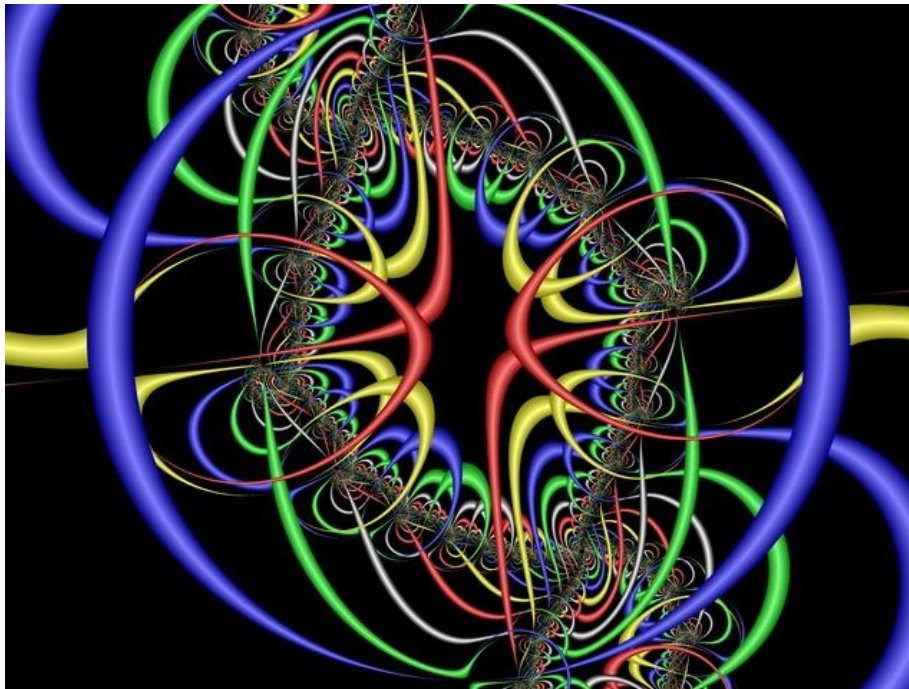


cse5441 - parallel computing

loop analysis and transformation



locality

review

- clustering references to individual variables is good
(temporal locality)
- stride-1 reference patterns are good
(spatial locality)

matrix summation

let:

- cold cache, $\text{sizeof}(\text{int}) = 4$ bytes, 16-byte cache blocks
- fully associative LRU cache
- $C = \text{sizeof}(\text{int})(M)(N)$ (assume $M, N \geq 64$)
- sum, i, j in registers

```
int sumarrayrows(int a[M][N])
{
    int sum = 0;

    for (int i = 0; i < M; i++)
        for (int j = 0; j < N; j++)

            sum += a[i][j]

    return sum;
}
```

```
int sumarraycols(int a[M][N])
{
    int sum = 0;

    for (int j = 0; j < N; j++)
        for (int i = 0; i < M; i++)

            sum += a[i][j]

    return sum;
}
```

matrix summation

let:

- cold cache, $\text{sizeof}(\text{int}) = 4$ bytes, 16-byte cache blocks
- fully associative LRU cache
- $C = (.5)(M)(N)$ (assume $M, N \geq 64$)
- sum in register

```
int sumarrayrows(int a[M][N])
{
    int sum = 0;

    for (int i = 0; i < M; i++)
        for (int j = 0; j < N; j++)

            sum += a[i][j]

    return sum;
}
```

```
int sumarraycols(int a[M][N])
{
    int sum = 0;

    for (int j = 0; j < N; j++)
        for (int i = 0; i < M; i++)

            sum += a[i][j]

    return sum;
}
```

matrix summation

let:

- cold cache, $\text{sizeof}(\text{int}) = 4$ bytes, 16-byte cache blocks
- fully associative LRU cache
- $C < (.5)(M)(B)$ (assume $M, N \geq 64$)
- sum in register

```
int sumarrayrows(int a[M][N])
{
    int sum = 0;

    for (int i = 0; i < M; i++)
        for (int j = 0; j < N; j++)

            sum += a[i][j]

    return sum;
}
```

```
int sumarraycols(int a[M][N])
{
    int sum = 0;

    for (int j = 0; j < N; j++)
        for (int i = 0; i < M; i++)

            sum += a[i][j]

    return sum;
}
```

it's your turn . . .

- let:
- cold cache, $\text{sizeof}(\text{int}) = 4$ bytes, 32-byte cache blocks
 - fully associative LRU cache
 - $C = 4(N)(P)$ (in bytes)
 - matrix allocated as contiguous memory (one "new")
 - $P \% 8 = 0$, $P < M < NP/8$
 - assume $N, M, P > \text{trivial}$

```
int summat_1(int a[M][N][P])
{
    int sum = 0;    // register

    for (int i = 0; i < M; i++)
        for (int j = 0; j < N; j++)
            for (int k = 0; k < P; k++)

                sum += a[ i ][ j ][ k ]

    return sum;
}
```

```
int summat_2(int a[M][N][P])
{
    int sum = 0;    // register

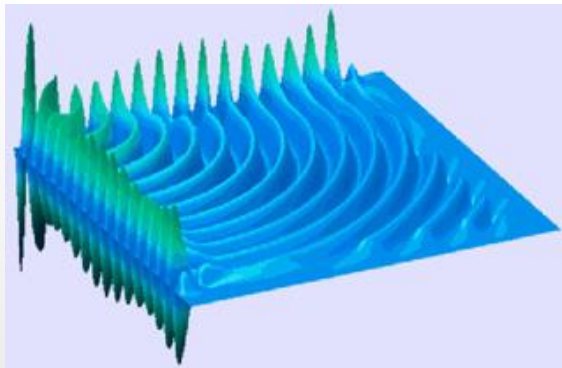
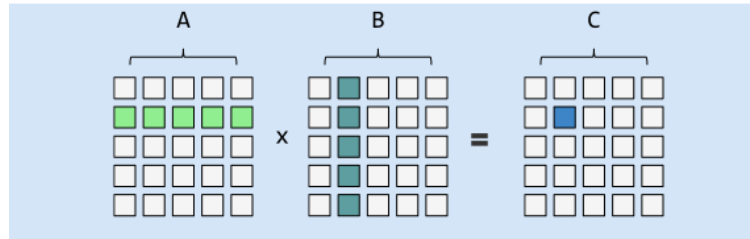
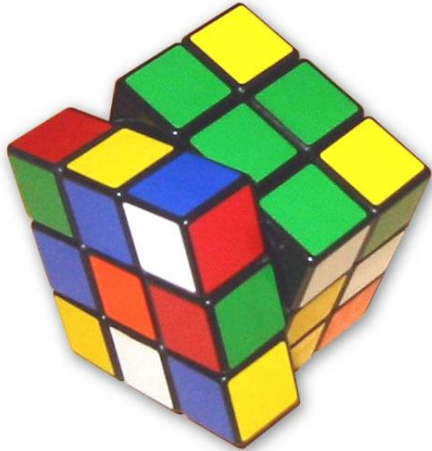
    for (int k = 0; k < P; k++)
        for (int j = 0; j < N; j++)
            for (int i = 0; i < M; i++)

                sum += a[ i ][ j ][ k ]

    return sum;
}
```

