

# Sparse Matrix Transposition for GPUs

Massimiliano Incudini VR433300  
Michele Penzo VR439232

**Sommario**—L'obiettivo principale di questo progetto è stato quello di implementare alcune metodologie proposte per effettuare *Sparse Matrix Transposition* su *Gpu*. Sono stati analizzati alcuni algoritmi, descritti in sezione III, partendo dall'algoritmo seriale, passando a cuSPARSE per finire con l'implementazione degli algoritmi descritti in [1]. Infine vengono esposti i risultati, tratte le conclusioni e proposti alcuni possibili miglioramenti.

## I. INTRODUZIONE E MOTIVAZIONI

Sempre più applicazioni computazionali in ambito scientifico necessitano di algoritmi che compiano operazioni applicabili su matrici sparse. Si parla di semplici operazioni di algebra lineare, di moltiplicazione o di calcolo della trasposta come in questo caso.

Il problema analizzato, quello della trasposizione di matrici, si presta bene al calcolo parallelo per l'esecuzione in maniera più efficiente e veloce. Verranno quindi mostrate le basi per la rappresentazione, i problemi riscontrati durante lo sviluppo e analizzati alcuni algoritmi per il calcolo su GPU.

## II. RAPPRESENTAZIONE DELLE MATRICI

Una matrice sparsa non è altro che una matrice i cui valori sono per la maggior parte uguali a zero. La matrice in formato classico necessita di una quantità di memoria minima di  $m \times n$  elementi, ma essendo l'obiettivo quello di lavorare su matrici sparse non è stato necessario e utile memorizzare la matrice in formato denso.

Per rappresentare in modo efficace le matrici sparse senza troppo utilizzo di memoria sono state quindi introdotte ed utilizzate delle forme di rappresentazione matriciale che permettono il salvataggio di dati utilizzando quantitativi di memoria inferiori.

Di seguito vengono spiegate le due metodologie da noi utilizzate.

### A. Formato Csr

Il *compressed sparse row* è una rappresentazione di una matrice  $M$  basata su tre array monodimensionali, che rispettivamente contengono:

- 1)  $V$ : i valori non zero ( $nnz$ ),
- 2)  $COL\_INDEX$ : gli indici delle colonne dove si trovano gli elementi  $nnz$ ,
- 3)  $ROW\_INDEX$ : ha un elemento per ogni riga della matrice e rappresenta l'indice in  $V$  dove comincia la riga data.

I primi due array sono di dimensione  $nnz$ , mentre il terzo array è al massimo di dimensione  $m$ .

### B. Formato Csc

Questa metodologia per la rappresentazione è simile alla precedente citata *Csr*, a differenza che i valori vengono letti prima per colonna. Di conseguenza, un indice di riga viene memorizzato per ogni valore e lo stesso viene fatto per i puntatori di colonna.

### C. Da Csr a Csc

Per il problema della trasposta di matrice è stato quindi utile introdurre entrambe le rappresentazioni. Infatti, ogni algoritmo descritto in sezione III, necessita di sei array per effettuare il calcolo della trasposta e dare l'output nella tipologia corretta. Abbiamo quindi:

- in input il formato *Csr*:  $csrRowPtr$ ,  $csrColIdx$ ,  $csrVal$ ;
- in output il formato *Csc*:  $cscColPtr$ ,  $cscRowIdx$ ,  $cscVal$ .

In base a come vengono create le matrici, se in modo casuale oppure se lette da file, vengono effettuate delle operazioni preliminari descritte dalla procedura in sezione IV che portano ad ottenere gli array in input e in output nel formato corretto per effettuarne il controllo di correttezza.

## III. METODOLOGIE ANALIZZATE

Gli algoritmi paralleli che calcolano la trasposta di una matrice sparsa descritti in [1] vengono confrontati nelle tempistiche rispetto all'algoritmo seriale e rispetto alle funzioni di libreria cuSPARSE fornite da NVidia come parte della loro piattaforma Cuda SDK.

Per convenzione la matrice è in formato *csr*. I nomi delle componenti della matrice sono chiamati  $csr\_row\_ptr$  (lunghezza  $m + 1$ ) per il vettore dei puntatori ad inizio riga,  $csr\_col\_idx$  (lunghezza  $nnz$ ) per il vettore degli indici di colonna e  $csr\_val$  (lunghezza  $nnz$ ) per il vettore dei valori.

