =$ read B(, 0) once: #n memory access
 + read entire A(,): # n^2 memory access (due to cache misses)
 + read/write C00 once: #2 memory access
 = (n^2+n+2)

$f_{C00} = n$ multiplications + n additions = $2n$

There are n^2 elements in matrix multiplication result, so total is

$m_{\text{total}} = n^2(n^2+n+2) \approx O(n^4)$ $f_{\text{total}} = 2n^3 \rightarrow q = f_{\text{total}}/m_{\text{total}} = 1/n$

Better approach

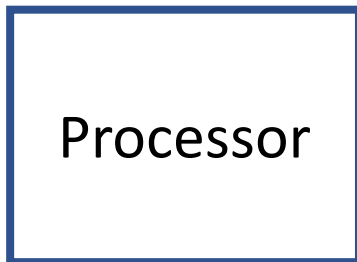
A(,) in row-major and B(,) in column-major orders

A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

Logical view

Memory layout

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33



-	-	-	-
-	-	-	-

Cache (fast but small)

A(,) in row-major and B(,) in column-major orders

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

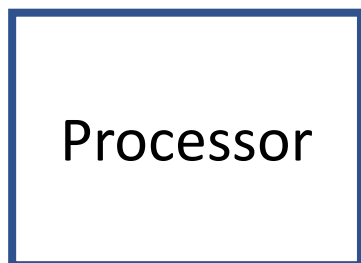
Spatial locality

Memory layout

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00*B00 + A01*B10 + A02*B20 + A03*B30$

1. Read A00 into fast memory
2. Read B00 into fast memory



B00	B10	B20	B30
A00	A01	A02	A03

Cache (fast but small)

A(,) in row-major and B(,) in column-major orders

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	
A30	A31	A32	A33	B30	B31	B32	B33	

Memory layout

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

Compute $C00 = A00*B00 + A01*B10 + A02*B20 + A03*B30$

3. Compute $C00 = A00*B00$
4. Compute $C00 = C00 + A01*B10$
5. and the rest



B00	B10	B20	B30
A00	A01	A02	A03

Cache (fast but small)

Processor finds all required elements in cache 😊
[Very few cache miss]

[Row]-[Column] major orders: #Comm and #Comp

A00	A01	A02	A03	B00	B01	B02	B03	Logical view
A10	A11	A12	A13	B10	B11	B12	B13	
A20	A21	A22	A23	B20	B21	B22	B23	Memory layout
A30	A31	A32	A33	B30	B31	B32	B33	

A00	A01	A02	A03	A10	A11	A12	A13	A20	A21	A22	A30	A30	A31	A32	A33
B00	B10	B20	B30	B01	B11	B21	B31	B02	B12	B22	B32	B03	B13	B23	B33

To compute $C_{00} = A_{00} * B_{00} + A_{01} * B_{10} + A_{02} * B_{20} + A_{03} * B_{30}$

$$m_{C00} = \begin{aligned} & \text{read } A(0,) \text{ once: } \#n \text{ memory access} \\ & + \text{read } B(,0) \text{ once: } \#n \text{ memory access} \\ & + \text{read/write } C00 \text{ once: } \#2 \text{ memory access} \\ & = (2n+2) \end{aligned}$$

$$f_{\text{C00}} = n \text{ multiplications} + n \text{ additions} = 2n$$

There are n^2 elements in matrix multiplication result, so total is

$$m_{\text{total}} = n^2(2n+2) \approx O(2n^3) \quad f_{\text{total}} = 2n^3 \quad \rightarrow \quad q = f_{\text{total}}/m_{\text{total}} \approx 1$$

³⁴
n times better!