what are the access stride, hit rate and cache contents at end of loop nest?

matrix multiplication



```
{(-0.4 X02MPS - 0.32 X14MPS - 0.6 X23MPS - 0.48 X36MPS + 10. X39MPS,
-1. X01MPS + 1. X02MPS + 1. X03MPS == 0. && -1.06 X01MPS + 1. X04MPS == 0. &&
1. X01MPS ≤ 80. && -1. X02MPS + 1.4 X14MPS ≤ 0. &&
-1. X06MPS - 1. X07MPS - 1. X08MPS - 1. X09MPS + 1. X14MPS + 1. X15MPS == 0. &&
-1.06 X06MPS - 1.06 X07MPS - 0.96 X08MPS - 0.86 X09MPS + 1. X16MPS == 0. &&
1. X06MPS - 1. X10MPS ≤ 80. && 1. X07MPS - 1. X11MPS ≤ 0. && 1. X08MPS - 1. X12MPS ≤ 0. &&
1. X09MPS - 1. X13MPS ≤ 0. && -1. X22MPS + 1. X23MPS + 1. X24MPS + 1. X25MPS == 0. &&
-0.43 X22MPS + 1. X26MPS == 0. && 1. X22MPS ≤ 500. && -1. X23MPS + 1.4 X36MPS ≤ 0. &&
-0.43 X28MPS - 0.43 X29MPS - 0.39 X30MPS - 0.37 X31MPS + 1. X38MPS == 0. &&
1. X28MPS + 1. X29MPS + 1. X30MPS + 1. X31MPS - 1. X36MPS + 1. X37MPS + 1. X39MPS == 44. &&
1. X28MPS - 1. X32MPS ≤ 500. && 1. X29MPS - 1. X33MPS ≤ 0. &&
1. X30MPS - 1. X34MPS ≤ 0. && 1. X31MPS - 1. X35MPS ≤ 0. &&
2.364 X10MPS + 2.386 X11MPS + 2.408 X12MPS + 2.429 X13MPS - 1. X25MPS + 2.191 X32MPS +
2.219 X33MPS + 2.249 X34MPS + 2.279 X35MPS ≤ 0. && -1. X03MPS + 0.109 X22MPS ≤ 0. &&
-1. X15MPS + 0.109 X28MPS + 0.108 X29MPS + 0.108 X30MPS + 0.107 X31MPS ≤ 0. &&
0.301 X01MPS - 1. X24MPS ≤ 0. &&
0.301 X06MPS + 0.313 X07MPS + 0.313 X08MPS + 0.326 X09MPS - 1. X37MPS ≤ 0. &&
1. X04MPS + 1. X26MPS ≤ 310. && 1. X16MPS + 1. X38MPS ≤ 300. && X01MPS ≥ 0 && X02MPS ≥ 0 &&
X03MPS ≥ 0 && X04MPS ≥ 0 && X06MPS ≥ 0 && X07MPS ≥ 0 && X08MPS ≥ 0 && X09MPS ≥ 0 &&
X10MPS ≥ 0 && X11MPS ≥ 0 && X12MPS ≥ 0 && X13MPS ≥ 0 && X14MPS ≥ 0 && X15MPS ≥ 0 &&
X16MPS ≥ 0 && X22MPS ≥ 0 && X23MPS ≥ 0 && X24MPS ≥ 0 && X25MPS ≥ 0 && X26MPS ≥ 0 &&
X28MPS ≥ 0 && X29MPS ≥ 0 && X30MPS ≥ 0 && X31MPS ≥ 0 && X32MPS ≥ 0 && X33MPS ≥ 0 &&
X34MPS ≥ 0 && X35MPS ≥ 0 && X36MPS ≥ 0 && X37MPS ≥ 0 && X38MPS ≥ 0 && X39MPS ≥ 0},
{X01MPS, X02MPS, X03MPS, X04MPS, X06MPS, X07MPS, X08MPS, X09MPS, X10MPS,
X11MPS, X12MPS, X13MPS, X14MPS, X15MPS, X16MPS, X22MPS, X23MPS,
X24MPS, X25MPS, X26MPS, X28MPS, X29MPS, X30MPS, X31MPS, X32MPS,
X33MPS, X34MPS, X35MPS, X36MPS, X37MPS, X38MPS, X39MPS}}
```

<http://www.stoimen.com/blog/wp-content/uploads/2012/11/3.-Matrix-Multiplication.png>

https://reference.wolfram.com/language/tutorial/Files/ConstrainedOptimizationLinearProgramming.en/0_19.gif

