

# Programmazione su schede GPU in CUDA

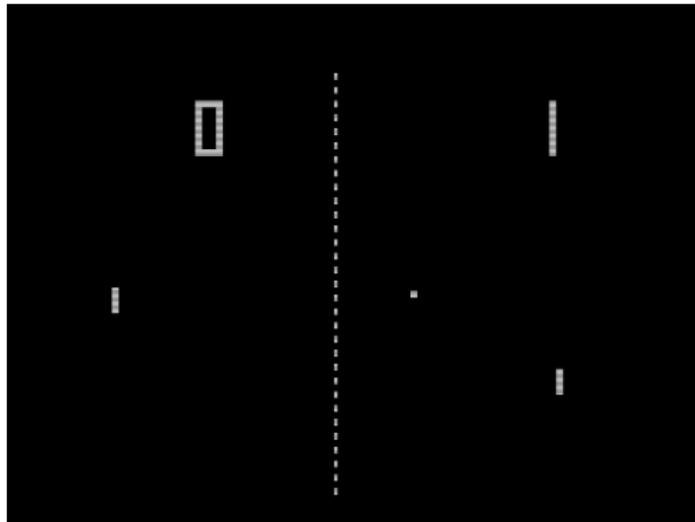
Marco Sansottera



UNIVERSITÀ DEGLI STUDI DI MILANO  
FACOLTÀ DI SCIENZE E TECNOLOGIE

Università degli Studi di Milano  
Roma, 12.04.2016

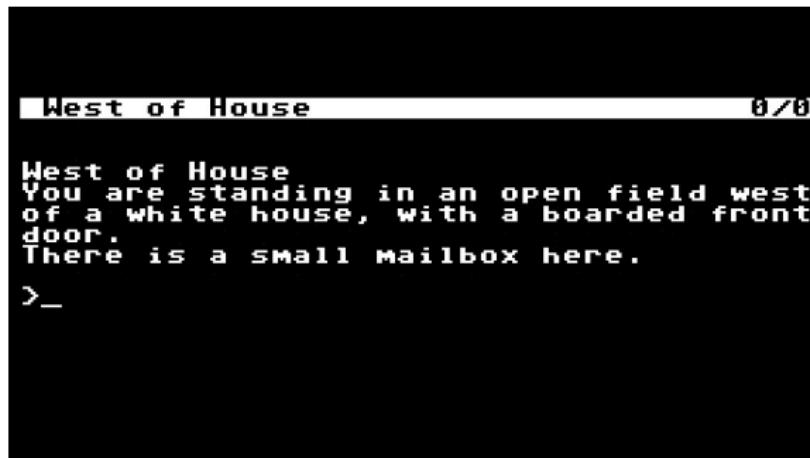
# Pong 1966



Le prime immagini di un videogioco di ping-pong appartengono al 1966 (Sanders Associates).

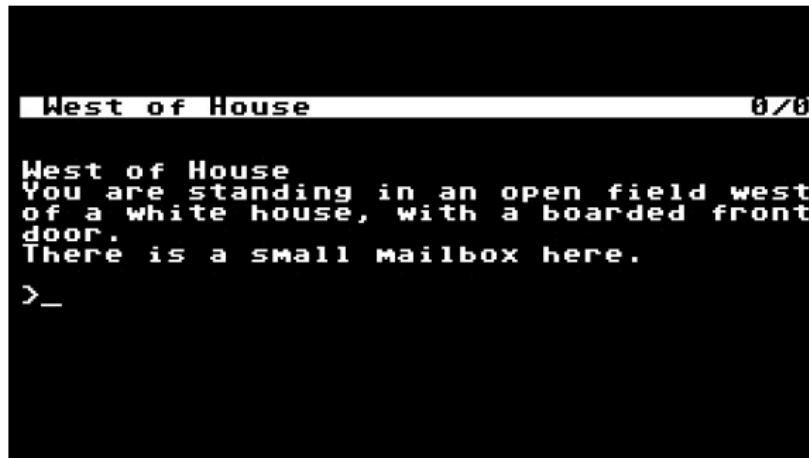
Nello stesso periodo Nolan Bushnell, neolaureato e futuro fondatore della *Atari*, realizza PONG. La versione domestica deve attendere il 1975.

# Zork 1977



Che ci crediate o no, questo è stato uno dei primi videogiochi di successo per personal computer! Zork è un'avventura testuale ideata da alcuni componenti del *Dynamic Modelling Group* (MIT) tra il 1977 e il 1979.

# Zork 1977



Che ci crediate o no, questo è stato uno dei primi videogiochi di successo per personal computer! Zork è un'avventura testuale ideata da alcuni componenti del *Dynamic Modelling Group* (MIT) tra il 1977 e il 1979.

dimenticavo . . . “open mailbox”.

# Wild Gunman 1974



Se si era abbastanza fortunati da vivere vicino a qualche sala giochi le cose non erano poi così male! *Nintendo* stava rivoluzionando la grafica!

# Wild Gunman 1974



Se si era abbastanza fortunati da vivere vicino a qualche sala giochi le cose non erano poi così male! *Nintendo* stava rivoluzionando la grafica!

*"si devono usare le mani, allora è un gioco da bambini"*

21.10.2015 — Ritorno al futuro II

# Super Mario 1985



Processori limitati ed altre restrizioni imposte dall'hardware delle console obbligavano all'uso di patterns ripetitivi.

Le nuvole e i cespugli sono la stessa immagine!

# De gustibus



Stessa storia, stesso posto . . . ma la grafica



GTA



Final  
Fantasy



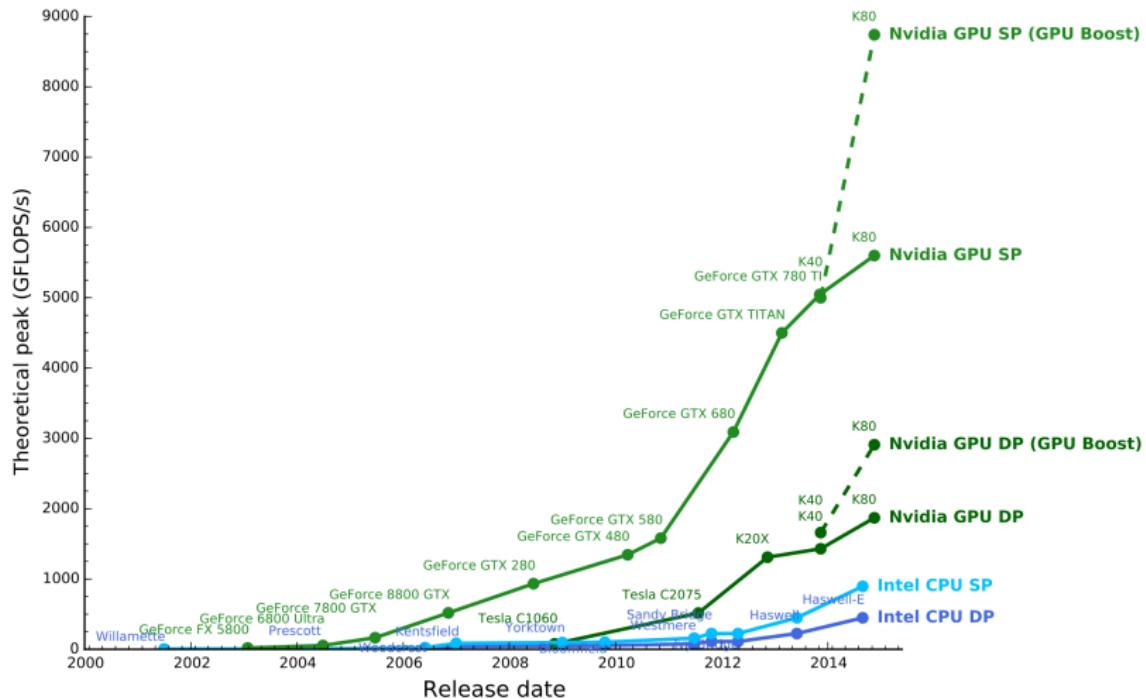
# Cosa ha permesso questa evoluzione?

