

Progetto relativo al corso di Programmazione su
architetture parallele
a.a. 2021-2022

Gianluca Macrì
mat. 142213

Indice

1	Problema affrontato	2
2	Metodologia adottata per la soluzione	4
3	Risultati sperimentali	9
4	Osservazioni finali	11
A	Miniguia all'utilizzo	12
B	Tabelle dati	13

1 Problema affrontato

Lo scopo di questo lavoro è stato quello di implementare mediante il linguaggio di programmazione CUDA una versione parallela dell'*Auction algorithm*, originariamente proposto da Bertsekas (1979), per la risoluzione del problema del massimo matching su grafi. In particolare, analogamente all'approccio seguito da Vasconcelos e Rosenhahn (2009), ci si limiterà al caso semplificato dei grafi bipartiti per il quale è sempre possibile ottenere una soluzione ottimale, quando questa esista.

Per grafo bipartito si intende un grafo $G = (V, E)$ con $X, Y \subseteq V$ partizione dei nodi di V , ossia $X \cup Y = V$ e $X \cap Y = \emptyset$, e con archi che collegano esclusivamente nodi appartenenti a partizioni differenti, i.e. $E \subseteq X \times Y$. Un matching, o accoppiamento, M è quindi definito come un sottoinsieme di archi $M \subseteq E$ tale che per ciascun nodo $v_1 \in V$ esiste al più un arco in M incidente in esso, ossia $\forall v_1 \in V ((\exists v_2 \in V : (v_1, v_2) \in M) \implies \neg (\exists v_3 \in V : v_3 \neq v_2 (v_1, v_3) \in M))$.

Quando si parla di *matching problem*, si tende tipicamente ad associare un'interpretazione semantica ai due insiemi X ed Y , ad esempio nel caso dell'*Auction algorithm* indicandoli rispettivamente come partecipanti ad un asta e lotti (od oggetti) in vendita. In tale scenario, inoltre, si associano agli archi dei pesi a_{ij} per $(i, j) \in E$ che, corrispondono al valore che ciascuna persona attribuisce soggettivamente ai diversi lotti, e si cerca di associare ad ogni persona esattamente un oggetto in modo tale da massimizzare la funzione di profitto $\sum_{(i,j) \in M} a_{ij}$. Chiamamente, detti m ed n rispettivamente in numero di persone (o la cardinalità di X) e il numero di oggetti (o la cardinalità di Y), dovrà valere $m \leq n$ affinché una soluzione possa esistere.

Per risolvere tale problema l'*Auction algorithm* prevede quindi di "simulare" lo svolgimento di un'asta nella quale a ciascun oggetto j viene associato un prezzo p_j (inizialmente 0) e le persone competono al rialzo per aggiudicarsi il lotto da cui poter ottenere il miglior profitto. In particolare l'algoritmo procede per round successivi, ciascuno suddiviso in due fasi distinte: la fase di offerta (o *bidding phase*) e quella di assegnamento (*assignment phase*). All'inizio della k -ma iterazione, ciascuna persona i non ancora assegnata ad alcun oggetto confronterà i prezzi correnti p_j con i valori che associa a ciascun lotto e proporrà un incremento $\delta(i)$ per l'oggetto j_i che ritiene più conveniente.

Formalmente, detti

$$\begin{aligned} N(i) &= \{j \mid (i, j) \in E\} \\ j_i^k &= \operatorname{argmax}_{j \in N(i)} (a_{ij} - p_j^k) \\ g_1^k(i) &= \max_{j \in N(i)} (a_{ij} - p_j^k) = a_{ij_i^k} - p_{j_i^k}^k \\ g_2^k(i) &= \begin{cases} \max_{j \in N(i), j \neq j_i^k} (a_{ij} - p_j^k) & \text{se } |N(i)| > 1 \\ -\infty & \text{altrimenti} \end{cases} \end{aligned}$$

dove l'apice k indica il round, definiamo $\delta^k(i) = g_1^k(i) - g_2^k(i) + \epsilon$ come l'incremento di offerta sull'oggetto j_i^k che la persona i propone all'iterazione k -ma. Detto in altre parole la persona i al round k propone un rialzo sul prezzo dell'oggetto per cui avrebbe attualmente il maggior margine di guadagno, pari alla differenza tra i due maggiori guadagni possibili, più un valore ϵ positivo per evitare di avere incrementi nulli nel caso di pareggi.