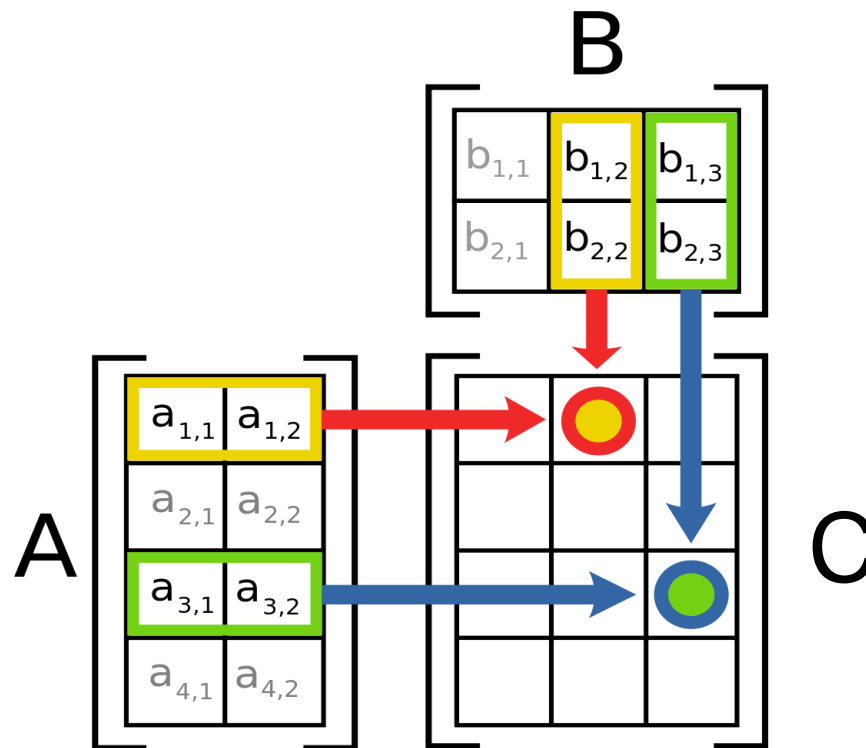
http://chemwiki.ucdavis.edu/@api/deki/files/9875/Rubiks_cube.jpg?size=bestfit&width=255&height=220&revision=1

<http://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-728-applied-quantum-and-statistical-physics-fall-2006/6-728f06.jpg>

<https://www.google.com/search?q=illustration&tbn=isch&tbo=u&source=univ&sa=X&ved=0CFsQsARqFQoTCI-y2JvM-8YCFRZ9iAodXTgNWQ&biw=1920&bih=969&tbn=isch&q=linear+transformation&imgc=TIIPHaDabGJoM%3A>