## **Acceleratore hardware**

Un acceleratore è uno speciale componente hardware progettato per migliorare sensibilmente le prestazioni del computer in un settore specifico, tipicamente per le operazioni di calcolo in virgola mobile e le elaborazioni grafiche.

Gli acceleratori vengono spesso utilizzati nei videogiochi per aumentarne il livello di realismo grafico e la fluidità. Questo è possibile in quanto l'acceleratore hardware è stato progettato specificamente per le operazioni più spinte nel settore specifico a cui si rivolge.

# CPU vs. GPU



# Prima legge di Moore

*La complessità di un microcircuito, misurata ad esempio tramite il numero di transistori per chip, raddoppia ogni 18 mesi.*

## **Limite della prima legge di Moore**

La legge di Moore ha un limite imposto dalla fisica: esiste un limite fisico per la riduzione delle dimensioni dei transistor, al di sotto dei quali si genererebbero effetti indesiderati di natura quantistica nei circuiti elettronici.

Raggiunto quel limite (ci siamo quasi) l'unico modo possibile e praticabile per aumentare le prestazioni di calcolo è rappresentato dalla tecnologia multicore.

### **GPU NVIDIA vs. Intel Xeon Phi**

Tianhe-2 (1) e Stampede (7) utilizzano acceleratori Intel Xeon Phi.  
Titan (2) e Piz Daint (6) utilizzano acceleratori GPU NVIDIA.

### **Alcuni dati interessanti**

Ad oggi, un totale di 90 sistemi della lista TOP500 utilizzano acceleratori e/o coprocessori. A Novembre 2014 i sistemi erano 75.

52 di questi sistemi utilizzano GPU NVIDIA, 4 usano ATI Radeon e 35 sfruttano la tecnologia Intel MIC (Xeon Phi).

### **GPU NVIDIA & Intel Xeon Phi**

Quattro sistemi usano sia GPU Nvidia che Intel Xeon Phi.

# CPU vs. GPU



Le CPU sono dispositivi *general purpose*: in linea di principio devono essere in grado di eseguire *qualunque* operazione.

Le GPU sono dispositivi *special purpose*: sono ottimizzate per processare (tramite operazioni elementari) grandi quantità di dati in parallelo.

Diversi scopi . . . diversi design.

**CPU**: minimizzare la latenza, l'intervallo di tempo che impiega la CPU ad eseguire una particolare operazione.

**GPU**: massimizzare il throughput, la capacità di elaborazione dei dati.

Diversi scopi . . . diversi design.

**CPU**: minimizzare la latenza, l'intervallo di tempo che impiega la CPU ad eseguire una particolare operazione.

**GPU**: massimizzare il throughput, la capacità di elaborazione dei dati.

CPU



GPU



# CPU & GPU ⇒ CUDA

CUDA (Compute Unified Device Architecture) è un'architettura hardware per l'elaborazione parallela creata da NVIDIA.

CUDA è un modello di programmazione che permette di aumentare le prestazioni di calcolo sfruttando la potenza di calcolo delle GPU.

I linguaggi di programmazione disponibili nell'ambiente di sviluppo per CUDA, costituiscono delle estensioni dei linguaggi di programmazione più diffusi. Il principale è CUDA-C (C con estensioni NVIDIA), altri sono estensioni di Python, Fortran, Java e MATLAB.

In questo laboratorio utilizzeremo **esclusivamente** CUDA-C.

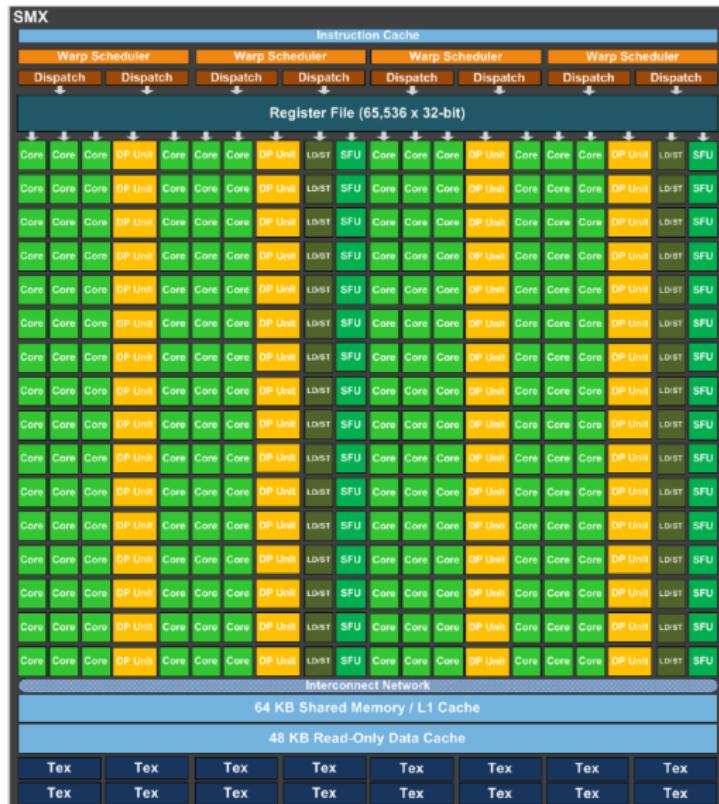
# NVIDIA GeForce GTX 980



NVIDIA Kepler GK110



# Architettura dello Streaming Multiprocessor (SMX)



# Single thread vs. Threads

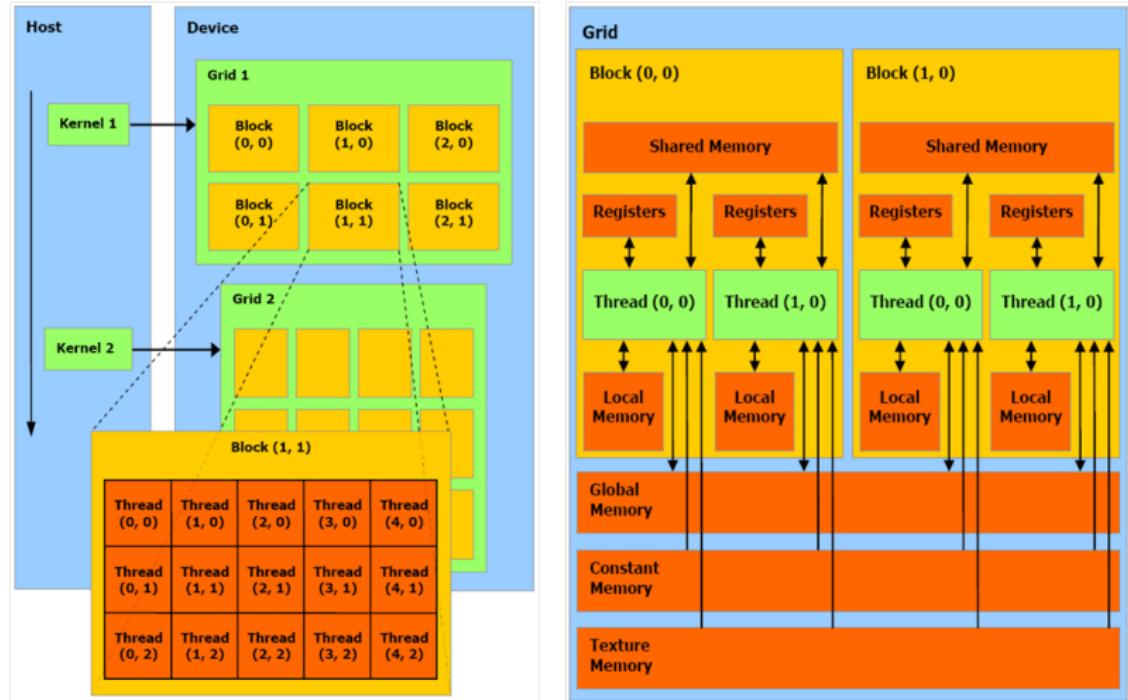
**CPU (single core)**: un solo processore.