+ exploit temporal locality of $A(i,)$

[Row]-[Column] major orders: #Comm and #Comp

A00	A01	A02	A03		B00	B01	B02	B03
A10	A11	A12	A13		B10	B11	B12	B13
A20	A21	A22	A23		B20	B21	B22	B23
A30	A31	A32	A33		B30	B31	B32	B33

Compute $C00 = A00*B00+A01*B10+A02*B20+A03*B30$

$C01 = A00*B01+A01*B11+A02*B21+A03*B31$

$C02 = A00*B02+A01*B12+A02*B22+A03*B32$

$C03 = A00*B03+A01*B13+A02*B23+A03*B33$

Compute one
row of C

$$\begin{aligned} m_{C(0,)} &= \text{read } A(0,) \text{ once: } \#n \text{ memory access} \\ &+ \text{read } B(i) \text{ } n \text{ times: } \#n^2 \text{ memory access} \\ &+ \text{read/write } C(0,) \text{ once: } \#2n \text{ memory access} \\ &= (n^2 + 3n) \end{aligned}$$

$$f_{C00} = n^2 \text{ multiplications} + n^2 \text{ additions} = 2n^2$$

Total cost for matrix multiplication

$$m_{\text{total}} = n(n^2 + 3n) \approx O(n^3) \quad f_{\text{total}} = 2n^3 \quad \rightarrow \quad q = f_{\text{total}}/m_{\text{total}} = 2$$

even better!

What assumptions did we make?

A00	A01	A02	A03		B00	B01	B02	B03
A10	A11	A12	A13		B10	B11	B12	B13
A20	A21	A22	A23		B20	B21	B22	B23
A30	A31	A32	A33		B30	B31	B32	B33

1. Cache memory is sufficiently large to store one row of A
2. and one column of B
3. Cache memory access has 0 overhead

In a real system, the advantage will be smaller compared to our hypothetical machine

Achieving $q > 2$

Partition matrix into blocks

$$\begin{bmatrix} \mathbf{A00} & \mathbf{A01} \\ \mathbf{A10} & \mathbf{A11} \end{bmatrix} \leftarrow \begin{bmatrix} A00 & A01 & A02 & A03 \\ A10 & A11 & A12 & A13 \\ A20 & A21 & A22 & A23 \\ A30 & A31 & A32 & A33 \end{bmatrix}$$

where

$$\mathbf{A00} = \begin{bmatrix} A00 & A01 \\ A10 & A11 \end{bmatrix}$$

$$\mathbf{A01} = \begin{bmatrix} A02 & A03 \\ A12 & A13 \end{bmatrix}$$

etc.

Bold font is used to represent a sub-matrix.

Multiplication after partitioning

A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

A00	A01	B00	B01
A10	A11	B10	B11

Then matrix multiplication result is

$$\begin{vmatrix} c00 & c01 \\ c10 & c11 \end{vmatrix} = \begin{vmatrix} A00*B00+A01*B10 & A00*B01+A01*B11 \\ A10*B00+A11*B10 & A10*B01+A11*B11 \end{vmatrix}$$

Multiplication after partitioning

A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

$$\begin{array}{|c|c|} \hline \mathbf{A00} & \mathbf{A01} \\ \hline \mathbf{A10} & \mathbf{A11} \\ \hline \end{array}
 \begin{array}{|c|c|} \hline \mathbf{B00} & \mathbf{B01} \\ \hline \mathbf{B10} & \mathbf{B11} \\ \hline \end{array}
 =
 \begin{array}{|c|c|} \hline \mathbf{A00*B00+A01*B10} & \mathbf{A00*B01+A01*B11} \\ \hline \mathbf{A10*B00+A11*B10} & \mathbf{A10*B01+A11*B11} \\ \hline \end{array}$$

Assume we can fit three such sub-matrices in the cache.



Cache (fast but small)

Multiplication after partitioning

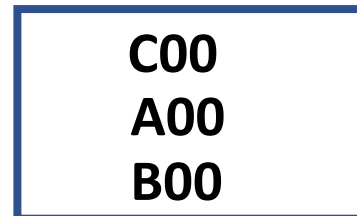
A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

A00	A01
A10	A11

B00	B01
B10	B11

$$\begin{array}{|c|c|} \hline C00 & C01 \\ \hline C10 & C11 \\ \hline \end{array} = \begin{array}{|cc|} \hline A00*B00+A01*B10 & A00*B01+A01*B11 \\ \hline A10*B00+A11*B10 & A10*B01+A11*B11 \\ \hline \end{array}$$

Assume we can fit three such sub-matrices in the cache.