matrix multiplication

$$C = A \times B$$
$$C_{ij} = \sum_k A_{ik} * B_{kj}$$



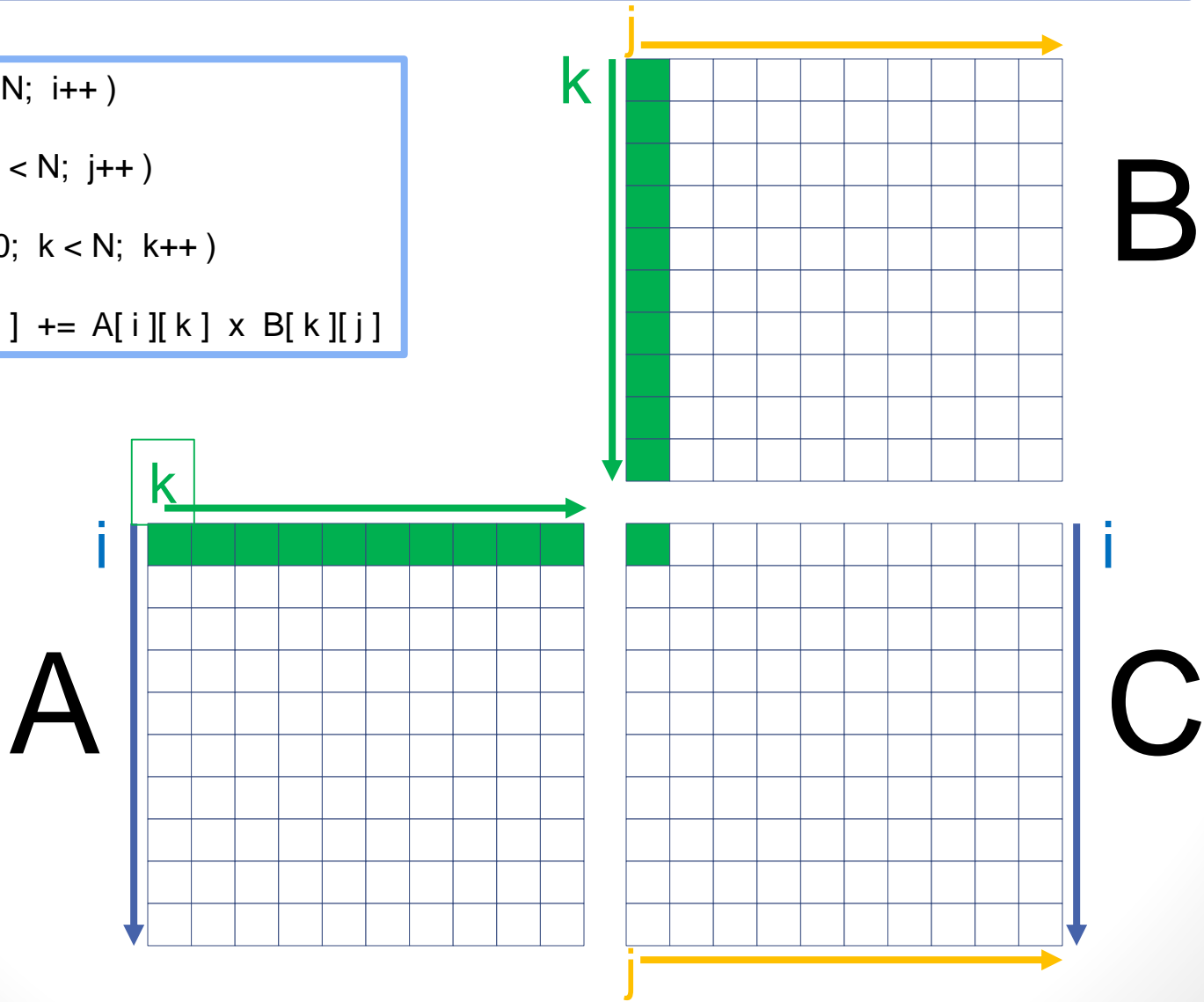
matrix multiplication

$$C = A \times B$$
$$C_{ij} = \sum_i \sum_j \sum_k A_{ik} * B_{kj}$$

```
for ( i = 0; i < N; i++ )  
    for ( j = 0; j < N; j++ )  
        for ( k = 0; k < N; k++ )  
            C[ i ][ j ] += A[ i ][ k ] x B[ k ][ j ]
```

matrix multiplication

```
for ( i = 0; i < N; i++ )  
  for ( j = 0; j < N; j++ )  
    for ( k = 0; k < N; k++ )  
      C[ i ][ j ] += A[ i ][ k ] x B[ k ][ j ]
```



m-mult cache performance

inner-loop analysis method

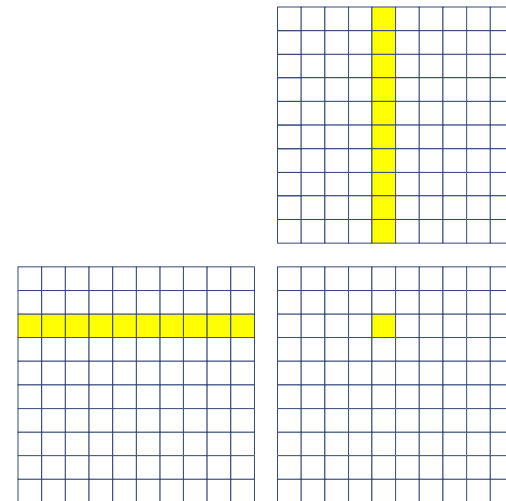
```
I for ( i = 0; i < N; i++ )
J   for ( j = 0; j < N; j++ )
K     for ( k = 0; k < N; k++ )
      C[i][j] += A[i][k] x B[k][j]
```

let: N is very large, $1/N \approx 0$
 $N \gg E$
fully associative cache
LRU
 $B = 4 * \text{sizeof}(\text{int})$

consider access pattern of inner loop:

- k is iterating
- what is stride of A ?
- what is stride of B ?
- what is stride of C ?

A	B	C
0.25	+ 1.00	+ 0.00 = 1.25



inner-loop iteration pattern

this is the “misses per iteration” of the inner loop

matrix multiplication

$$C = A \times B$$
$$C_{ij} = \sum_i \sum_j \sum_k A_{ik} * B_{kj}$$

```
for ( i = 0; i < N; i++ )  
    for ( j = 0; j < N; j++ )  
        for ( k = 0; k < N; k++ )  
            C[ i ][ j ] += A[ i ][ k ] x B[ k ][ j ]
```

matrix multiplication

disclaimer

WLOG, all examples are integer

m-mult cache performance

JKI

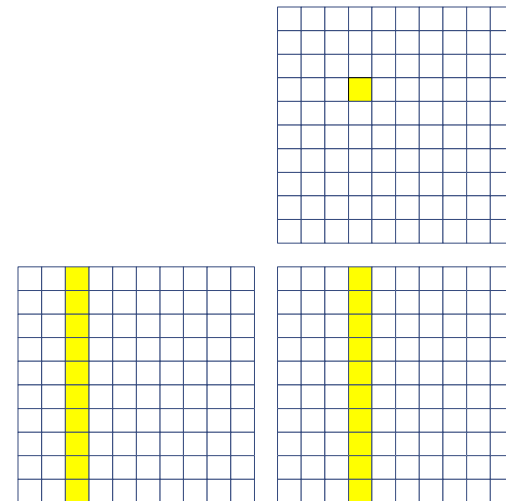
```
J for ( j = 0; j < N; j++ )
K   for ( k = 0; k < N; k++ )
    for ( i = 0; i < N; i++ )
        C[i][j] += A[i][k] x B[k][j]
```

let: N is very large, $1/N \approx 0$
 $N \gg E$
fully associative cache
LRU
 $B = 4 * \text{sizeof}(\text{int})$

consider access pattern of inner loop:

- i is iterating
- what is stride of A ?
- what is stride of B ?
- what is stride of C ?

A	B	C
1.00	+ 0.00	+ 1.00 = 2.00



inner-loop iteration pattern

m-mult cache performance

IKJ

I
K
J

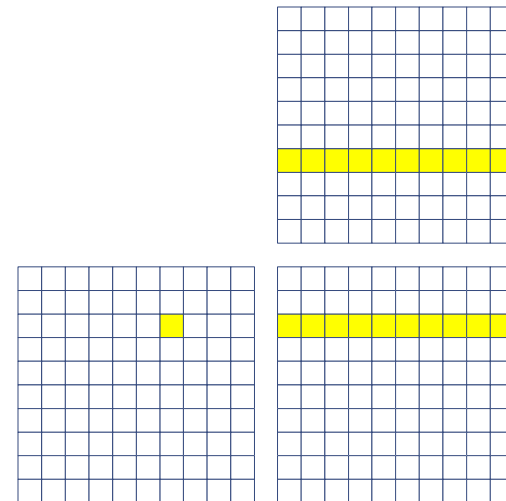
```
for ( i = 0; i < N; i++ )  
    for ( k = 0; k < N; k++ )  
        for ( j = 0; j < N; j++ )  
            C[i][j] += A[i][k] x B[k][j]
```

let: N is very large, $1/N \approx 0$
N >> E
fully associative cache
LRU
 $B = 4 * \text{sizeof}(\text{int})$

consider access pattern of inner loop:

- j is iterating
- what is stride of A?
- what is stride of B?
- what is stride of C?