Derivati tali valori, si conclude la fase di offerta e inizia quella di assegnamento in cui ogni oggetto viene assegnato al miglior offerente del round. In particolare, definendo $P^k(j)$ come l'insieme degli offerenti per j al round k , ogni oggetto j per il quale vale $P^k(j) \neq \emptyset$ è assegnato all'offerente $i_j^k = \operatorname{argmax}_{i \in P^k(j), j_i^k = j} \delta^k(i)$, lasciando disaccoppiato l'eventuale possessore precedente, e il prezzo di j viene incrementato diventando $p^{k+1}(j) = p^k + \delta^k(i_j)$. L'algoritmo si conclude non appena tutte le persone risultino assegnate ad un oggetto.

Bertsekas (1979) dimostra che, nel caso in cui una soluzione esista, la terminazione di tale procedura avviene in un numero finito di passi grazie a due elementi. Per prima cosa il fatto che ogni oggetto, una volta che viene assegnato ad una qualche persona, non può più tornare ad essere libero, facendo quindi crescere in maniera monotona il numero di persone accoppiate. A ciò si unisce il costante incrementare dei prezzi per gli oggetti che ricevono almeno un'offerta, aspetto che porta gli offerenti a considerarli via via meno vantaggiosi. Inoltre, si può verificare che nel caso in cui $a_{ij} \in \mathbb{N}$, utilizzare un valore $\epsilon < \frac{1}{\min(m,n)}$ dove $m = |X|$ e $n = |Y|$, permette all'algoritmo di raggiungere una soluzione ottimale.

Nel frammento 1 si propone lo pseudocodice per l'*auction algorithm* in versione seriale, in cui si utilizzano una matrice A per rappresentare i pesi degli archi tra persone ed oggetti, utilizzando $-\infty$ per gli archi inesistenti, e diversi vettori per rappresentare i costrutti necessari. Rispetto a questi ultimi, si noti come, a differenza della descrizione precedente dove si utilizzava δ^k per indicare gli incrementi proposti dalle diverse persone all'iterazione k -ma, nello pseudocodice si è preferito mantenere tale informazione catalogata per oggetto nel vettore Δ_p , così da velocizzare la fase di assegnamento.

Algorithm 1 Auction algorithm, versione seriale

```

1: function AUCTION( $A, m, n, \epsilon$ ) ▷  $m$  è il numero di persone,  $n$  quello degli oggetti
2:    $peopleMatched \leftarrow 0, matching \leftarrow \mathbf{0}$ 
3:    $p \leftarrow \mathbf{0}$  ▷ prezzi per oggetto
4:    $\Delta_p \leftarrow \mathbf{0}$  ▷ incrementi dei prezzi per oggetto
5:    $roundBidders \leftarrow \mathbf{0}$  ▷ migliori offerenti per oggetto nel round
6:   while  $peopleMatched < m$  do
7:     UPDATE_ROUND_BIDS( $p, matching, roundBidders, \Delta_p, A, m, n, \epsilon$ )
8:     UPDATE_MATCHING( $\Delta_p, roundBidders, p, matching, m, peopleMatched$ )
9:   return  $match$ 
10:
11: procedure UPDATE_ROUND_BIDS( $p, matching, roundBidders, \Delta_p, A, m, n, \epsilon$ )
12:    $\Delta_p \leftarrow \mathbf{0}, roundBidders \leftarrow \mathbf{0}$ 
13:   for person in  $\{1..m\}$  do
14:     if  $matching[person] = 0$  then ▷ solo le persone non accoppiate partecipano
15:        $bid, bidTarget \leftarrow \text{GET\_PERSON\_ROUND\_BID}(p, person, A, m, n, \epsilon)$ 
16:       if  $bid > \Delta_p[bidTarget]$  then ▷ salva solo info su bid massime per oggetto
17:          $\Delta_p[bidTarget] \leftarrow bid$ 
18:          $roundBidders[bidTarget] \leftarrow person$ 
19:
20: function GET_PERSON_ROUND_BID( $p, person, A, m, n, \epsilon$ )
21:    $maxProfit \leftarrow -\infty, sndMaxProfit \leftarrow -\infty, target \leftarrow -1$ 
22:   for object in  $\{1..n\}$  do
23:     if  $A[person, object] \neq -\infty$  then ▷ archi non esistenti hanno peso  $-\infty$ 
24:        $profit = A[person, object] - p[object]$ 
25:       if  $profit > maxProfit$  then
26:          $sndMaxProfit \leftarrow maxProfit$ 
27:          $maxProfit \leftarrow profit$ 
28:          $target \leftarrow object$ 
29:       else if  $profit > sndMaxProfit$  then
30:          $sndMaxProfit \leftarrow profit$ 
31:   return ( $maxProfit - sndMaxProfit + \epsilon$ ),  $target$ 