### A. Trasposta seriale

L'algoritmo si sviluppa nel seguente modo:

- 1) si applica la funzione *istogramma*  $csr\_col\_idx$  che calcola le occorrenze di ogni colonna, applicando *scan* si ottiene il vettore  $csc\_col\_ptr$  (lunghezza  $n + 1$ ) che conterrà i puntatori agli elementi di inizio riga trasposta;
- 2) allochiamo il vettore  $csc\_row\_idx$  contenente gli indici di riga (lunghezza  $nnz$ );
- 3) allochiamo il vettore  $csc\_val$  contenente gli elementi non nulli (lunghezza  $nnz$ );
- 4) per ogni riga  $i \in [0, m)$  processiamo gli elementi corrispondenti all' $i$ -esima riga, legati alle posizioni  $j \in [csr\_row\_ptr[i], csr\_row\_ptr[i + 1])$

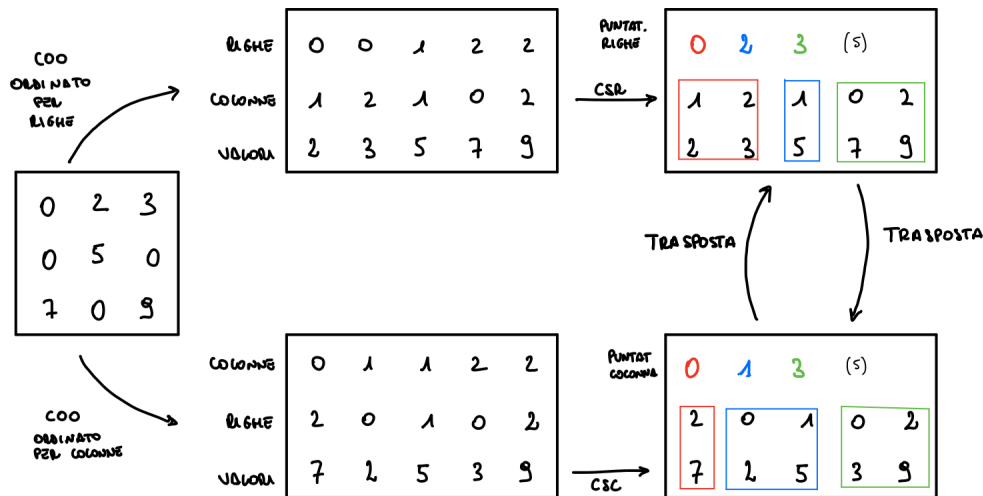


Figura 1. Trasformazione da formato esteso a CSR, oppure CSC

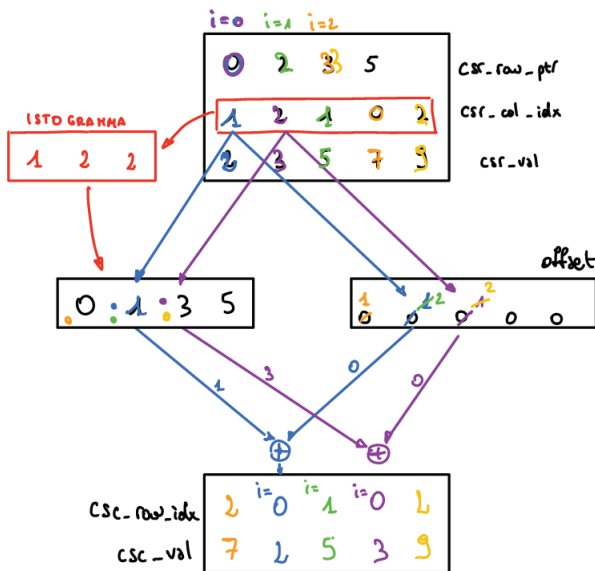


Figura 2. Algoritmo seriale

- la locazione loc del  $j$ -esimo elemento ordinata per righe è  $csc\_col\_ptr[j] + offset[j]$  dove  $offset[j]$  è un contatore incrementato ogni volta che aggiungiamo un elemento della colonna  $csc\_col\_ptr[j]$ ;
- l'indice di riga dell'elemento loc-esimo è  $i$ ;
- il valore dell'elemento loc-esimo è  $csr\_val[j]$ .

L'esecuzione dell'algoritmo seriale è illustrato in Figura 2.

### B. ScanTrans

L'algoritmo si sviluppa nel seguente modo:

- il vettore  $csr\_row\_ptr$  viene espanso per ottenere il vettore  $csr\_row\_idx$  attraverso la procedura *pointers-to-indexes* che effettua l'operazione inversa

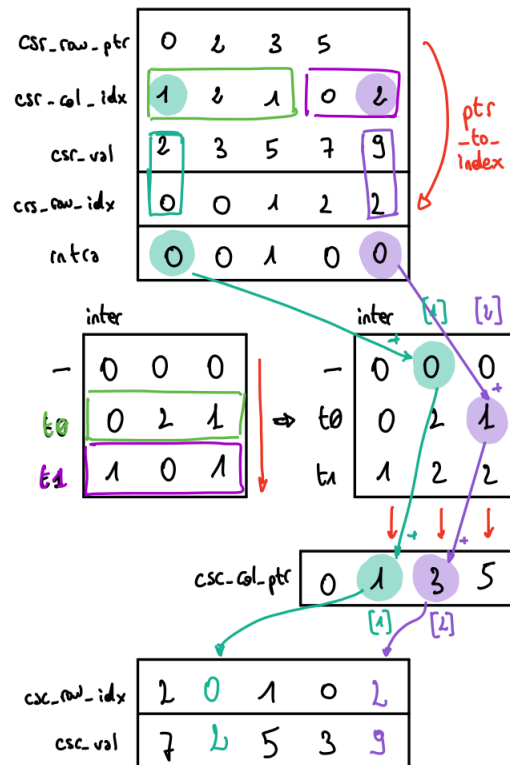


Figura 3. Algoritmo ScanTrans

dell'istogramma;

ora la struttura ( $csr\_row\_idx$ ,  $csr\_col\_idx$ ,  $csr\_val$ ) rappresenta la matrice in formato COO ordinato per righe;