Cache (fast but small)

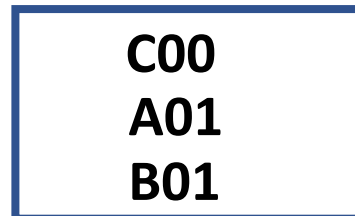
We can compute **entire** sub-matrix operation smoothly!
C00=A00*B00

Multiplication after partitioning

A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

$$\begin{array}{|c|c|} \hline A00 & A01 \\ \hline A10 & A11 \\ \hline \end{array}
 \begin{array}{|c|c|} \hline B00 & B01 \\ \hline B10 & B11 \\ \hline \end{array}
 =
 \begin{array}{|c|c|} \hline C00 & C01 \\ \hline C10 & C11 \\ \hline \end{array}
 =
 \begin{array}{|c|c|} \hline A00*B00+A01*B10 & A00*B01+A01*B11 \\ \hline A10*B00+A11*B10 & A10*B01+A11*B11 \\ \hline \end{array}$$

Assume we can fit three such sub-matrices in the cache.



Cache (fast but small)

We can compute **entire** sub-matrix operation smoothly!

$$C00 = C00 + A01 * B10$$

Multiplication after partitioning

A00	A01	A02	A03	B00	B01	B02	B03
A10	A11	A12	A13	B10	B11	B12	B13
A20	A21	A22	A23	B20	B21	B22	B23
A30	A31	A32	A33	B30	B31	B32	B33

A00	A01
A10	A11

B00	B01
B10	B11

$$\begin{vmatrix} c00 & c01 \\ c10 & c11 \end{vmatrix} = \begin{vmatrix} A00*B00+A01*B10 & A00*B01+A01*B11 \\ A10*B00+A11*B10 & A10*B01+A11*B11 \end{vmatrix}$$

What can we do if these sub-matrices do not fit in Cache?

Answer: partition into multiple smaller sub-matrices such that they can be placed in the cache.

(this is the idea behind Tiled Matrix Multiplication)

Tiled matrix multiplication

Consider A,B,C to be N-by-N matrices of b-by-b subblocks where $b=n / N$ is called the **block size**

```
for(i = 0; i<N; i++)
```

```
  for(j = 0; j<N; j++)
```

```
    {read block C(i,j) into fast memory}
```

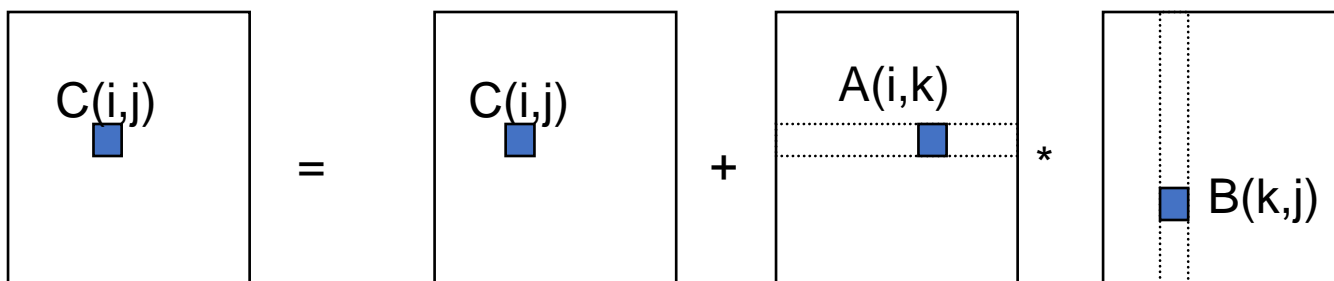
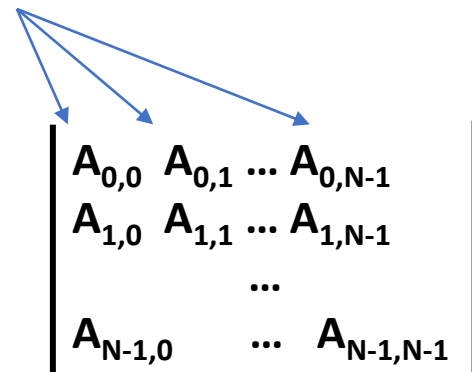
```
    for(k = 1; k<N; k++)
```

```
      {read block A(i,k) into fast memory}
```

```
      {read block B(k,j) into fast memory}
```

```
      C(i,j) = C(i,j) + A(i,k) * B(k,j) {do a matrix multiply on blocks}
```

```
    {write block C(i,j) back to slow memory}
```