**CPU (multi-core)**: più processori (2, 4, 8, 16).

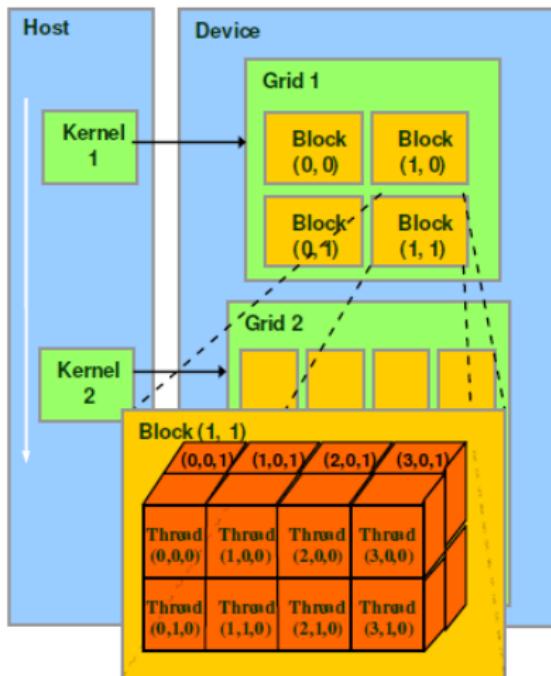
**CPU (cluster)**: più macchine.

**CPU + GPU**: la CPU, sfrutta i *microprocessori* della GPU.

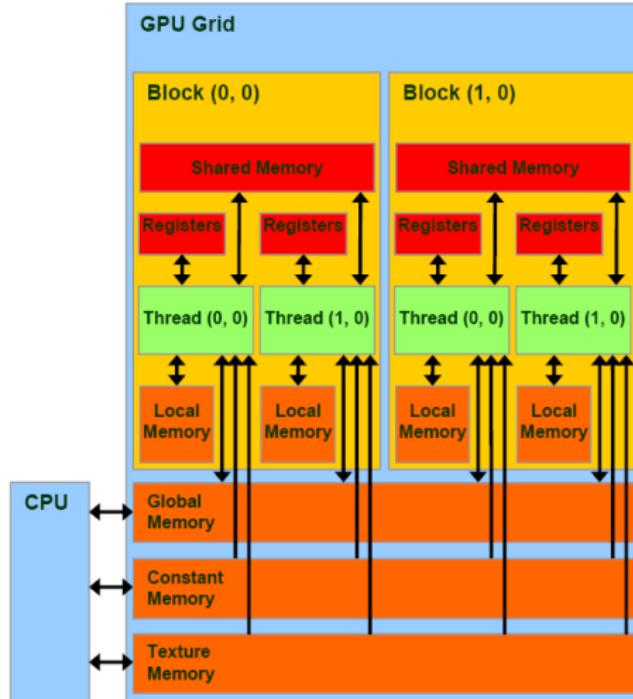
# Organizzazione dei threads



# Grid, Blocks e Threads



# Gestione delle memorie



# NVIDIA Kepler GK110

**13 Streaming Multiprocessor (SMX)**

CUDA Compute Capability **3.5**

## **SMX**

- 192 single-precision CUDA cores
- 64 double-precision units
- 32 special function units (SFU)
- 32 load/store units (LD/ST)

# NVIDIA Tesla K20c

- Name: Tesla K20c
- Compute capability: 3.5
- Global memory available on device in bytes: 5032706048
- Constant memory available on device in bytes: 65536
- Maximum pitch in bytes allowed by memory copies: 2147483647
- Shared memory available per block in bytes: 49152
- Size of L2 cache in bytes: 1310720
- 32-bit registers available per block: 65536

# NVIDIA Tesla K20c

- Clock frequency in kilohertz: 705500
- Number of multiprocessors on device: 13
- Maximum number of threads per block: 1024
- Maximum size of each dimension of a block: (1024, 1024, 64)
- Maximum size of each dimension of a grid: (2147483647, 65535, 65535)
- Warp size in threads: 32
- Maximum resident threads per multiprocessor: 2048

# CUDA-C: alcuni concetti di base

Qualificatori per le **routines**:

```
--global__ void routine_da_CPU_a_GPU() { }  
--device__ float routine_da_GPU_a_GPU() { }
```

Qualificatori per le **variabili**:

```
--constant__ float vettore_memoria_costant[32] ;  
--shared__ float vettore_memoria_shared[32] ;
```

**Configurare** l'esecuzione:

```
dim3 grid_dim(10, 5);  
dim3 block_dim(4, 8, 8);  
my_kernel <<< grid_dim, block_dim >>> ();
```

Variabili **built-in**:

```
dim3 blockDim;  
dim3 blockDim;  
dim3 blockIdx;  
dim3 threadIdx;  
void __syncthreads();
```

# Gestione delle variabili sul DEVICE (GPU)

Dichiarazione	Memoria	Visibilità	Durata
int var	registro	thread	thread
int var[10]	local	thread	thread
<code>__shared__ int var</code>	shared	block	block
<code>__global__ int var</code>	global	grid	programma

# Ok, ma come funziona?

Bene, ora siamo pronti a *sporcarci* le mani!

Vediamo come sfruttare la GPU per accelerare alcune semplici routine (*kernel*) tramite CUDA.

Le implementazioni presentate **NON** sono ottimali!

## my\_cuda.h

```
1 void HandleError(cudaError_t err, const char *file,
2     int line)
3 {
4     if (err != cudaSuccess)
5     {
6         printf("%s in %s at line %d\n",
7             cudaGetErrorString( err ), file, line);
8         exit( EXIT_FAILURE );
9     }
10 #define HANDLE_ERROR( err ) (HandleError( err ,
11     __FILE__ , __LINE__ ))
```

## somma\_CPU.cu |

```
1 #include <stdio.h>
2 #include <stdlib.h>
3
4 /* */
5 void sum(int x, int y, int *z);
6
7 /* ### */
8 int main()
9 {
10     int x = 17, y = 13, *z;
11
12     z = (int *) malloc(sizeof(int));
13
14     sum(x, y, z);
15
16     printf("%d + %d = %d\n", x, y, *z);
17     free(z);
18     return 0;
19 }
```

## somma\_CPU.cu II

```
20  /* ### */
21 void sum(int x, int y, int *z)
22 {
23     *z = x+y;
24     return;
25 }
```

# somma\_GPU.cu |

```
1 #include <stdio.h>
2 #include <cuda.h>
3 #include "my_cuda.h"
4
5 /* */
6 __global__ void sum(int x, int y, int *z);
7
8 /* ### */
9 int main()
10 {
11     int x = 17, y = 13, *z;
12     int *d_z;
13
14     z = (int *) malloc(sizeof(int));
15     HANDLE_ERROR(cudaMalloc((void**)&d_z, sizeof(int)));
16
17     sum<<<1, 1>>>(x, y, d_z);
18 }
```