- decidiamo un numero arbitrario di thread  $K$  che lavoreranno ciascuna su blocchi di  $K/nnz$  elementi. Allochiamo una matrice *inter* di dimensioni  $(K+1) \times n$  ed un vettore *intra* di dimensione  $nnz$ . La prima manterrà nella  $(i+1)$ -esima riga l'istogramma parziale relativo

all' $i$ -esimo blocco. La seconda mantiene un offset di colonna di ciascun elemento;

- 3) i  $K$  thread eseguono l'*istogramma* ciascuna sul proprio blocco, la matrice *inter* viene riempita;
- 4) viene applicata l'operazione *scan* ad ogni colonna della matrice *inter*;
- 5) il vettore *csc\_col\_ptr* viene ottenuto copiando l'ultima riga di *inter* ed applicando al vettore risultante l'operazione *scan*;
- 6) riordiniamo gli elementi nei vettori *csr\_row\_idx*, *csr\_val* ottenendo quindi *csc\_row\_idx*, *csc\_val*; la nuova posizione dell'elemento  $i$ -esimo viene calcolata nel seguente modo:

- $b = i / \text{nnz}$ ;
- $c = \text{csc\_col\_idx}[i]$ ;
- $l = \text{csc\_col\_ptr}[c] + \text{inter}[bn + c] + \text{inter}[c]$ ;
- $\text{csc\_val}[l] \leftarrow \text{csr\_val}[j]$ ;
- $\text{csc\_row\_idx}[l] \leftarrow \text{csr\_row\_idx}[j]$ .

Il funzionamento dell'algoritmo è illustrato in Figura 3.

### C. MergeTrans

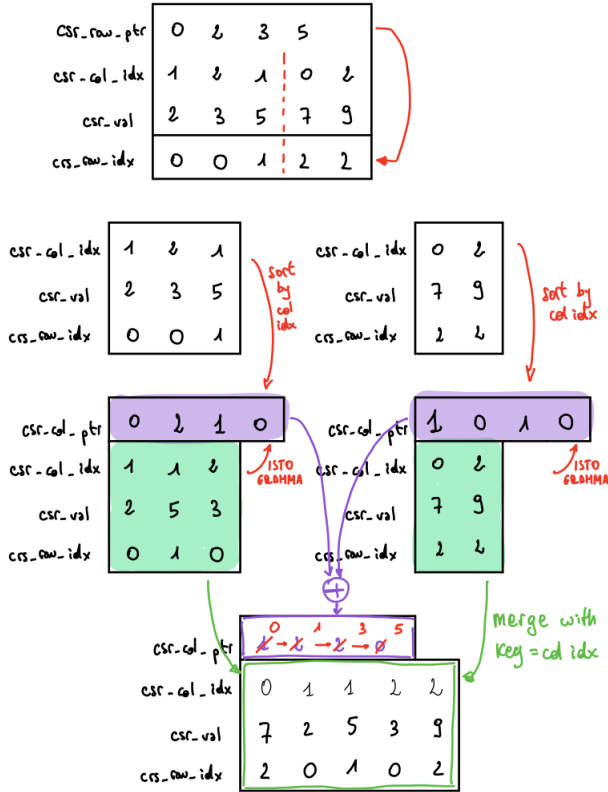


Figura 4. Algoritmo MergeTrans

Il precedente algoritmo esegue nello step (6) un insieme di accessi casuali in memoria che non permettono di sfruttare in modo adeguato la shared memory.

Per mitigare questo aspetto negativo viene introdotta una variante di questo algoritmo che, se implementato correttamente, permette ai thread di lavorare sempre con un numero costante, contiguo di elementi.

L'algoritmo considerato prevede due passi importanti: *sort* e *merge*.

Inizialmente sono stati creati gli indici di riga a partire dai puntatori delle colonne e su questi ultimi è stato fatto un sort su piccole porzioni di array, mantenendo quindi i vari blocchi disordinati tra di loro ma con gli elementi ordinati.

Successivamente è stato utilizzato il merge ricorsivo partendo dai blocchi più piccoli e unendoli in blocchi sempre più grandi. Per funzionare questo processo necessita dell'utilizzo di due buffer di memoria che contengono gli elementi appena ordinati.

Infine dai puntatori delle colonne vengono estratti gli indici e viene fatta la scan che ritorna il risultato in formato *csc*.

Il funzionamento dell'algoritmo è illustrato in Figura 4.