Tiled matrix multiplication: #Comm and #Comp

$$\begin{aligned} m &= N \cdot n^2 && \text{read blocks of B } N^3 \text{ times } (N^3 \cdot b^2 = N^3 \cdot (n/N)^2 = N \cdot n^2) \\ &+ N \cdot n^2 && \text{read blocks of A } N^3 \text{ times} \\ &+ 2n^2 && \text{read and write each block of C once} \\ &= (2N + 2) \cdot n^2 \end{aligned}$$

$$f = 2n^3$$

So computational intensity $q = f / m = 2n^3 / ((2N + 2) \cdot n^2)$
 $\approx n / N = b$ for large n

- By choosing $b > 2$ we can achieve $q > 2$
- Additionally q is proportional to block size b

Advantages of tiled matrix multiplication

Computational intensity $q = f / m \approx b$ for large n

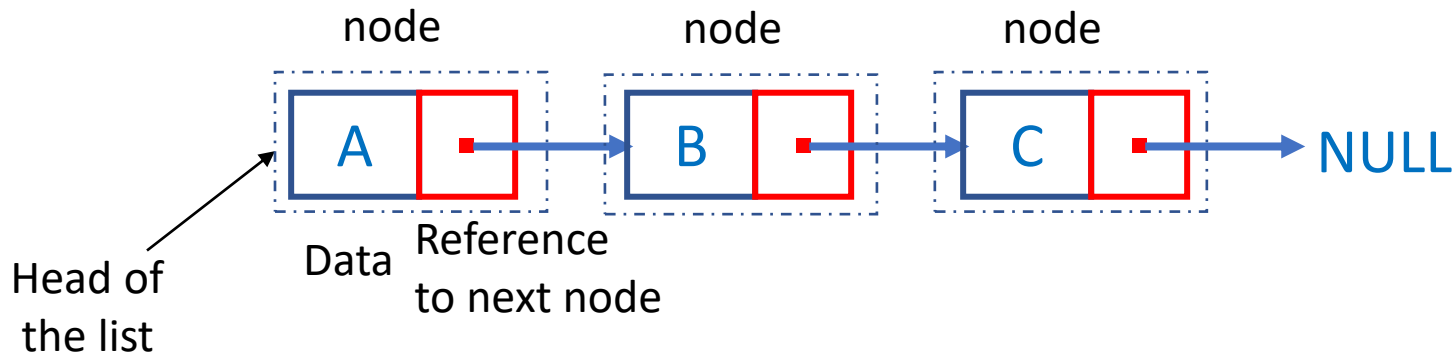
Advantages:

- Enhancement in computational intensity
- Flexibility depending on size of fast memory in machine (choose block size b according to system cache)

Linked List (Recap from last week)

A 'linked list' is a

- linear collection of data elements called 'nodes'
- each node points to the next node in the list
- unlike arrays, linked list nodes are not stored at contiguous locations; they are linked using pointers as shown below.



**Can we have a cache-efficient linked-list?
(or any other data structures)
Try yourself**