A	B	C
0.00	+ 0.25	+ 0.25 = 0.50



inner-loop iteration pattern

m-mult cache performance

summary

let: N is very large, $1/N \approx 0$
 $N \gg E$
fully associative cache
LRU
 $B = 4 * \text{sizeof}(\text{float})$

I
J
K

```
for ( i = 0; i < N; i++ )  
    for ( j = 0; j < N; j++ )  
        for ( k = 0; k < N; k++ )  
            C[i][j] += A[i][k] x B[k][j]
```

A	B	C	
0.25	+ 1.00	+ 0.00	= 1.25

J
K
I

```
for ( j = 0; j < N; j++ )  
    for ( k = 0; k < N; k++ )  
        for ( i = 0; i < N; i++ )  
            C[i][j] += A[i][k] x B[k][j]
```

A	B	C	
1.00	+ 0.00	+ 1.00	= 2.00

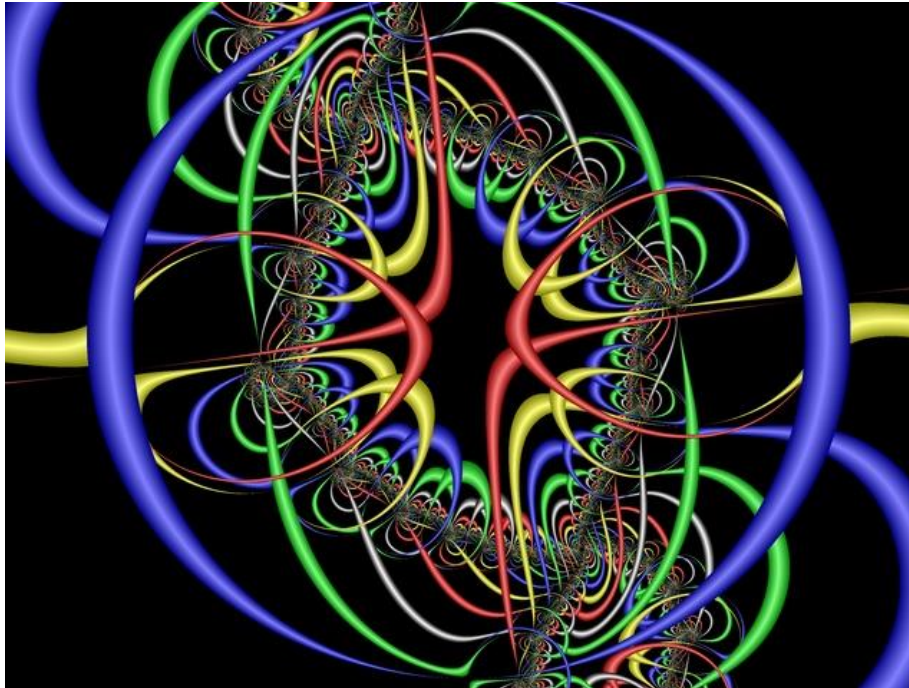
I
K
J

```
for ( i = 0; i < N; i++ )  
    for ( k = 0; k < N; k++ )  
        for ( j = 0; j < N; j++ )  
            C[i][j] += A[i][k] x B[k][j]
```

A	B	C	
0.00	+ 0.25	+ 0.25	= 0.50

cse5441 - parallel computing

loop analysis and transformation



m-mult cache performance

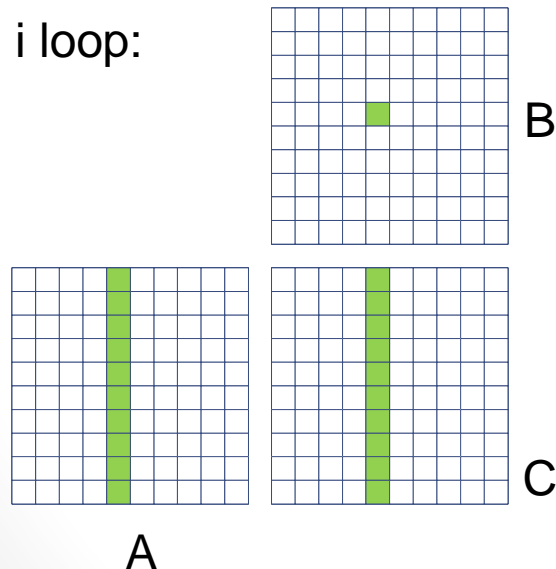
total miss analysis method

```

I   for ( i = 0; i < N; i++ )
J   for ( j = 0; j < N; j++ )
K   for ( k = 0; k < N; k++ )
    C[i][j] += A[i][k] x B[k][j]
    
```

let: N is very large, $1/N \approx 0$
 $N \gg 4E$
 fully associative cache
 LRU
 $B = 4 * \text{sizeof}(\text{int})$
 all register variables are held in cache

i loop:



	A	B	C
I	N	N	N
J	N	N	$\frac{N}{B}$
K	$\frac{N}{B}$	N	1
	$\frac{N^3}{B}$	N^3	$\frac{N^2}{B}$

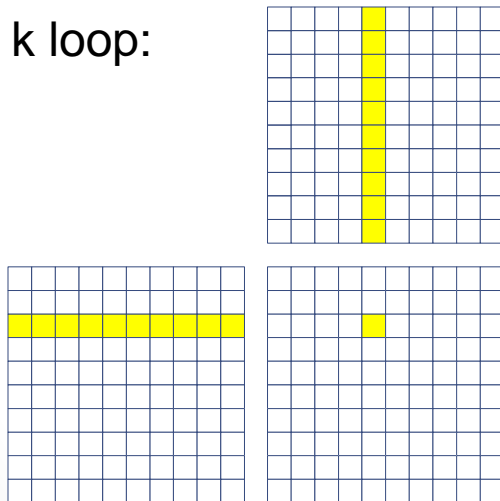
m-mult cache performance

total miss analysis method

J for (j = 0; j < N; j++)
 I for (i = 0; i < N; i++)
 K for (k = 0; k < N; k++)
 C[i][j] += A[i][k] x B[k][j]

let: N is very large, $1/N \approx 0$
 $N \gg 4E$
 fully associative cache
 LRU
 $B = 4 * \text{sizeof(int)}$
 all register variables are held in cache

k loop:



	A	B	C
J	N	N	N
I	N	N	N
K	$\frac{N}{B}$	N	1
	$\frac{N^3}{B}$	N^3	N^2

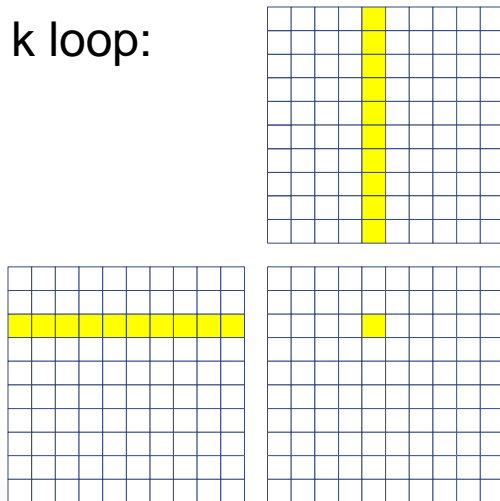
m-mult cache performance

total miss analysis method, example 2 (IJK)

I for (i = 0; i < N; i++)
 J for (j = 0; j < N; j++)
 K for (k = 0; k < N; k++)
 C[i][j] += A[i][k] x B[k][j]

let: $C = (N^2 + 4N + \epsilon) * \text{sizeof}(\text{element})$
 fully associative cache
 LRU
 $B = 4 * \text{sizeof}(\text{element})$
 all register variables are held in cache

k loop:



	A	B	C
I	N	1	N
J	1	$\frac{N}{B}$	$\frac{N}{B}$
K	$\frac{N}{B}$	N	1
	$\frac{N^2}{B}$	$\frac{N^2}{B}$	$\frac{N^2}{B}$

it's your turn ...

```
for ( ? = 0; ? < N; ?++ )  
    for ( ? = 0; ? < N; ?++ )  
        for ( ? = 0; ? < N; ?++ )  
            C[ i ][ j ] += A[ i ][ k ] x B[ k ][ j ]
```

let: $C = 3N * \text{sizeof}(\text{element}) + \epsilon$
fully associative cache
LRU
 $B = 8 * \text{sizeof}(\text{element})$

- what is an optimal loop ordering?
- what is the hit rate of your solution?
- what analysis method did you use?

loop permutation

summary

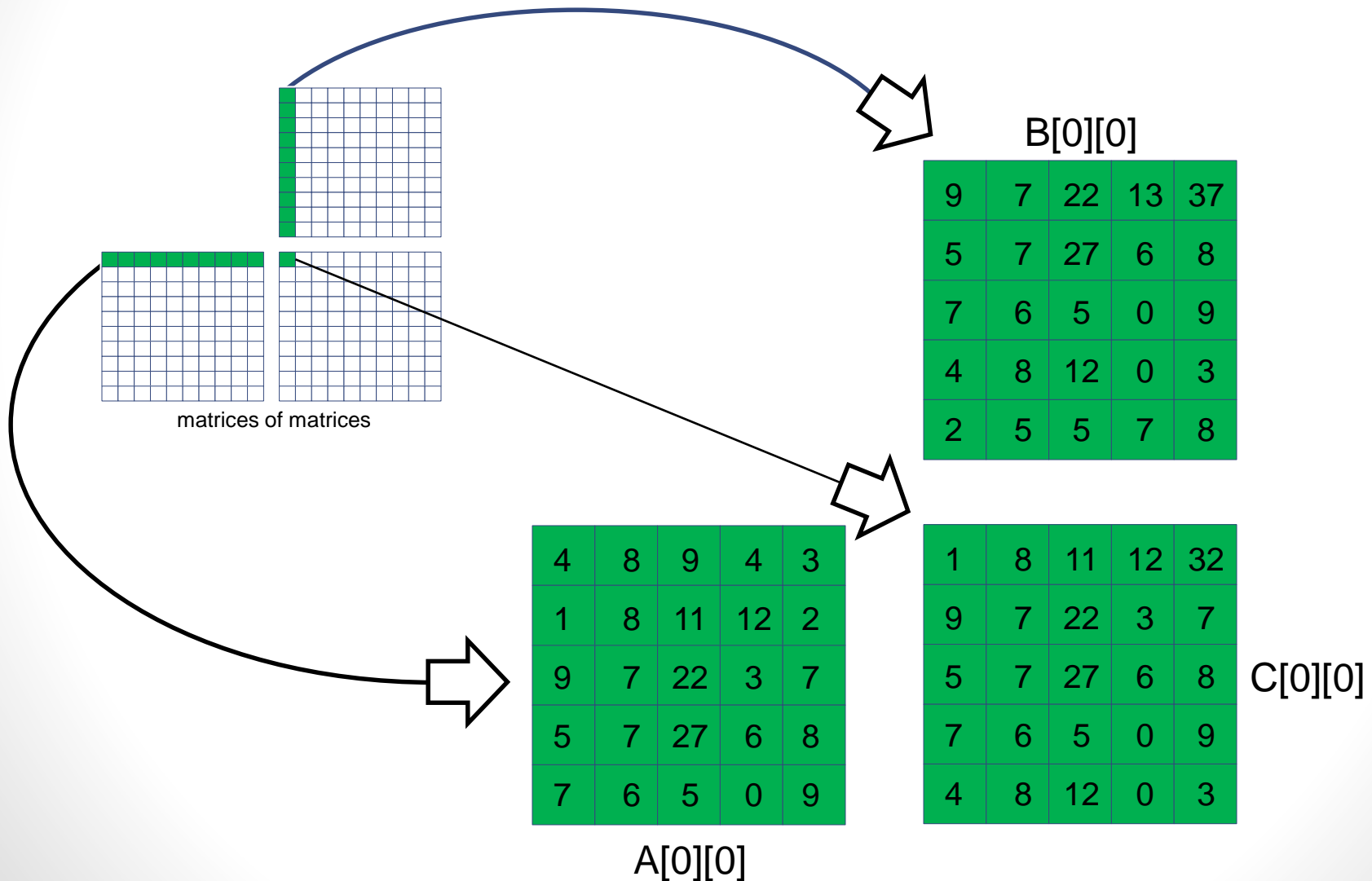
cache-friendly code will:

- consider loop permutations
- strive for stride-1 access
(or at least $< B$)
- minimize long strides

variation in loop permutation:

- affect the cache miss ratio
- affect the minimum cache needed
for optimal performance

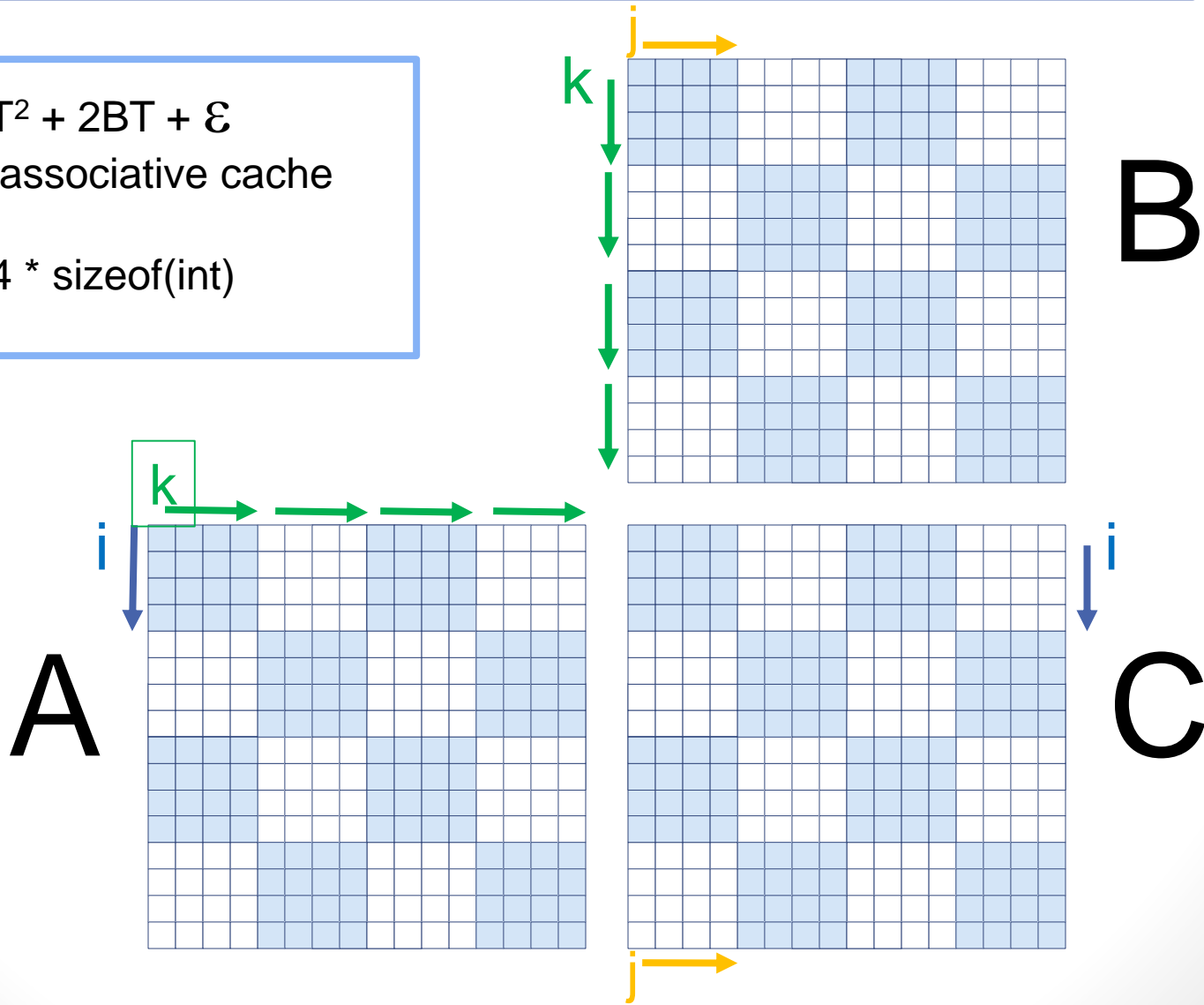
blocking (tiling)



blocking (tiling)

let: $C = T^2 + 2BT + \varepsilon$
fully associative cache
LRU
 $B = 4 * \text{sizeof}(\text{int})$

IJK loop



it's your turn ...

```
I   for ( i = 0; i < N; i++ )
K   for ( k = 0; k < N; k++ )
J   for ( j = 0; j < N; j++ )
    C[ i ][ j ] += A[ i ][ k ] x B[ k ][ j ]
```

- $N = 2000$
- $\text{sizeof}(\text{int}) = 4$
- $b = 5$
- $E = 1$
- $C = 16k$ (4096 elements)
- $A[i][k]$ moved to register (temporal locality)

what tile size would you use?

inner loop unrolling

ORIGINAL

```
for ( i = 0; i < n; i++ )  
    for ( j = 0; j < n; j++ )  
        y[i] = y[i] + a[ i ][ j ] * x[ j ];
```