### D. NVidia cuSPARSE

I risultati temporali di tutti gli algoritmi precedentemente descritti saranno confrontati con l'implementazione di NVidia della trasposta della matrice sparsa.

La libreria cuSPARSE è parte di CUDA Toolkit, insieme di librerie utilizzate per effettuare operazioni tra vettori e matrici attraverso diversi formati.

Versioni diverse di CUDA Toolkit implementano funzioni con nomi diversi e funzionamenti diversi. In particolare per Cuda Toolkit 9.0 possiamo trasporre la matrice sparsa attraverso la procedura *cusparscScsr2csc*.

Nella versione Cuda Toolkit 10.2 viene esposto il metodo *cusparscCsr2cscEx2* che effettua la trasposta in due possibili modi a seconda se il parametro "algoritmo" è valorizzato con la costante *CUSPARSE\_CSR2CSC\_ALG1* oppure *CUSPARSE\_CSR2CSC\_ALG2*. Le due versioni differiscono per le tempistiche e per il consumo di memoria (maggiore memoria necessaria per la seconda implementazione). Le tempistiche sono riportate in Tabella 15.

La libreria è dettagliatamente documentata in [2]. A seconda della versione di *CUDA* dobbiamo modificare lo script di compilazione. In particolare il *Makefile* del progetto permette di compilare per due architetture specifiche, una con CUDA 9 e l'altra CUDA 10/11.

## IV. PROCEDURE

I due algoritmi ScanTrans e MergeTrans vengono scomposti in diversi componenti, ognuno dei quali viene valutato nelle performance e testato separatamente.

### A. Scan

Questa operazione prende in input un vettore  $A = (a_0, a_1, \dots, a_n)$  e ritorna un vettore  $B = (I, a_0, a_0 \oplus a_1, \dots, a_0 \oplus a_1 \oplus \dots \oplus a_{n-1})$  con  $\oplus$  un'operazione binaria il cui elemento identità è  $I$ . Nel nostro caso l'operazione è la somma.

L'algoritmo apparentemente sembra difficile da parallelizzare in quanto il risultato di ogni elemento dipende da tutti i gli elementi precedenti. Diverse soluzioni sono state proposte tra cui l'*algoritmo di Blelloch*. Il suo funzionamento in due fasi è illustrato in Figura 5 e dettagliatamente discusso in [3].

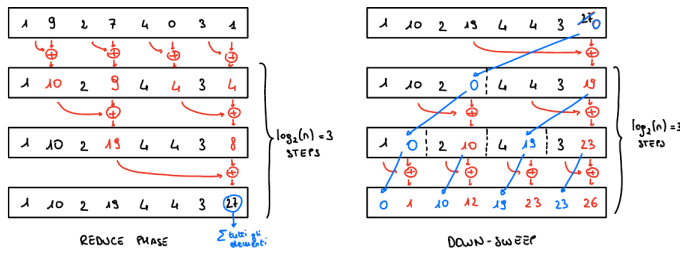


Figura 5. Algoritmo di Blelloch

L'implementazione prevede che se l'intero vettore riesce ad essere memorizzato all'interno della shared memory di  $N$  elementi, allora possiamo calcolare scan con una singola chiamata a kernel.

Nel caso questo non sia possibile, l'operazione di scan viene segmentata, applicata separatamente a blocchi di  $N$  elementi. Successivamente si mantiene un vettore di somme (vettore degli ultimi elementi del blocco), si applica ricorsivamente *scan* su esso e si sommano gli offset ottenuti all'intero vettore di partenza.

### B. Segmented sort

Questa operazione prende in input un vettore di lunghezza  $n$  ed un intero `BLOCK_SIZE`. Il vettore viene diviso in segmenti di lunghezza `BLOCK_SIZE`. Gli elementi di ogni segmento vengono permutati in modo che siano ordinati stabilmente. L'intero segmento deve rientrare nella shared memory.

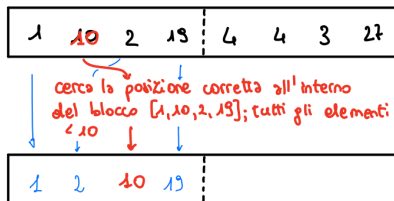


Figura 6. Segmented Sort

Una volta caricato il blocco in shared memory, la  $i$ -esima thread del blocco è incaricata di trovare la posizione corretta dell' $i$ -esimo elemento all'interno del segmento, ed assegnarlo a tale posizione. L'algoritmo viene illustrato in Figura 6.

L'ordinamento deve essere stabile quindi la posizione dell' $i$ -esimo elemento di valore  $y$  è dato dal numero di elementi  $< y$ , sommati al numero degli elementi  $= y$  per indici  $< i$ .