```

```

32: procedure UPDATE_MATCHING( $\Delta_p, roundBidders, p, matching, n, peopleMatched$ )
33:   for object in  $\{1..n\}$  do
34:     if  $\Delta_p[object] > 0$  then ▷ l'oggetto ha ricevuto almeno un incremento
35:        $p[object] \leftarrow p[object] + \Delta_p[object]$ 
36:        $previousHighestBidder \leftarrow person \text{ s.t. } matching[person] = object$ 
37:       if  $previousHighestBidder \leftarrow 0$  then
38:          $peopleMatched \leftarrow peopleMatched + 1$ 
39:       else
40:          $matching[previousHighestBidder] \leftarrow 0$ 
41:          $matching[roundBidders[object]] \leftarrow object$ 

```

2 Metodologia adottata per la soluzione

Per implementare una versione parallela si è partito dall'osservazione che nell'*auction algorithm* ogni round dipende strettamente dai risultati del precedente, richiedendo una successione seriale. Similmente ogni fase di assegnamento può essere eseguita solo in seguito alla conclusione della fase di offerta che la precede. Infatti, per poter poter determinare il miglior offerente per un certo oggetto è necessario attendere che tutte le persone non assegnate abbiano proposto i propri incrementi.

Nonostante ciò, il lavoro all'interno di ciascuna fase può essere parallelizzato in maniera abbastanza naturale, visto che ciascun offerente sceglie su che oggetto puntare in modo indipendente dagli altri, così come l'assegnamento dei lotti non condiziona i rimanenti. Di conseguenza si è pensato di suddividere il lavoro su due kernel principali, uno per ciascuna fase dell'algoritmo: *bidding kernel* e *assignment kernel*, delegando alla CPU il compito di avviare l'esecuzione di un round finché necessario.

Analogamente al precedente caso seriale, come input si considera di ricevere una matrice di pesi non negativi $A \in \mathbb{Z}^{m \times n}$, dove si utilizza un valore predefinito equiparabile a $-\infty$ per gli archi assenti, e il valore di ϵ . Inoltre, si utilizza una variabile intera, letta da *host* e scritta dal *device* per indicare il numero delle persone accoppiate in un certo round. La funzione ritorna quindi l'assegnamento trovato quando questo esiste.

L'algoritmo 2 descrive in maniera semplificata gli aspetti principali della proposta parallela sviluppata, cercando di trascurare gli aspetti più tecnici. Di seguito si analizzano più nel dettaglio le diverse funzioni, indicando eventuali modifiche implementative presenti nel codice e motivando le scelte effettuate.

Algorithm 2 Auction algorithm, versione parallela

```

1: function AUCTION_PARALLEL( $A, m, n, \epsilon$ )
2:    $pplMatched \leftarrow 0$ 
3:   while  $pplMatched < m$  do
4:      $personBiddedForObj \leftarrow \mathbf{False}$ 
5:     kernel exec BIDDING( $A, p, m, n, \epsilon, objRecBidFrom, bidInc$ )
6:     kernel exec ASSIGNMENT( $p, match, pplMatched, m, n, objRecBidFrom, bidInc$ )
7:   return match

```

Internamente, la funzione “AUCTION_PARALLEL”, utilizza diversi vettori e matrici per rappresentare le strutture dati utilizzate dall'algoritmo. La maggior parte di queste viene utilizzata esclusivamente dalla GPU e dunque viene allocata esclusivamente in VRAM. Fanno eccezione il contatore di persone correntemente accoppiate e il vettore corrispondente all'accoppiamento, entrambi interi, che devono essere inizializzati e letti da *host*. Per tali strutture

viene quindi utilizzata della memoria *page-locked* così da velocizzare i trasferimenti. Oltre a ciò *host* inizializza un secondo array di interi di lunghezza n che rappresenta l'assegnamento dal punto di vista degli oggetti, utilizzato nell'implementazione della fase di assegnamento per questioni di efficienza (si veda linea 65 dello pseudo codice).

