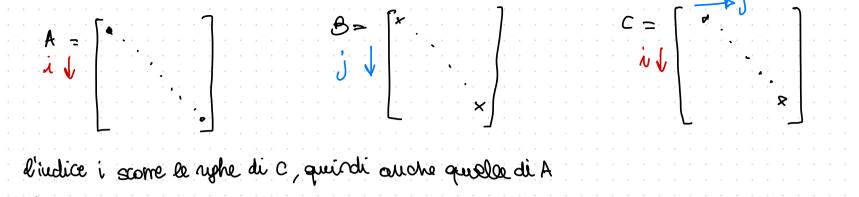
# Esercizio Cache

Andrea Bartolini – <a href="mailto:a.bartolini@unibo.it">a.bartolini@unibo.it</a>

## FP Example: Array Multiplication

```
moltiplicazione di A e B e sonno 8UC
     C = C + A \times B
         All 32 × 32 matrices, 64-bit double-precision elements 30 32 × 32 20 ment
         DGEMM (Double precision GEneral Matrix Multiply) (e ciascun elemento e
     C code: l'implementazione u Cde questa opera composto de 64 bét
zione e con 3 cicli for investort, dove a non-
      void mm (double c[][]; il ciclo+esterno scome la rughe delle
         double a[][], double b[][]) { mothice C, (2) for size_t i, j, k; Some levolonne dolla mothice for (i = 0; i < 32; i = i + 1) C e is either k scorne for (j = 0; j < 32; j = j + 1) for (k = 0; k < 32; k = k + 1) be ushe dolla m. A element of (k = 0; k < 32; k = k + 1)
                                                                colonne della m.B.
                   c[i][j] = c[i][j]
                                   + a[i][k] * b[k][j];
                                                                          Suppongo che gli udinitti base di
A,B,C sutnovino oll'inferno dei
       lacktriangle Addresses of c, a, b in x10, x11, x12, and
                                                               &6[0](0) registri K10, ×11, ×12.
contains \downarrow, \downarrow, k in x5, x6, x7
```



l'india j some le alonne della motrisse C, a le colonne della B l'indice u mon e in c, in quanto e quello che pernette di fare le sperazioni riga x colonna.

tutti di elementi di una rupe di A Cx) vengano Scotter

x colcolore l'élevents ijesmo di C belsono no prender tutte gli elevents della rupe di A empliphicarli per tutte gle clevrent della colonna di B. Injetti velo travo cone vudice di colonna di D. 2 come vudice di rupe di B.

### FP Example: Array Multiplication

RISC-V code:

il codice assembly precarica una Serve di immediation Come 32 ( dimensione dei 3 cicli for Cic 32, je 32, ke 32)

```
load immediate
                             // x28 = 32 (row size/loop end)
               x28,32
                                                                  inizializzo
               x5,0
                             // i = 0; initialize 1st for loop
                                                                  1,8, Kato.
               x6,0
                             // j = 0; initialize 2nd for loop
               \times 7,0
                             // k = 0; initialize 3rd for loop
                             // x30 = i * 2**5  (size of row of c) with
               x30, x5, 5
 a fine 1 slli
                             // x30 = i * size(row) + j & sommano, out offset di
                x30,x30,x6
         add
                             // x30 = byte offset of [i][i] colonno
         slli
                x30, x30, 3
                             // x30 = byte address of c[i][j] shifts x = x^2
                x30, x10, x30
         add
                             // fo = c[i][i] & 10 load dell' elemento c (i)(j)
               (f0)0(x30)
         fld
                             // x29 = k * 2**5  (size of row of b) whice or
         slli
               x29, x7, 5
                             // x29 = k * size(row) + j skdis bit
    Assambadd
                x29,x29,x6
                x29,x29,3
                             // x29 = byte offset of [k][j]
         add
                x29, x12, x29
                             // x29 = byte address of b[k][j]
osco gas
th'aicidli
                f1,0(x29) // f1 = b[k][j]
          ACESSOAB
```