## somma\_GPU.cu II

```
19 HANDLE_ERROR(cudaMemcpy(z, d_z, sizeof(int),
20                         cudaMemcpyDeviceToHost));
21
22 printf("%d + %d = %d\n", x, y, *z);
23
24 free(z);
25 cudaFree(d_z);
26
27 return (0);
28 */
29 __global__ void sum(int x, int y, int *z)
30 {
31     *z = x+y;
32     return;
33 }
```

## somma\_vettori\_CPU.cu |

```
1 #include <stdio.h>
2
3 #define SIZE 10000000
4
5 /* */
6 void add(double *x, double *y, double *z);
7
8 /* ### */
9 int main()
10 {
11     long int i, istep;
12     double *x, *y, *z;
13
14     x = (double *) malloc(SIZE*sizeof(double));
15     y = (double *) malloc(SIZE*sizeof(double));
16     z = (double *) malloc(SIZE*sizeof(double));
17
18
19
```

## somma\_vettori\_CPU.cu //

```
20  for(i = 0; i < SIZE; i++) {
21      x[i] = (double)rand() / RAND_MAX * 10;
22      y[i] = (double)rand() / RAND_MAX * 10;
23  }
24
25  add(x, y, z);
26
27  istep = SIZE / 10;
28  for(i = 0; i < SIZE; i += istep) {
29      printf("[%ld] %le + %le = %le\n", i, x[i], y[i], z[i]);
30  }
31
32  free(x);
33  free(y);
34  free(z);
35
36  return 0;
37 }
38 }
```

## somma\_vettori\_CPU.cu III

```
39 | /* ### */
40 | void add(double *x, double *y, double *z)
41 | {
42 |     int id = 0;
43 |
44 |     while(id < SIZE)  {
45 |         z[id] = x[id]+y[id];
46 |         id += 1;
47 |     }
48 |     return;
49 | }
```

## somma\_vettori\_GPU.cu |

```
1 #include <stdio.h>
2 #include <cuda.h>
3 #include "my_cuda.h"
4
5 #define SIZE 10000000
6
7 /* */
8 __global__ void add(double *x, double *y, double *z);
9
10 /* ### */
11 int main()
12 {
13     long int i, istep;
14     double *x, *y, *z;
15     double *d_x, *d_y, *d_z;
16
17     x = (double *) malloc(SIZE*sizeof(double));
18     y = (double *) malloc(SIZE*sizeof(double));
19     z = (double *) malloc(SIZE*sizeof(double));
```

## somma\_vettori\_GPU.cu //

```
20
21 HANDLE_ERROR(cudaMalloc((void**)&d_x, SIZE*sizeof(
22     double)));
22 HANDLE_ERROR(cudaMalloc((void**)&d_y, SIZE*sizeof(
23     double)));
23 HANDLE_ERROR(cudaMalloc((void**)&d_z, SIZE*sizeof(
24     double)));
24
25 for(i = 0; i < SIZE; i++) {
26     x[i] = (double)rand() / RAND_MAX * 10.;
27     y[i] = (double)rand() / RAND_MAX * 10.;
28 }
29
30 HANDLE_ERROR(cudaMemcpy(d_x, x, SIZE*sizeof(double),
31     cudaMemcpyHostToDevice));
31 HANDLE_ERROR(cudaMemcpy(d_y, y, SIZE*sizeof(double),
32     cudaMemcpyHostToDevice));
32
33 add<<<2048*10,1024>>>(d_x, d_y, d_z);
34
```

## somma\_vettori\_GPU.cu III

```
35     HANDLE_ERROR(cudaMemcpy(z, d_z, SIZE*sizeof(double),  
36                      cudaMemcpyDeviceToHost));  
  
37     istep = SIZE/10;  
38     for(i = 0; i < SIZE; i += istep) {  
39         printf("[%ld] %le + %le = %le\n", i, x[i], y[i], z  
        [i]);  
40    }  
  
41  
42    free(x);  
43    free(y);  
44    free(z);  
  
45  
46    cudaFree(d_x);  
47    cudaFree(d_y);  
48    cudaFree(d_z);  
  
49  
50    return(0);  
51}  
52}
```

## somma\_vettori\_GPU.cu IV

```
53  /* ### */
54  __global__ void add(double *x, double *y, double *z)
55  {
56      int id = blockIdx.x*blockDim.x+threadIdx.x;
57      while(id < SIZE)  {
58          z[id] = x[id]+y[id];
59          id += blockDim.x*gridDim.x;
60      }
61      return;
62 }
```

# prodotto\_scalare.cu |

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <cuda.h>
4 #include "my_cuda.h"
5
6 #define MIN(a,b) ((a)<(b)?(a):(b))
7 #define N 10000000
8 #define NBLOCKS 1024
9 #define NTHREADS 512
10
11 typedef double DOUBLE;
12
13 __global__ void scalar_product(DOUBLE *x, DOUBLE *y,
14                               DOUBLE *z);
15
16
17
18
```

## prodotto\_scalare.cu //

```
19  /* ### */
20 int main()
21 {
22     int i;
23     DOUBLE *x, *y, *ztmp, z, zcheck;
24     DOUBLE *d_x, *d_y, *d_ztmp;
25
26     x = (DOUBLE *) malloc(N*sizeof(DOUBLE));
27     y = (DOUBLE *) malloc(N*sizeof(DOUBLE));
28     ztmp = (DOUBLE *) malloc(NBLOCKS*sizeof(DOUBLE));
29
30     HANDLE_ERROR(cudaMalloc((void**)&d_x, N*sizeof(
31         DOUBLE)));
32     HANDLE_ERROR(cudaMalloc((void**)&d_y, N*sizeof(
33         DOUBLE)));
34     HANDLE_ERROR(cudaMalloc((void**)&d_ztmp, NBLOCKS*
35         sizeof(DOUBLE)));
```

## prodotto\_scalare.cu III

```
36     for(i = 0; i < N; ++i) {
37         x[i] = (DOUBLE) i;
38         y[i] = i*2.;
39     }
40
41     HANDLE_ERROR(cudaMemcpy(d_x, x, N*sizeof(DOUBLE),
42                             cudaMemcpyHostToDevice));
42     HANDLE_ERROR(cudaMemcpy(d_y, y, N*sizeof(DOUBLE),
43                             cudaMemcpyHostToDevice));
44
44     scalar_product<<<NBLOCKS, NTHREADS>>>(d_x, d_y,
45                                         d_ztmp);
46
46     HANDLE_ERROR(cudaMemcpy(ztmp, d_ztmp, NBLOCKS*sizeof
47                           (DOUBLE), cudaMemcpyDeviceToHost));
48
48     z = 0;
49     for(i = 0; i < NBLOCKS; ++i)
50         z += ztmp[i];
51
```

## prodotto\_scalare.cu IV