UNROLLED (for simplicity, let $j \% 4 = 0$)

```
for ( i = 0; i < n; i++ )  
    for ( j = 0; j < n; j += 4 )  
        y[i] = y[i] + a[ i ][ j ] * x[ j ];  
        y[i] = y[i] + a[ i ][ j+1 ] * x[ j+1 ];  
        y[i] = y[i] + a[ i ][ j+2 ] * x[ j+2 ];  
        y[i] = y[i] + a[ i ][ j+3 ] * x[ j+3 ];
```

```
for ( i = 0; i < n; i++ )  
    for ( j = 0; j < n; j += 4 )  
        y[i] = y[i] + a[ i ][ j ] * x[ j ];  
                + a[ i ][ j+1 ] * x[ j+1 ];  
                + a[ i ][ j+2 ] * x[ j+2 ];  
                + a[ i ][ j+3 ] * x[ j+3 ];
```

outer loop unrolling (unroll / jam)

ORIGINAL

```
for ( i = 0; i < 2n; i++ )  
    for ( j = 0; j < m; j++ )  
        loop-body(j, i)
```

UNROLLED

```
for ( i = 0; i < 2n; i += 2 )  
    for ( j = 0; j < m; j++ )  
        loop-body(j, i)  
    for ( j = 0; j < m; j++ )  
        loop-body(j, i+1)
```

preserves
execution order

but doesn't
accomplish
anything ...

UNROLLED / JAMMED

```
for ( i = 0; i < 2n; i += 2 )  
    for ( j = 0; j < m; j++ )  
        loop-body(j, i)  
        loop-body(j, i+1)
```

unrolling benefits

but ... does not
preserve
access order

outer loop unrolling (unroll / jam)

example

ORIGINAL

```
for ( i = 0; i < 4n; i++ )  
    for ( j = 0; j < 4n; j++ )  
        y[ i ] = y[ i ] + a[ i ][ j ] * x[ j ]
```

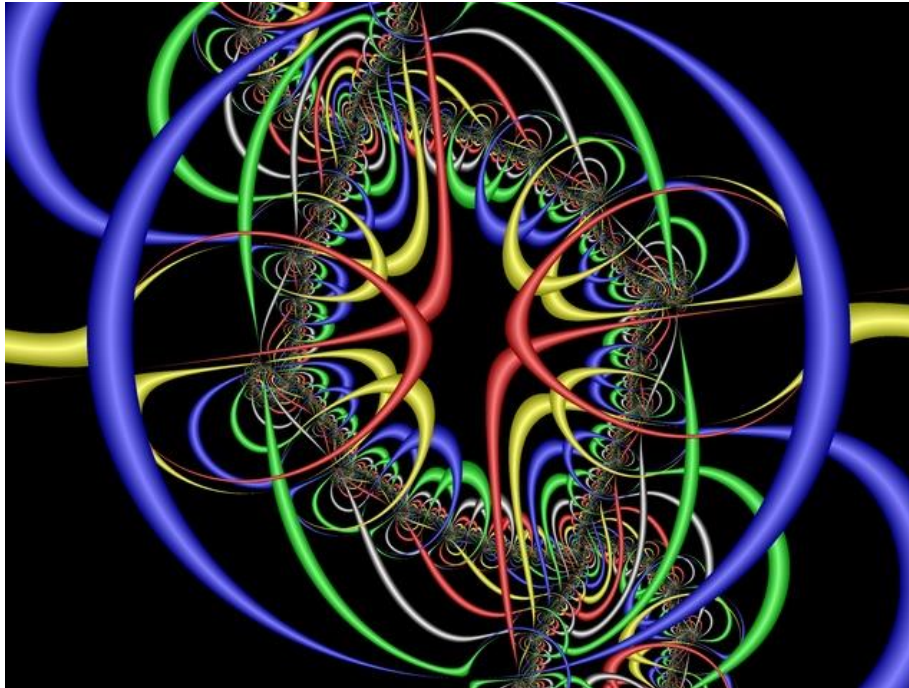
UNROLLED / JAMMED

```
for ( i = 0; i < 4n; i += 4 )  
    for ( j = 0; j < 4n; j++ )  
        y[ i ]    = y[ i ]    + a[ i ][ j ]    * x[ j ]  
        y[ i+1 ] = y[ i+1 ] + a[ i+1 ][ j ] * x[ j ]  
        y[ i+2 ] = y[ i+1 ] + a[ i+2 ][ j ] * x[ j ]  
        y[ i+3 ] = y[ i+1 ] + a[ i+3 ][ j ] * x[ j ]
```

does not preserve
access order

cse5441 - parallel computing

loop analysis and transformation



loops - answers

slide 6: summat_1

- access stride is 1 element, or 4 bytes
- hit rate is $7/8 = 87.5\%$
- $a[M-1][*][*]$ will be in cache

summat_2

- stride is $N * P * |\text{element}|$
- hit rate is 0%
- $a[*][(N - (NP/*M)) \dots N-1][(P-8) \dots (P-1)]$ in cache

loops - answers

slide 21: IKJ will be best, need total miss analysis to differentiate between IKJ and KIJ

inner loop analysis method:

**K: A= .125, B= 1.000, C= 0, total 1.125

**J: A= 0, B= .125, C= .125, total 0.250

**I: A= 1.000, B= 0, C= 1.000, total 2.000

	A	B	C
I	N	N	N
K	$\frac{N}{B}$	N	1
J	1	$\frac{N}{B}$	$\frac{N}{B}$
	$\frac{N^2}{B}$	$\frac{N^3}{B}$	$\frac{N^2}{B}$

	A	B	C
K	N	N	N
I	N	1	N
J	1	$\frac{N}{B}$	$\frac{N}{B}$
	N^2	$\frac{N^2}{B}$	$\frac{N^3}{B}$

loops - answers

slide 25: the loop is already set up optimally $3N^2/B$)

if we allow K and J to iterate, we would need:

A – N elements

B – N^2 elements

C – N elements

so, we pick a tile size such that $T^2 + 2T < C$, also with the restriction that the tile size is a multiple of $B / |\text{element}|$ so that we align with the cache blocks (over- or under-running would waste space and could introduce contention)

solving for T can be done numerically or by "guess and test"

$$24^2 + 2(24) = 624 < 4096/4$$

$$32^2 + 2(32) = 1088 \nless 4096/4$$