Riguardo al *bidding kernel*, visto che ciascuna persona non accoppiata deve considerare tutti i diversi oggetti per determinare il rialzo più vantaggioso, si è deciso di utilizzare delle riduzioni parallele lungo le righe della matrice dei valori A , diminuiti del prezzo corrente per ciascun oggetto in modo da considerare i profitti. In particolare, si è utilizzata una griglia bidimensionale di blocchi, nella quale la prima dimensione viene utilizzata nella riduzione di una singola persona, mentre la seconda permette di distinguere tra persone diverse. Rispetto allo pseudocodice, in sede implementativa si è preferito limitare il numero di blocchi lanciati, utilizzando dei valori multipli del numero di *Streaming Multiprocessor* (SM) della scheda, iterando sulla dimensione della griglia (*grid stride*) in entrambe le dimensioni per gestire istanze di dimensione arbitraria. Tale aspetto permette, inoltre, di gestire la riduzione lungo le righe di A mediante un numero determinato di fasi, come descritto nel seguito.

```

8: procedure BIDDING( $A, p, m, n, \epsilon, objRecBidFrom, bidInc$ )
9:   for unmatched person in parallel on blockIdx.y do
10:     $object \leftarrow blockIdx.x * blockDim.x + threadIdx.x$ 
11:    if  $object \geq m$  then  $\triangleright$  si assume un numero adeguato di thread lungo l'asse x
12:      return
13:    if  $A[person, object] \neq -\infty$  then
14:       $profit \leftarrow A[person, object] - p[object]$ 
15:       $twoMaxPair \leftarrow (profit, -\infty, object)$   $\triangleright$  dato con campi max, sndMax e maxIdx
16:    else
17:       $twoMaxPair \leftarrow (-\infty, -\infty, -1)$ 
18:     $twoMaxPair \leftarrow BLOCK\_PARALLEL\_REDUCTION\_PAIR(twoMaxPair)$ 
19:    if ( $threadIdx.x = 0$ )
20:       $aux[person, blockIdx.x] = twoMaxPair$ 
21:    all threads of the block are synchronized
22:    if blockIdx.x is the last block of the reduction then
23:      if  $threadIdx.x < gridDim.x$  then  $\triangleright$  assumiamo  $gridDim.x \leq threadIdx.x$ 
24:         $twoMaxPair \leftarrow aux[person, threadIdx.x]$ 
25:      else
26:         $twoMaxPair \leftarrow (-\infty, -\infty, -1)$ 
27:       $twoMaxPair \leftarrow BLOCK\_PARALLEL\_REDUCTION\_TWO\_PAIR(twoMaxPair)$ 
28:      if ( $threadIdx.x = 0$ )
29:         $bidInc \leftarrow twoMaxPair.max - twoMaxPair.sndMax + \epsilon$ 
30:         $objRecBidFrom[twoMaxPair.maxIdx, person] \leftarrow True$ 

```

Nella realizzazione della riduzione parallela, che calcola sia i due valori massimi sia l'indice del massimo, corrispondente con l'oggetto per cui una persona incrementerà l'offerta, si è optato per uno schema in due fasi che aggregasse inizialmente le informazioni di ciascun blocco per poi far ripetere la riduzione ad un unico blocco così da ottenere il risultato finale. Al fine di svolgere entrambe nella stessa funzione, evitando il lancio di più kernel, si è adoperato un contatore ausiliario posto in memoria globale che indicasse per ogni persona il numero di blocchi che che avessero concluso la prima fase. Per la gestione di quest'ultimo si è utilizzata l'operazione atomica *atomicInc*, trattandosi di una variabile condivisa con scritture concorrenti.

Coerentemente con lo pseudo codice, *BLOCK_PARALLEL_REDUCTION_TWO_PAIR* utilizza internamente due chiamate a *WARP_PARALLEL_REDUCTION_TWO_PAIR*, sfruttando

la memoria condivisa dal blocco per memorizzare i risultati intermedi. Avendo optato per l'utilizzo del tipo *double* per la gestione degli array relativi agli incrementi, così da limitare possibili errori dovuti all'aritmetica a virgola mobile, si è cercato di strutturare la *shared* in modo tale da prevenire l'insorgere di conflitti di banco. In particolare sono stati utilizzati diversi vettori di interi "spacchettando" e "re-impacchettando" i valori di *max* e *sndMax* tramite appositi costrutti forniti da CUDA.