```
52     zcheck = (N-1.)*(N)*(2.*N-1)/3. ;
53     printf("GPU value %24.16le\nExpected value %24.16le\
54             n", z, zcheck);
55
56     HANDLE_ERROR(cudaFree(d_x));
57     HANDLE_ERROR(cudaFree(d_y));
58     HANDLE_ERROR(cudaFree(d_ztmp));
59
60     free(x);
61     free(y);
62     free(ztmp);
63
64
65
66
67
68
69
70 }
```

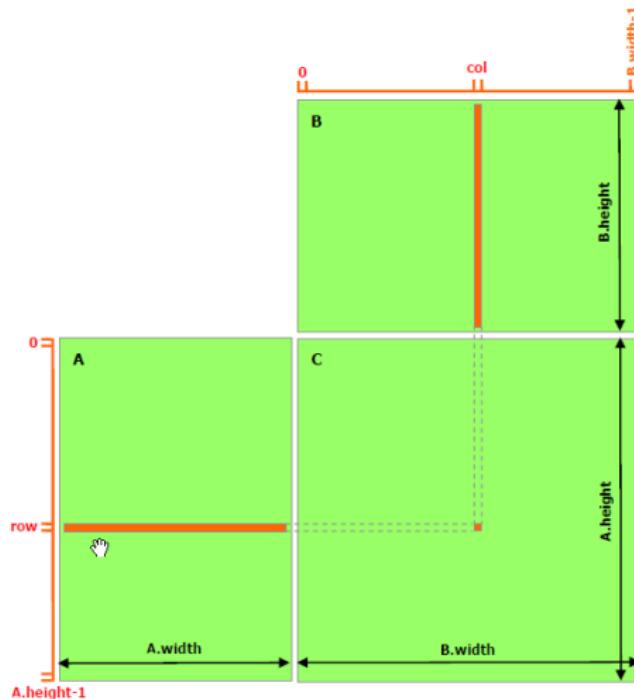
## prodotto\_scalare.cu V

```
71 /* ### */
72 __global__ void scalar_product(DOUBLE *x, DOUBLE *y,
73                                DOUBLE *z)
74 {
75     __shared__ float ztmp[NTHREADS];
76     int id, tid;
77     unsigned int j;
78     DOUBLE tmp;
79
80     id = threadIdx.x+blockIdx.x*blockDim.x;
81     tid = threadIdx.x;
82
83     tmp = 0.;
84     while(id < N) {
85         tmp += x[id]*y[id];
86         id += blockDim.x*gridDim.x;
87     }
88
89 }
```

## prodotto\_scalare.cu VI

```
90     ztmp[tid] = tmp;
91
92     __syncthreads();
93
94     j = blockDim.x;
95     while((j >= 1) != 0) {
96         if(tid < j)
97             ztmp[tid] += ztmp[tid+j];
98         __syncthreads();
99     }
100
101    if(tid == 0)
102        z[blockIdx.x] = ztmp[tid];
103
104    return;
105 }
```

# Prodotto di Matrici (global memory)



# prodotto\_matrici\_CPU-1.cu |

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4
5 typedef struct {
6     unsigned int h;
7     unsigned int w;
8     unsigned int offset;
9     double *el;
10 } MATRIX;
11
12 #define NTEST 10
13
14
15
16
17
18
19
```

## prodotto\_matrici\_CPU-1.cu II

```
20 /* ### */
21 int matrix_create_CPU(unsigned int w, unsigned int h,
22 MATRIX *a);
23 void matrix_destroy_CPU(MATRIX *a);
24 void matrix_randomfill_CPU(MATRIX *a);
25 int matrix_product_CPU(MATRIX *a, MATRIX *b, MATRIX *c
26 );
27 void matrix_export_octave(MATRIX *a, char *varname,
28 char *filename);
29 void matrix_product_octave(char *filename);
30 void matrix_clear_octave(char *filename);
31
32
33
34
35
36
```

## prodotto\_matrici\_CPU-1.cu III

```
37  /* ### */
38  int main()
39  {
40      MATRIX a, b, c;
41      unsigned int l, n;
42      clock_t start0, end0, delta0;
43      double msec;
44      char filename[256];
45      FILE *ofp;
46
47      sprintf(filename, "matrix_product_CPU-1_times.dat");
48      ofp = fopen(filename, "w");
49
50      for(l = 2; l <= 2048; l *= 2) {
51
52          if(matrix_create_CPU(l, l, &a) != 0) exit(-1);
53          if(matrix_create_CPU(l, l, &b) != 0) exit(-1);
54
55
56
```

## prodotto\_matrici\_CPU-1.cu IV

```
57     /*  */
58     delta0 = 0.;
59     for(n = 0; n < NTEST; ++n) {
60
61         sprintf(filename, "octave_matrixprodCPU1_%04u
62             -%02u.txt", l, n);
63
64         matrix_randomfill_CPU(&a);
65         matrix_randomfill_CPU(&b);
66
67         start0 = clock();
68         if(matrix_product_CPU(&a, &b, &c) != 0) exit(-1)
69     ;
70         end0 = clock();
71         delta0 += end0-start0;
72
73         if(l <= 100 && n == 0) {
74             matrix_clear_octave(filename);
75             matrix_export_octave(&a, "A", filename);
76             matrix_export_octave(&b, "B", filename);
```

# prodotto\_matrici\_CPU-1.cu V

```
75     matrix_export_octave(&c, "C", filename);
76     matrix_product_octave(filename);
77 }
78
79     matrix_destroy_CPU(&c);
80 }
81
82 matrix_destroy_CPU(&a);
83 matrix_destroy_CPU(&b);
84
85 msec = (double) delta0*1000./CLOCKS_PER_SEC/NTEST;
86
87 fprintf(ofp, "%08u\t%24.16le\n", l, msec);
88 }
89
90 fclose(ofp);
91
92 return 0;
93 }
```

# prodotto\_matrici\_CPU-1.cu VI

```
95  /* ### */
96 int matrix_create_CPU(unsigned int h, unsigned int w,
97   MATRIX *a)
98 {
99   size_t size;
100
101   a->h = h;
102   a->w = w;
103   a->offset = w;
104   size = h*w*sizeof(double);
105   a->el = malloc(size);
106
107   if(a->el == NULL) {
108     perror("ERROR: malloc failed (matrix_create_CPU)\n");
109     return -1;
110   }
111
112   return 0;
}
```

## prodotto\_matrici\_CPU-1.cu VII

```
113 /* ### */
114 void matrix_destroy_CPU(MATRIX *a)
115 {
116     a->h = 0;
117     a->w = 0;
118     a->offset = 0;
119     free(a->el);
120     a->el = NULL;
121     return;
122 }
123 /* ### */
124 void matrix_randomfill_CPU(MATRIX *a)
125 {
126     size_t i, size;
127
128     size = a->w*a->h;
129     for(i = 0; i < size; ++i)
130         a->el[i] = drand48();
131     return;
132 }
```

## prodotto\_matrici\_CPU-1.cu VIII

```
133 /* ### */
134 int matrix_product_CPU(MATRIX *a, MATRIX *b, MATRIX *c
135 )
136 {
137     unsigned int i, j, l;
138     double tmp;
139
140     if(a->w != b->h) {
141         perror("ERROR: incompatible matrices sizes (
142             matrix_product_CPU)\n");
143         return -1;
144     }
145
146     if(matrix_create_CPU(a->h, b->w, c) != 0) {
147         perror("ERROR: matrix creation for the product
148             failed (matrix_product_CPU)\n");
149         return -2;
150     }
```

# prodotto\_matrici\_CPU-1.cu IX

```
150     for(i = 0; i < c->h; ++i)
151         for(j = 0; j < c->w; ++j) {
152             tmp = 0.;
153             for(l = 0; l < a->w; ++l)
154                 tmp += a->el[i*a->w+l]*b->el[l*b->w+j];
155             c->el[i*c->w+j] = tmp;
156         }
157
158     return 0;
159 }
160 /* ### */
161 void matrix_export_octave(MATRIX *a, char *varname,
162                           char *filename)
163 {
164     int i, j;
165     FILE *ofp;
166
167     ofp = fopen(filename, "a");
168
169     fprintf(ofp, "%s", varname);
```