#### FP Example: Array Multiplication

```
add x29, x29, x7   // x29 = i * size(row) + k
        slli x29, x29, 3 // x29 = byte offset of [i][k]
         add x29, x11, x29 // x29 = byte address of a[i][k]
         fld f2,0(x29) // f2 = a[i][k]
         fmul.d f1, f2, f1 // f1 = a[i][k] * b[k][j]
        fadd.d f0, f0, f1 // f0 = c[i][j] + a[i][k] * b[k][j]
        addi x7, x7, 1  // k = k + 1 Quant k=32 & Store
        bltu x7, x28, L3 // if (k < 32) go to L3
         fsd f0,0(x30) // c[i][j] = f0
       addi x6,x6,1 // j=j+1-b Incremento \hat{j} function of \frac{1}{2} bltu x6,x28,L2 // \frac{1}{2} // \frac{1}{2}
         bltu x5,x28,L1 // if (i < 32) go to L1
judinizzo assoluto dib
```

- La cpu dispone di una cache dati a 2 vie da 16KiB complessivi e linee da 64 byte, gestita con stato MESI e con politica di scrittura Write-Around in caso di miss.
- Si considerino A,B,C immagazzinate in memoria a partire rispettivamente dagli
- indirizzi: 0x0100 8000, 0x0100 A000, 0x0100 C000

  1) Si disegni la mappa della memoria più divizzi in memorio degli elementi delle matrici
  (selo i primi egei nelimi.
- 2) Si analizzi la dinamica della cache dati, e, tenendo ben presente che il sistema ha un solo caching agent, si risponda in modo preciso, schematico, conciso e tabellare ai seguenti quesiti:
  - Quali sono gli indici di set e linea, e i tag associati ad A,B,C per il primo e l'ultimo blocco contenenti le matrici?
  - Quante linee di cache occuperanno nel loro insieme? Possono i due vettori e la variabile stare per intero e simultaneamente in cache? se tutte e3 le matrice mescaro ad esser contenute in cache
- 3) Si consideri la dinamica della cache nel calcolo della prima iterazione i=0, j=0, k=0 nel calcolo del primo elmento di C, disegnando lo stato MESI, il contenuto della cache e il valore del bit LRU dopo ogni operazione elementare (Load e Store) e si indichi il numero di accessi, il numero di miss e il numero di cicli di writeback e gli eventuali cicli di write allocate.
- 4) Si indichi il numero di accessi, il numero di miss e il numero di cicli di writeback e gli eventuali cicli di write allocate, nonché lo stato MESI della cache al termine del calcolo del primo elemento di C[0][0]. (si riporti il contenuto del primo ed ultimo set della cache)

faceudo variane k da 0 a 32.

ce istruzioni che perturbano

lo stato delle cadne sono FP Example: Array Multiplication iskutioni di lood/ store. Nel codia: la store e' l' asseguamento.

```
C = C + A \times B
```

i, j, k in x5, x6, x7

le cood cisoro quoendo si All 32 × 32 matrices, 64-bit double-precision elements legono: swedt elements

DGEMM (Double precision GEneral Matrix Multiply)

C code:

```
void mm (double c[][],
         double a[][], double b[][]) {
  size t i, j, k;
  for (i = 0; i < 32; i = i + 1)
                                           qui alobiamo 32 store e 32.2+1 land
    for (j = 0; j < 32; j = j + 1)
                                            quindi a agui ciclo obbiamo: 32×32 local in c
      for (k = 0; k < 32; k = k + 1)
        c[i][j] = c[i][j]
                   + a[i][k] * b[k][j];
                                                32×32 store we
                                                32×32×32×2 road too a e b.
  Addresses of c, a, b in x10, x11, x12, and
```

mmo: i registri x sono registri interi

da quanti byte e conporto un sugalo elevento delle tre matrici? 8 byte mce's sous double se prende a (0) [0), quali indivissi u memoria occupera (partida por 64 bit da 0x 0100 8000)? da 0×0100 8000 Juna 0×0100 8007 c[i][j] + a[i][k] \* b[k][j]; | valore di i e di i c[i][j] = c[i][j]

rendo i (che è l'indice di rige) e la shisto a se di 56t, molt picandolo per 32. Avestoci porto nella porte di memorie associata a quella riga, poi sei somo j e ci porto or un certo valore. Poi devo dirgi che croscum elemento è 8 byte, quindi moltiplico per 8.