La dimensione ideale del blocco pari a 128 elementi (caso interi a 32bit), ed è stata trovata empiricamente:

Thread per blocco	Performance (ms)
128	2135.58
256	2139.62
512	2156.61

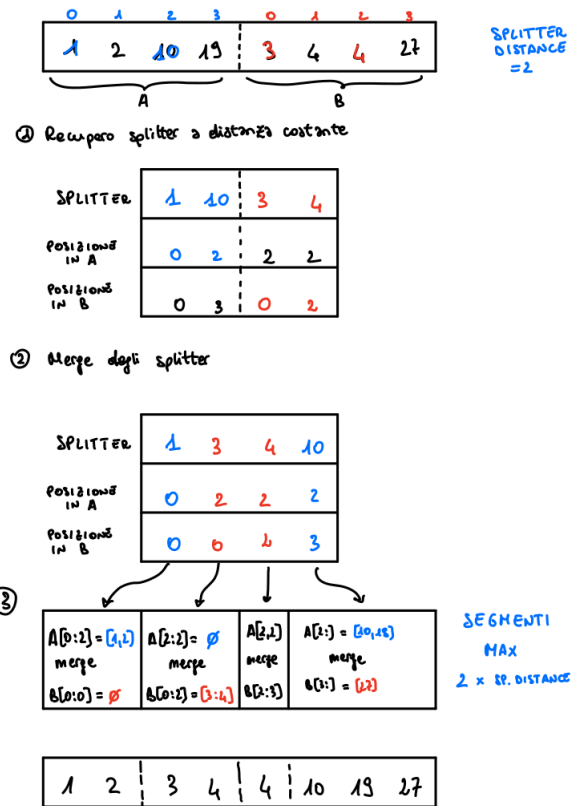
Figura 7. Performance su array di  $2 \cdot 10^7$  elementi

Figura 8. Merge big

### C. Merge

L'operazione di *merge* trasforma un vettore diviso in segmenti di dimensione `BLOCK_SIZE` nel quale gli elementi di ogni segmento sono ordinati, in un vettore diviso in segmenti ordinati di dimensione  $2 * \text{BLOCK\_SIZE}$ , ognuno dei quali è l'unione di una coppia di blocchi contigui.

Differenziamo il caso in cui il blocco di dimensione `BLOCK_SIZE` rientri o meno nella shared memory.

### D. Merge small

In questo caso una coppia di blocchi rientra completamente nella shared memory. In modo analogo a quanto fatto per l'operazione di *segmented sort*, abbiamo un numero di thread per blocco pari a `BLOCK_SIZE` nel quale l' $i$ -esimo thread è incaricato di calcolare la posizione dell' $i$ -esimo elemento.

In questo caso la posizione dell' $i$ -esimo elemento del blocco di sinistra è  $i + j$  con  $j$  posizione dell'elemento all'interno del blocco di destra, trovato attraverso una ricerca binaria in quanto i blocchi sono ordinati (differentemente da quanto avviene per il *segmented sort*). Analogamente lo stesso avviene per gli elementi del blocco di destra.

A parità di valore, gli elementi del blocco di sinistra hanno indice minore di quelli di destra.

### E. Merge big

Questo secondo algoritmo è illustrato in Figura 8 ed originariamente preso da [4]. Applicando sempre *merge small*, oltre che rinunciare alla shared memory, ogni thread del blocco

dovrebbe lavorare su più di un elemento, aumentando la complessità della procedura.

L'algoritmo proposto funziona nel seguente modo:

- 1) recuperiamo dal vettore degli elementi segnastop detti *splitter* presi a distanza costante e pari ad  $SP\_DIST$ ; per ogni coppia di segmenti da unire ho una coppia di blocchi di splitter, trovo quindi la posizione di ogni splitter all'interno del blocco di destra ( $A$ ) e di sinistra ( $B$ );
- 2) applico *merge* ricorsivamente sugli array di splitter;
- 3) gli indici associati agli splitter dividono la coppia di blocchi in modo tale da poter effettuare tanti merge indipendenti da quanti sono gli splitter. Ogni *merge* indipendente considererà al massimo  $2 * SP\_DIST$  elementi (numero costante che scegliamo tale che rientri nella shared memory).

La dimensione ideale del blocco pari a 256 elementi (caso interi a 32bit), ed è stata trovata empiricamente:

Thread per blocco	Performance (ms)
128	708.87
256	701.54
512	715.71

Figura 9. Performance su array di  $2 \cdot 10^7$  elementi

Mentre il valore migliore per la distanza degli splitter ( $SP\_DIST$ ) è di 128.

#### F. Istogramma / Index to pointers

Dato un vettore  $A$  di  $n$  elementi compresi tra 0 ed  $m - 1$  riceviamo un vettore di  $m$  elementi nel quale l' $i$ -esima cella contiene la frequenza con cui il valore  $i$  è presente in  $A$ .

L'algoritmo si sviluppa in due fasi:

- 1) il vettore  $A$  viene diviso in  $N$  segmenti di lunghezza omogenea, ogni segmento viene processato da un blocco di thread che mantiene l'istogramma parziale;
- 2) gli istogrammi parziali vengono poi uniti attraverso un'operazione di prefix scan che si effettua "in verticale". Tale operazione può essere ottenuta:
  - trasponendo la matrice degli istogrammi parziali ed applicando  $N$  volte *prefix sum*; oppure
  - attraverso  $N$  blocchi di thread che attraversano sequenzialmente il vettore colonna degli istogrammi parziali.
- 3) si applica *prefix sum* sul risultato dell'operazione precedente.