Nella funzione `WARP_PARALLEL_REDUCTION_PAIR` si richiamano quindi più volte delle funzioni di *shuffle* per realizzare efficientemente la riduzione a livello dei *warp*.

```

31: function BLOCK_PARALLEL_REDUCTION_TWO_PAIR(twoMaxPair)
32:   lane  $\leftarrow$  threadIdx.x mod warpSize, wid = threadIdx.x div warpSize
33:   twoMaxPair  $\leftarrow$  WARP_PARALLEL_REDUCTION_TWO_PAIR(twoMaxPair)
34:   if lane = 0 then
35:     sharedAux[wid] = twoMaxPair ▷ sharedAux vettore in memoria shared
36:   all threads of the block are synchronized
37:   if threadIdx.x < blockIdx.x div warpSize then
38:     twoMaxPair  $\leftarrow$  sharedAux[wid]
39:   else
40:     twoMaxPair  $\leftarrow$   $(-\infty, -\infty, -1)$ 
41:   if wid = 0 then
42:     twoMaxPair  $\leftarrow$  WARP_PARALLEL_REDUCTION_TWO_PAIR(twoMaxPair)
43:   return twoMaxPair
44:
45: function WARP_PARALLEL_REDUCTION_PAIR(twoMaxPair)
46:   candidate  $\leftarrow$   $(-\infty, -\infty, -1)$ 
47:   for offset  $\leftarrow$  1 to warpSize multiplying by 2 do
48:     candidate  $\leftarrow$  _shfl_xor_sync with offset for all twoMaxPair components
49:     if candidate.max > twoMaxPair.max then
50:       twoMaxPair.sndMax  $\leftarrow$  max(twoMaxPair.max, candidate.sndMax)
51:       twoMaxPair.max  $\leftarrow$  candidate.max
52:       twoMaxPair.maxIdx  $\leftarrow$  candidate.Idx
53:     else if candidate.max > twoMaxPair.sndMax then
54:       twoMaxPair.sndMax  $\leftarrow$  candidate.max
55:   return twoMaxPair

```

Una migliaia che dovrebbe permettere di bilanciare maggiormente il carico tra i diversi SM, aumentando al contempo la coalescenza negli accessi agli array ausiliari, utilizzati nel *bidding kernel* per memorizzare i risultati intermedi della riduzione parallela, consiste nell'utilizzo di una coda contenente le persone non accoppiate in un certo round. Tale struttura viene implementata attraverso un contatore *unmatchedPeople* e un array di interi di dimensione *m* dove le prime posizioni contengono gli indici delle persone attualmente disaccoppiate. Inizialmente queste corrispondono con l'intero insieme dei partecipanti all'asta, inizializzato da *host*, mentre successivamente la coerenza della struttura verrà mantenuta dai kernel sulla base dell'accoppiamento.

Terminata la fase di offerta, si passano i risultati sui rialzi e relativi target all'*assignment kernel* che avrà la responsabilità di aggiornare conseguentemente l'assegnamento e i prezzi degli oggetti. Similmente a quanto fatto nella fase di assegnamento anche questo kernel esegue una riduzione parallela per determinare il miglior offerente per ciascuno oggetto. In questo caso però si è preferito l'utilizzo di una griglia monodimensionale, sempre di dimensione multipla del numero di SM, dove ciascun blocco viene legato a oggetti distinti rispetto alla dimensione *y*

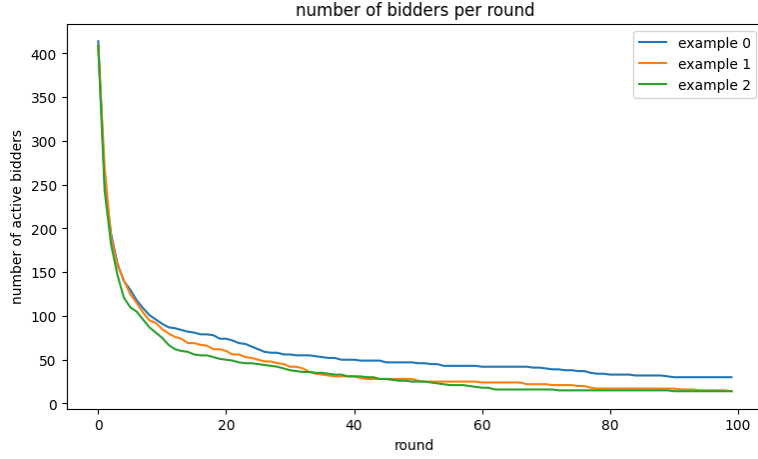


Figura 1: Numero di offerenti attivi durante i primi 100 round di alcune istanze contenenti 1108 persone ed altrettanti oggetti con densità 1.

(utilizzata al posto della x per una migliore distinzione semantica), mentre la riduzione avviene tra i thread dell'unico blocco sulla dimensione x .