# prodotto\_matrici\_CPU-1.cu X

```
169     fprintf(ofp, " = [");
170
171     for(i = 0; i < a->h; ++i) {
172         for(j = 0; j < a->w; ++j)
173             fprintf(ofp, "%24.16le%s", a->el[i*a->w+j], (j ==
174             a->w-1)?";\n":",");
175     }
176     fprintf(ofp, "] ;\n");
177
178     fclose(ofp);
179
180     return;
181 }
182 /* ### */
183 void matrix_product_octave(char *filename)
184 {
185     FILE *ofp;
186
187     ofp = fopen(filename, "a");
188     fprintf(ofp, "C1 = A*B;\n");
```

# prodotto\_matrici\_CPU-1.cu XI

```
188     fprintf(ofp, "diff=C-C1;\n");
189     fclose(ofp);
190
191     return;
192 }
193 /* #### */
194 void matrix_clear_octave(char *filename)
195 {
196     FILE *ofp;
197
198     ofp = fopen(filename, "w");
199     fprintf(ofp, "# Created by matrix_product_CPU-1\n");
200     fclose(ofp);
201
202     return;
203 }
```

## prodotto\_matrici\_CPU-2.cu |

```
1  /* #### */
2  int matrix_transpose_CPU(MATRIX *a, MATRIX *at)
3  {
4      int i, j;
5
6      if(matrix_create_CPU(a->w, a->h, at) != 0) {
7          perror("ERROR: matrix creation for the transpose
8          failed (matrix_transpose_CPU)\n");
9          return -1;
10     }
11
12     for(i = 0; i < at->h; ++i)
13         for(j = 0; j < at->w; ++j)
14             at->el[i*at->w+j] = a->el[j*a->w+i];
15
16     return 0;
17 }
18 }
```

## prodotto\_matrici\_CPU-2.cu II

```
19  /* ### */
20  int matrix_product_CPU(MATRIX *a, MATRIX *b, MATRIX *c
21      )
22  {
23      unsigned int i, j, l;
24      MATRIX bt;
25      double tmp;
26
27      if(a->w != b->h) {
28          perror("ERROR: incompatible matrices sizes (
29                  matrix_product_CPU)\n");
30          return -1;
31      }
32
33      if(matrix_create_CPU(a->h, b->w, c) != 0) {
34          perror("ERROR: matrix creation for the product
35                  failed (matrix_product_CPU)\n");
36          return -2;
37      }
38
39      for(i = 0; i < a->h; i++) {
40          for(j = 0; j < b->w; j++) {
41              c->data[i][j] = 0;
42              for(l = 0; l < a->w; l++)
43                  c->data[i][j] += a->data[i][l] * b->data[l][j];
44          }
45      }
46
47      return 0;
48  }
```

## prodotto\_matrici\_CPU-2.cu III

```
36     if(matrix_transpose_CPU(b, &bt) != 0) {
37         perror("ERROR: matrix transpose for the product
38             failed (matrix_product_CPU)\n");
39         return -3;
40     }
41
42     for(i = 0; i < c->h; ++i)
43         for(j = 0; j < c->w; ++j) {
44             tmp = 0.;
45             for(l = 0; l < a->w; ++l)
46                 tmp += a->el[i*a->w+l]*bt.el[j*bt.w+l];
47                 c->el[i*c->w+j] = tmp;
48             }
49
50             matrix_destroy_CPU(&bt);
51
52             return 0;
53 }
```

# prodotto\_matrici\_GPU-1.cu |

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include "my_cuda.h"
5
6 typedef struct {
7     unsigned int h;
8     unsigned int w;
9     unsigned int offset;
10    double *el;
11 } MATRIX;
12
13 #define NTEST 10
14 #define NBLOCKS 1024
15 #define NTHREADS 256
16
17
18
19
```

## prodotto\_matrici\_GPU-1.cu II

```
20 /* ### */
21 int matrix_create_CPU(unsigned int w, unsigned int h,
22     MATRIX *a);
23 void matrix_destroy_CPU(MATRIX *a);
24 void matrix_randomfill_CPU(MATRIX *a);
25 int matrix_product_CPU(MATRIX *a, MATRIX *b, MATRIX *c
26 );
27 void matrix_export_octave(MATRIX *a, char *varname,
28     char *filename);
29 void matrix_product_octave(char *filename);
30 void matrix_clear_octave(char *filename);
31 void matrix_create_GPU(unsigned int h, unsigned int w,
32     MATRIX *d_a);
33 void matrix_destroy_GPU(MATRIX *d_a);
34 void matrix_copy_to_GPU(MATRIX *a, MATRIX *d_a);
35 void matrix_copy_from_GPU(MATRIX *d_a, MATRIX *a);
36 int matrix_product_GPU(MATRIX *a, MATRIX *b, MATRIX *c
37 );
38 __global__ void ker_matrix_product(MATRIX d_a, MATRIX
39     d_b, MATRIX d_c);
```

## prodotto\_matrici\_GPU-1.cu III

```
34  /* ### */
35  int main()
36  {
37      MATRIX a, b, c;
38      unsigned int l, n;
39      clock_t start0, end0, delta0;
40      double msec;
41      char filename[256];
42      FILE *ofp;
43
44      sprintf(filename, "matrix_product_GPU-1_times.dat");
45      ofp = fopen(filename, "w");
46
47      for(l = 2; l <= 2048; l *= 2) {
48
49          if(matrix_create_CPU(l, l, &a) != 0) exit(-1);
50          if(matrix_create_CPU(l, l, &b) != 0) exit(-1);
51
52
53 }
```

## prodotto\_matrici\_GPU-1.cu IV

```
54     /*  */
55     delta0 = 0.;
56     for(n = 0; n < NTEST; ++n) {
57
58         sprintf(filename, "octave_matrixprodGPU1_%04u
59             -%02u.txt", l, n);
60
61         matrix_randomfill_CPU(&a);
62         matrix_randomfill_CPU(&b);
63
64         start0 = clock();
65         if(matrix_product_GPU(&a, &b, &c) != 0) exit(-1)
66     ;
67         end0 = clock();
68         delta0 += end0-start0;
69
70         if(l <= 100 && n == 0) {
71             matrix_clear_octave(filename);
72             matrix_export_octave(&a, "A", filename);
73             matrix_export_octave(&b, "B", filename);
74         }
```

# prodotto\_matrici\_GPU-1.cu V

```
72 |     matrix_export_octave(&c, "C", filename);
73 |     matrix_product_octave(filename);
74 | }
75 |
76 |     matrix_destroy_CPU(&c);
77 | }
78 |
79 |     matrix_destroy_CPU(&a);
80 |     matrix_destroy_CPU(&b);
81 |
82 |     msec = (double) delta0*1000./CLOCKS_PER_SEC/NTEST;
83 |
84 |     fprintf(ofp, "%08u\t%24.16le\n", l, msec);
85 | }
86 |
87 | fclose(ofp);
88 |
89 | return 0;
90 | }
```

# prodotto\_matrici\_GPU-1.cu VI

```
92  /* ### */
93 void matrix_create_GPU(unsigned int h, unsigned int w,
94   MATRIX *d_a)
95 {
96   size_t size;
97
98   d_a->h = h;
99   d_a->w = w;
100  d_a->offset = w;
101  size = h*w*sizeof(double);
102  HANDLE_ERROR(cudaMalloc((void **) &(d_a->el), size));
103
104  return;
105
106
107
108
109
110 }
```

# prodotto\_matrici\_GPU-1.cu VII

```
111 /* ### */
112 void matrix_destroy_GPU(MATRIX *d_a)
113 {
114     d_a->h = 0;
115     d_a->w = 0;
116     HANDLE_ERROR(cudaFree(d_a->el));
117     d_a->el = NULL;
118 }
119
120
121
122
123
124
125
126
127
128
129
130
```

prodotto\_matrici\_GPU-1.cu VIII