L'algoritmo è illustrato in Figura 10.

La dimensione ideale del blocco pari a 32 elementi (caso interi a 32bit), ed è stata trovata empiricamente:

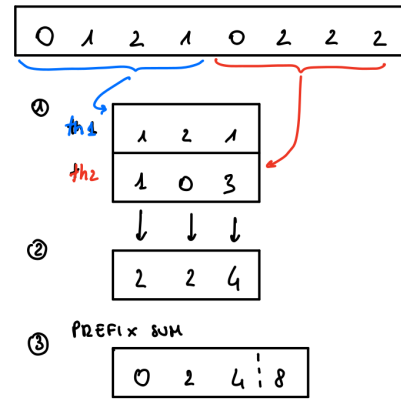


Figura 10. Index to pointers

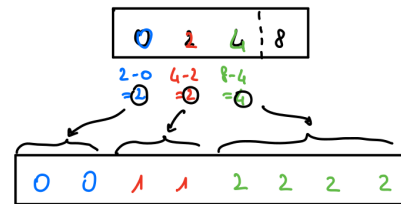


Figura 12. Pointers to index

Thread per blocco	Performance (ms)
1	1016.87
32	1013.65
64	1024.70
256	986.84

Figura 11. Performance su array di  $10^6$  elementi

#### G. Pointers to index

Operazione inversa di *index to pointers*. Il vettore risultante è ordinato. Può essere implementato assegnando ad ogni blocco di thread un elemento del vettore delle frequenze da espandere.

L'algoritmo è illustrato in Figura 12.

La dimensione ideale del blocco pari a 32 elementi (caso interi a 32bit), ed è stata trovata empiricamente:

Thread per blocco	Performance (ms)
1	2328.88
16	2375.74
32	2322.75
1024	3326.13

Figura 13. Performance su array di  $2 \cdot 10^7$  elementi

## V. STRUTTURA DELL'IMPLEMENTAZIONE

L'intera implementazione è scaricabile attraverso *git* dalla repository <https://github.com/michelepenzo/architetture-avanzate>.

La struttura delle directory del progetto è presente in Figura 14.

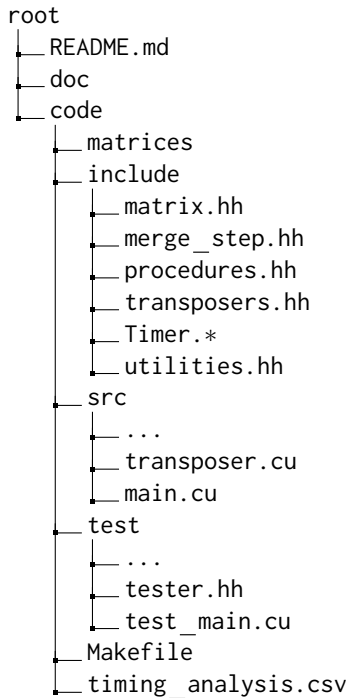


Figura 14. Struttura delle directory del progetto

La sottodirectory *doc* contiene questo stesso documento in formato *pdf* ed i rispettivi sorgenti *tex*.

La sottodirectory *code* contiene i sorgenti dell'applicativo principale e di quello secondario di test.

Lo scopo del primo è generare un file *csv* contenente le tempistiche e gli speedup di ogni algoritmo applicato sulle istanze di matrici in input generate casualmente oppure lette da file *mtx* (*market matrix*, una rappresentazione della matrice sparsa in formato COO attraverso file di testo).

Lo scopo del secondo applicativo è testare il corretto funzionamento di ogni componente del progetto. In particolare, sottopongo le stesse istanze di array o matrici (a seconda del componente che sto testando) sia alla funzione che ne implementa l'algoritmo parallelo, sia alla funzione che ne implementa l'algoritmo seriale. Mi aspetto che i risultati siano uguali per tutte le istanze.

In particolare:

- il file `include/matrix.hh` contiene le classi *FullMatrix* e *SparseMatrix* che si occupano di allocare nella memoria host lo spazio necessario a contenere la matrice date le sue specifiche ( $m$ ,  $n$ ,  $nnz$ ), sia come matrice estesa sia in formato *csr*. Inoltre contiene i metodi per inizializzare la matrice in modo casuale e per passare da un formato all'altro;
- il file `include/utills.hh` contiene i metodi di utilità quali le funzioni di stampa e di allocazione e deallocazione della memoria device;
- i file `include/procedures.hh` e `merge_step.hh` contengono le dichiarazioni delle procedure descritte nella Sezione IV. La maggior parte delle definizioni sono presenti nella sottodirectory *src*, nel caso del metodo *merge\_step* la definizione è scritta direttamente nell'header.

Questa scelta è conveniente in quanto tale funzione è definita rispetto ad un tipo generico, se la definizione fosse stata riportata nei file *cpp* avremmo dovuto indicare esplicitamente i tipi concreti per il quale vogliamo rendere disponibile il nostro metodo ([5]);

- il file `include/transposers.hh` e rispettivo sorgente `src/transposers.cu` contengono le dichiarazioni e definizioni dei metodi che effettuano la trasposta: seriale, parallela con *ScanTranse MergeTran* ed infine da libreria *cuSPARSE* con entrambi i possibili algoritmi;
- i file `Timer.*` contengono una classe di utilità *timer* che permette di cronometrare il tempo che occorre per eseguire un dato pezzo di codice;
- il file `src/main.cu` contiene l'applicativo principale che chiama i diversi metodi di trasposta sulla matrice fornita in input, ne cronometra le tempistiche e le stampa in output;
- il file `test/tester.hh` contiene la classe astratta *tester* che espone un metodo *test\_many\_instances* che chiama un metodo astratto *test* un numero arbitrario di volte, passandogli in input un intero che rappresenta il numero dell'istanza attuale che può essere usato per decidere la dimensione dell'istanza di test. Attualmente le istanze testate vanno da 1 a 20'000, poi da 20'000 a 20'000'000 con step  $\times 1.5$ ;
- i file `test/tester_*.hh` si occupano di testare un singolo componente, contengono ciascuna uno o più classi concrete che estendono la classe astratta *tester*;
- il file `test/test_main.cu` contiene l'applicativo di test che alloca oggetti delle varie classi *tester*, li avvia e ne stampa gli eventuali errori a video.

## VI. AVVIO DEGLI APPLICATIVI

L'applicativo principale può essere avviato con tre modalità diverse:

- senza parametri, genera una matrice  $500'000 \times 500'000$  con  $10'000'000$  elementi non nulli, ne valuta le tempistiche con i diversi algoritmi ritornando la media su un numero arbitrario di esecuzioni;
- con tre parametri interi  $m$ ,  $n$ ,  $nnz$ , genera una matrice avente le dimensioni ricevute in input e procede alla valutazione delle tempistiche come sopra;
- con un parametro stringa *filename*, legge da file una matrice che deve essere nel formato *mtx market matrix*.

Avviando l'applicativo principale attraverso il Makefile con *make run* viene avviato molteplici volte l'applicativo principale, ogni volta con un'istanza di matrice diversa in input, con lo scopo di popolare un documento *timing\_analysis.csv* contenente le tempistiche medie su diversi input.

L'applicativo di test si avvia con una sola modalità equivalentemente avviando il nome dell'applicativo senza parametro oppure attraverso il Makefile con *make test*. Su *stdout* viene stampato "no" se almeno un test ha mostrato anomalie, "ok" altrimenti.

## VII. RISULTATI SPERIMENTALI

Confrontiamo ora le performance delle varie implementazioni che seguono:

Nome	M	N	NNZ	Serial	ScanTrans	MergeTrans	cuSPARSE1	cuSPARSE2
language.mtx	399 130	399 130	1 216 334	55.75	114.71	197.05	122.02	18.36
webbase-1M.mtx	1 000 005	1 000 005	3 105 536	138.84	278.55	520.90	149.81	51.00
rajat21.mtx	411 676	411 676	1 893 370	77.97	147.64	306.56	127.31	26.77
ASIC_680k.mtx	682 862	682 862	3 871 773	154.37	265.47	844.76	155.31	56.74
memchip.mtx	2 707 524	2 707 524	14 810 202	594.18	994.62	2328.67	298.93	188.07
cant.mtx	62 451	62 451	2 034 917	73.83	114.40	248.49	127.26	31.11
FullChip.mtx	2 987 012	2 987 012	26 621 990	997.37	1543.02	9481.42	454.53	328.32
stomach.mtx	213 360	213 360	3 021 648	110.84	164.30	387.78	139.57	40.26
web-Google.mtx	916 428	916 428	5 105 039	399.63	382.37	3327.56	170.59	73.25
random	100 000	100 000	10 000 000	898.99	475.17	3056.12	210.57	130.25
random	100 000	100 000	10 000 000	902.23	475.72	3060.25	208.22	133.22
random	150 000	200 000	5 000 000	523.66	263.90	1351.83	161.90	72.70
random	150 000	200 000	5 000 000	527.15	262.88	1353.93	165.19	74.98
random	500 000	500 000	10 000 000	1380.96	532.76	2853.55	227.33	141.18

Figura 15. Risultati sperimentali - M, N, NNZ rispettivamente numero di righe, di colonne, di elementi non nulli della matrice. I tempi sono in ms.