```
131 /* ### */
132 void matrix_copy_to_GPU(MATRIX *a, MATRIX *d_a)
133 {
134     size_t size;
135
136     size = a->w*a->h*sizeof(double);
137     HANDLE_ERROR(cudaMemcpy(d_a->el, a->el, size,
138                           cudaMemcpyHostToDevice));
139     return;
140 }
141 /* ### */
142 void matrix_copy_from_GPU(MATRIX *d_a, MATRIX *a)
143 {
144     size_t size;
145
146     size = d_a->w*d_a->h*sizeof(double);
147     HANDLE_ERROR(cudaMemcpy(a->el, d_a->el, size,
148                           cudaMemcpyDeviceToHost));
149     return;
150 }
```

# prodotto\_matrici\_GPU-1.cu IX

```
149 /* ### */
150 int matrix_product_GPU(MATRIX *a, MATRIX *b, MATRIX *c
151 )
152 {
153     MATRIX d_a, d_b, d_c;
154
155     if(a->w != b->h) {
156         perror("ERROR: incompatible matrices sizes (
157             matrix_product_GPU)\n");
158         return -1;
159     }
160
161     if(matrix_create_CPU(a->h, b->w, c) != 0) {
162         perror("ERROR: matrix creation for the product
163             failed (matrix_product_GPU)\n");
164         return -2;
165     }
166
167     matrix_create_GPU(a->h, a->w, &d_a);
168     matrix_create_GPU(b->h, b->w, &d_b);
```

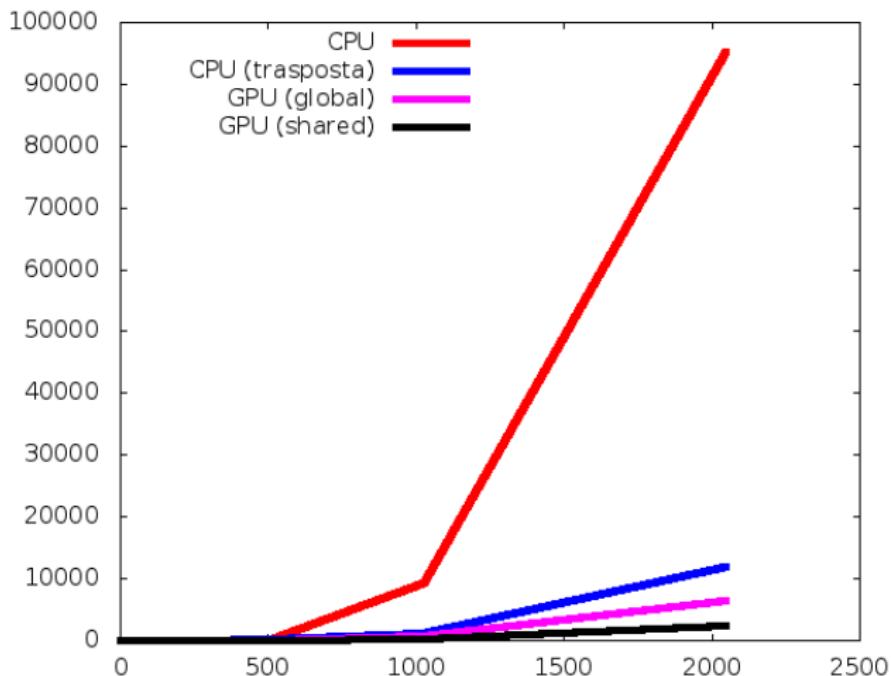
# prodotto\_matrici\_GPU-1.cu X

```
166     matrix_create_GPU(c->h, c->w, &d_c);  
167  
168     matrix_copy_to_GPU(a, &d_a);  
169     matrix_copy_to_GPU(b, &d_b);  
170  
171     ker_matrix_product<<< NBLOCKS , NTHREADS >>>(d_a, d_b  
172     , d_c);  
173  
174     matrix_destroy_GPU(&d_a);  
175     matrix_destroy_GPU(&d_b);  
176  
177     matrix_copy_from_GPU(&d_c, c);  
178  
179     matrix_destroy_GPU(&d_c);  
180  
181     return 0;  
182 }  
183  
184 }
```

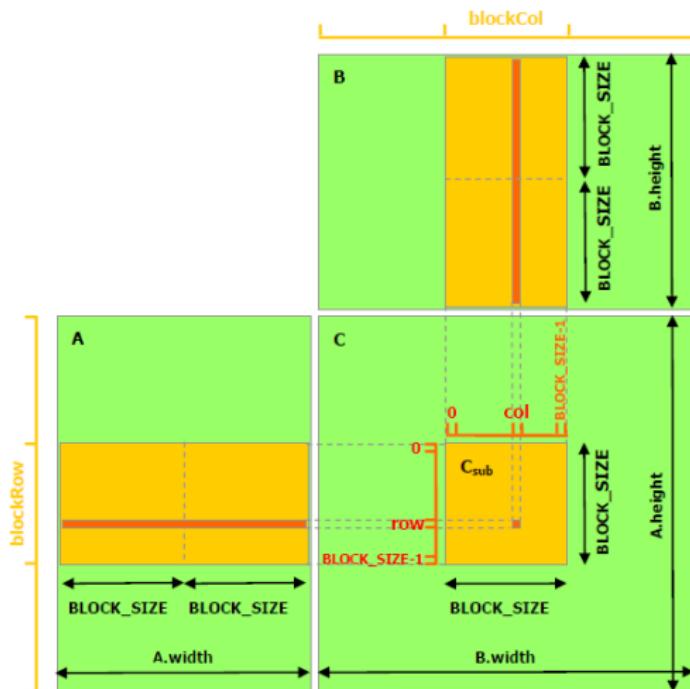
# prodotto\_matrici\_GPU-1.cu XI

```
185 /* ### */
186 __global__ void ker_matrix_product(MATRIX d_a, MATRIX
187     d_b, MATRIX d_c)
188 {
189     double tmp;
190     unsigned int idx, i, j, k, size;
191
192     idx = blockIdx.x*blockDim.x+threadIdx.x;
193     size = d_c.h*d_c.w;
194     while(idx < size) {
195         i = idx/d_c.w;    j = idx%d_c.w;
196         tmp = 0.;
197         for(k = 0; k < d_a.w; ++k)
198             tmp += d_a.el[i*d_a.w+k] * d_b.el[k*d_b.w+j];
199         d_c.el[idx] = tmp;
200         idx += blockDim.x*gridDim.x;
201     }
202 }
```

# Prodotto di Matrici (CPU vs. GPU)



# Prodotto di Matrici (shared memory)



# prodotto\_matrici\_GPU-1.cu |

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <time.h>
4 #include "my_cuda.h"
5
6 typedef struct {
7     unsigned int h;
8     unsigned int w;
9     unsigned int offset;
10    double *el;
11 } MATRIX;
12
13 #define NTEST 10
14 #define NBLOCKS 1024
15 #define NTHREADS 16
16
17
18
19
```

## prodotto\_matrici\_GPU-1.cu II

```
20
21  /* ### */
22  __device__ double get_elem(const MATRIX d_a, int row,
23      int col) {
24      return d_a.el[row*d_a.offset+col];
25  }
26  /* ### */
27  __device__ void set_elem(MATRIX d_a, int row, int col,
28      double value)
29  {
30      d_a.el[row*d_a.offset+col] = value;
31
32
33
34
35
36
37
```