Nome	M	N	NNZ	Serial	ScanTrans	MergeTrans	cuSPARSE1	cuSPARSE2
language.mtx	399 130	399 130	1 216 334	1.00	0.49	0.28	0.46	3.04
webbase-1M.mtx	1 000 005	1 000 005	3 105 536	1.00	0.50	0.27	0.93	2.72
rajat21.mtx	411 676	411 676	1 893 370	1.00	0.53	0.25	0.61	2.91
ASIC_680k.mtx	682 862	682 862	3 871 773	1.00	0.58	0.18	0.99	2.72
memchip.mtx	2 707 524	2 707 524	14 810 202	1.00	0.60	0.26	1.99	3.16
cant.mtx	62 451	62 451	2 034 917	1.00	0.65	0.30	0.58	2.37
FullChip.mtx	2 987 012	2 987 012	26 621 990	1.00	0.65	0.11	2.19	3.04
stomach.mtx	213 360	213 360	3 021 648	1.00	0.67	0.29	0.79	2.75
web-Google.mtx	916 428	916 428	5 105 039	1.00	1.05	0.12	2.34	5.46
random	100 000	100 000	10 000 000	1.00	1.89	0.29	4.27	6.90
random	100 000	100 000	10 000 000	1.00	1.90	0.29	4.33	6.77
random	150 000	200 000	5 000 000	1.00	1.98	0.39	3.23	7.20
random	150 000	200 000	5 000 000	1.00	2.01	0.39	3.19	7.03
random	500 000	500 000	10 000 000	1.00	2.59	0.48	6.07	9.78

Figura 16. Risultati sperimentali - Speedup

- seriale;
- parallela *scan trans*;
- parallela *merge trans*;
- cuSPARSE(entrambi gli algoritmi).

Le istanze su cui vengono eseguiti i vari algoritmi sono in parte generate in modo casuale (a partire dalle specifiche della matrice sparsa), in parte recuperate dal dataset "University of Florida sparse matrix collection" [6]. Tale dataset è stato usato per valutare le performance degli algoritmi in [1].

La macchina sul quale vengono eseguiti i vari algoritmi è equipaggiata con una scheda NVidia GeForce GTX 780 con Cuda Runtime 10.2.

I risultati sono visibili in Tabella 15 (tempistiche) e in Tabella 16 (speedup).

Analizzando i risultati notiamo che per istanze di dimensioni notevoli l'implementazione ScanTrans raggiunge speedup fino a  $\times 2.5$  rispetto all'algoritmo seriale, mentre le due implementazioni fornite da cuSPARSE raggiungono  $\times 6.0$  e  $\times 9.8$  di speedup.

L'implementazione MergeTrans risulta sempre meno efficiente della versione seriale dell'algoritmo. I possibili problemi relativi a questa implementazione sono discussi in Sezione VIII.

## VIII. CONSIDERAZIONI FINALI

Contrariamente a quanto asserito in [1], nel nostro caso ScanTrans si comporta meglio di MergeTrans. Proponiamo alcune idee per migliorare il risultato:

- utilizzare un algoritmo più efficiente per realizzare il merge quale il *Merge Path* ([7]);
- la procedura da noi descritte potrebbero essere sostituite con quelle fornite nella libreria *modern gpu* presente su Github; inoltre seguendo la descrizione delle implementazioni presenti in [8] ci avrebbe aiutato nella fase di ingegnerizzazione.
- la procedura *istogramma* (*index\_to\_pointers*) può essere migliorata nel seguente modo: attualmente una griglia di un unico thread calcola l'istogramma parziale sulla sua porzione di vettore; si potrebbero introdurre tante thread ognuna che legge l'intera porzione di blocco ma si occupa di processare gli elementi appartenenti ad una sola (o ad un piccolo insieme) di valori. Questo permetterebbe anche di caricare gli elementi in shared memory, eventualmente una porzione di blocco per volta.

## RIFERIMENTI BIBLIOGRAFICI

- [1] H. Wang, W. Liu, and K. Hou, "Parallel transposition of sparse data structures," *ICS '16*, 2016.

- [2] [Online]. Available: <https://docs.nvidia.com/cuda/cusparse/index.html#abstract>
- [3] M. Harris and J. D. Owens, "Chapter 39. parallel prefix sum (scan) with cuda." [Online]. Available: <https://developer.nvidia.com/gpugems/gpugems3/part-vi-gpu-computing/chapter-39-parallel-prefix-sum-scan-cuda>
- [4] "Bad to have one big merge." [Online]. Available: [https://www.youtube.com/watch?v=wx8LJIdJ\\_K8](https://www.youtube.com/watch?v=wx8LJIdJ_K8)
- [5] "Explicit template instantiation - when is it used?" [Online]. Available: <https://stackoverflow.com/questions/2351148/explicit-template-instantiation-when-is-it-used>
- [6] Y. H. Timothy A. Davis, "The university of florida sparse matrix collection," 2011. [Online]. Available: <https://sparse.tamu.edu/>
- [7] D. A. B. Oded Green, Robert McColl, "Gpu merge path - a gpu merging algorithm." [Online]. Available: <https://web.cs.ucdavis.edu/~amenta/f15/GPUmp.pdf>
- [8] [Online]. Available: <https://github.com/moderngpu/moderngpu>