# prodotto\_matrici\_GPU-1.cu III

```
38 /* ### */
39 __device__ MATRIX get_submatrix(MATRIX d_a, int row,
40     int col)
41 {
42     MATRIX d_sa;
43     d_sa.w = NTHREADS;
44     d_sa.h = NTHREADS;
45     d_sa.offset = d_a.offset;
46     d_sa.el = d_a.el+(d_a.offset*NTHREADS*row+NTHREADS*
47         col);
48     return d_sa;
49 }
50
51
52
53
54
55
```

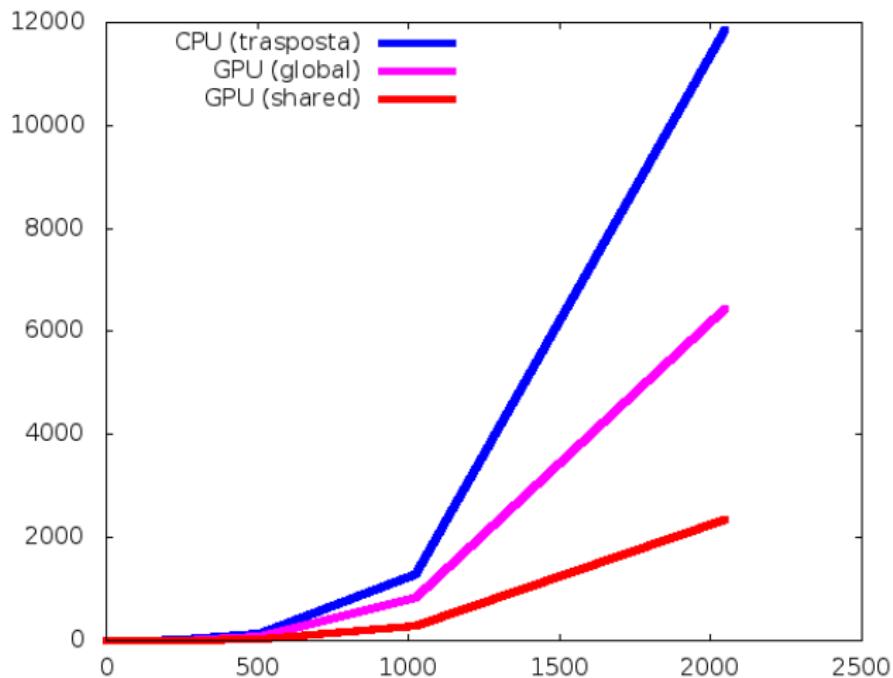
## prodotto\_matrici\_GPU-1.cu IV

```
56 /* ### */
57 __global__ void ker_matrix_product(MATRIX d_a, MATRIX
58   d_b, MATRIX d_c)
59 {
60   MATRIX d_asub, d_bsub, d_csub;
61   unsigned int i, j, row, col, blkrow, blkcol;
62   double cvalue;
63   __shared__ double asub_el[NTHREADS][NTHREADS];
64   __shared__ double bsub_el[NTHREADS][NTHREADS];
65
66   row = threadIdx.y;
67   col = threadIdx.x;
68
69   blkrow = blockIdx.y;
70   blkcol = blockIdx.x;
71
72   d_csub = get_submatrix(d_c, blkrow, blkcol);
73   cvalue = 0.;
```

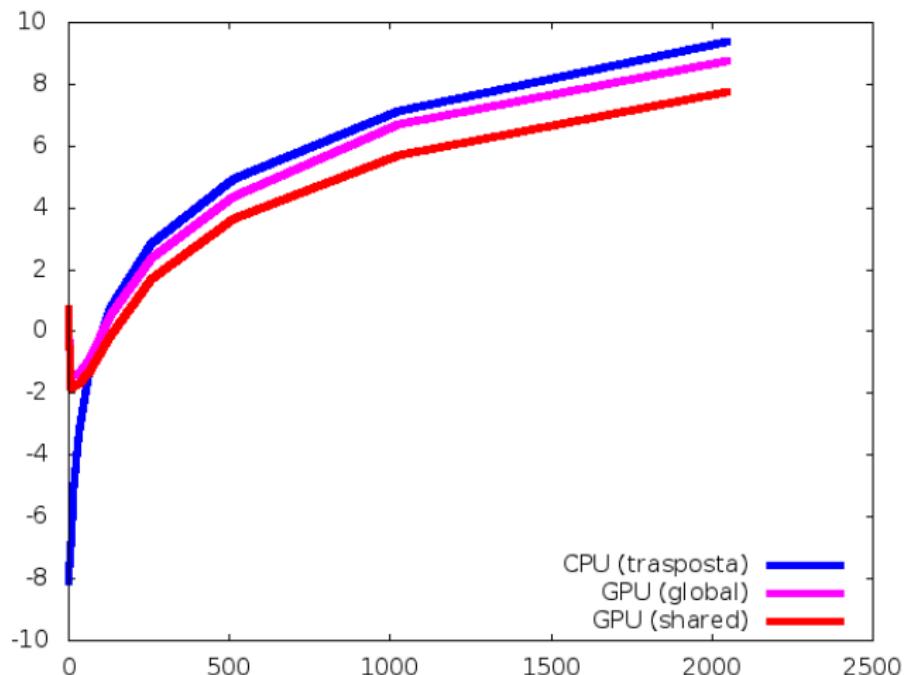
## prodotto\_matrici\_GPU-1.cu V

```
75     for (i = 0; i < d_a.w/NTHREADS; ++i) {
76         d_asub = get_submatrix(d_a, blkrow, i);
77         d_bsub = get_submatrix(d_b, i, blkcol);
78
79         asub_el[row][col] = get_elem(d_asub, row, col);
80         bsub_el[row][col] = get_elem(d_bsub, row, col);
81
82         __syncthreads();
83
84         for(j = 0; j < NTHREADS; ++j)
85             cvalue += asub_el[row][j]*bsub_el[j][col];
86         __syncthreads();
87     }
88
89     set_elem(d_csub, row, col, cvalue);
90
91     return;
92 }
```

# Prodotto di Matrici (CPU vs. global vs. shared)



# Prodotto di Matrici (CPU vs. global vs. shared)



# Non si vive di soli conti...

The screenshot shows a web browser displaying the [OpenSubdiv](http://graphics.pixar.com/opensubdiv/docs/intro.html) documentation. The URL in the address bar is `graphics.pixar.com/opensubdiv/docs/intro.html`. The page title is "Introduction". The top navigation bar includes links for "User Docs", "API Docs", "Release Notes", and "Forum", along with a "GITHUB" button. On the left, there is a sidebar with a search bar and a list of navigation links:

- > Introduction
- > License
- > Getting Started
- > Contributing
- > Building OpenSubdiv
- > Code Examples
- > Roadmap
- > References
- > Release 3.0
  - > Overview
  - > Porting Guide: 2.0 to 3.0
  - > Subdivision Compatibility
- > Subdivision Surfaces
  - > Introduction
  - > Topology
  - > Uniform
  - > Feature Adaptive
  - > Boundary Interpolation
  - > Face-Varying Interpolation
  - > Semi-Sharp Creases
  - > Modeling Tips
- > OpenSubdiv User Guide
  - > API Overview
  - > Sdc
  - > Vtr
  - > Far
    - > Topology Refiner
    - > Topology Refiner Factory
    - > Primvar Refiner
    - > Patch Table
    - > Stencil Table
  - > Osd
  - > Shader Interface
- > Tutorials
- > Historical But Relevant
  - > Hbr
  - > Using Hbr

---

# Grazie per l'attenzione!

*Domande?*

*Commenti?*

<http://www.mat.unimi.it/users/sansotte/cuda/>

---