Tale scelta è motivata dal fatto che, come si può notare dall'esempio riportato in figura 1, si è osservato che il numero di offerenti tende a decrescere in modo esponenziale al procedere dei round. Questo fenomeno fa sì che ciascun oggetto debba scegliere il miglior offerente tra un numero sempre minore di candidati, rendendo svantaggioso l'utilizzo di una riduzione in più passaggi.

In ogni caso, per garantire la gestione di un numero arbitrario di contendenti, si è utilizzata una fase preliminare di iterazione con *stride* pari alla dimensione x del blocco.

Analogamente a quanto fatto per l'*assignment kernel* anche in questo caso nell'implementazione si utilizza una coda per indicare quali oggetti abbiano ricevuto almeno un offerta e debbano quindi partecipare alla fase di assegnamento. La struttura dati scelta è sempre un array di interi, questa volta di dimensione n , associato ad un contatore, e viene gestita dal precedente kernel durante la fase di offerta.

```

56: procedure ASSIGNMENT( $p, match, pplMatched, m, n, objRecBidFrom, bidInc$ )
57:   for unmatched object in parallel on blockIdx.y do
58:      $maxPair \leftarrow (-\infty, -1)$  ▷ dato con campi max e maxIdx
59:     for  $person \leftarrow blockIdx.x * blockDim.x + threadIdx.x$  with x-blockstride do
60:       if  $objRecBidFrom[object, person]$  then
61:         if  $bidInc[person] > maxPair.max$  then
62:            $maxPair \leftarrow (maxPair.max, maxPair.maxIdx)$ 
63:        $maxPair \leftarrow BLOCK\_PARALLEL\_REDUCTION\_PAIR(maxPair)$ 
64:       if  $threadIdx = 0$  then
65:          $prevHighestBidder \leftarrow person$  s.t.  $match[person] = object$ 
66:         if  $prevHighestBidder = -1$  then
67:            $atomicAdd(pplMatched, 1)$ 
68:         else
69:            $match[previousHighestBidder] \leftarrow -1$ 
70:            $p[object] \leftarrow p[object] + maxPair.max$ 
71:            $match[object] \leftarrow maxPair.maxIdx$ 

```

BLOCK_PARALLEL_REDUCTION_PAIR e WARP_PARALLEL_REDUCTION_PAIR, come indicato anche nello pseudo codice, sono funzioni analoghe a quelle utilizzate nella fase di offerta, presentando le stesse peculiarità implementative e per tale ragione il loro pseudocodice viene omesso. La differenza principale sta nel fatto che, in questo caso, le strutture gestite comprendono unicamente l'informazione sul valore massimo e sul suo indice, corrispondenti rispettivamente al miglior rialzo e al relativo offerente.

```

72: function BLOCK_PARALLEL_REDUCTION_PAIR(MaxPair)
73:   analogous to BLOCK_PARALLEL_REDUCTION_TWO_PAIR considering
74:   just the max and maxIdx components, calls WARP_PARALLEL_REDUCTION_PAIR
75:   similarly to WARP_PARALLEL_REDUCTION_TWO_PAIR

```

Terminata la fase di riduzione, un unico thread per blocco si occupa di alzare il prezzo dell'oggetto al quale è legato e di aggiornare il vettore *match*, incrementando atomicamente il contatore di persone assegnate quando l'oggetto diventi assegnato per la prima volta. Alternativamente, il precedente possessore rimane disaccoppiato e viene inserito nella coda di offerenti per il round successivo.

Per completare tale coda, aggiungendo le persone attive che non si sono aggiudicate alcun oggetto viene quindi utilizzato un ultimo kernel, non riportato nello pseudocodice, che scandisce parallelamente la coda (vettore) degli offerenti del round corrente alla ricerca di quelli non assegnati.

Quest'ultima esecuzione può essere svolta in parallelo con la copia del contatore di persone assegnate, da *device* a *host*, così come accade tra l'inizializzazione di alcune strutture dati e l'esecuzione dei kernel. Per tale ragione l'implementazione proposta prevede l'utilizzo di più *stream* diversi e la loro sincronizzazione tramite eventi. A tal proposito, la figura 2 descrive schematicamente la ripartizione del lavoro all'interno di un singolo round con le relative dipendenze.

Oltre a ciò, visto che l'algoritmo prevede un'esecuzione di svariati round con la medesima struttura, si è realizzata una variazione del codice che utilizza il costrutto dei *CUDA graph* in modo da provare a ridurre l'overhead presente nel lancio dei kernel.

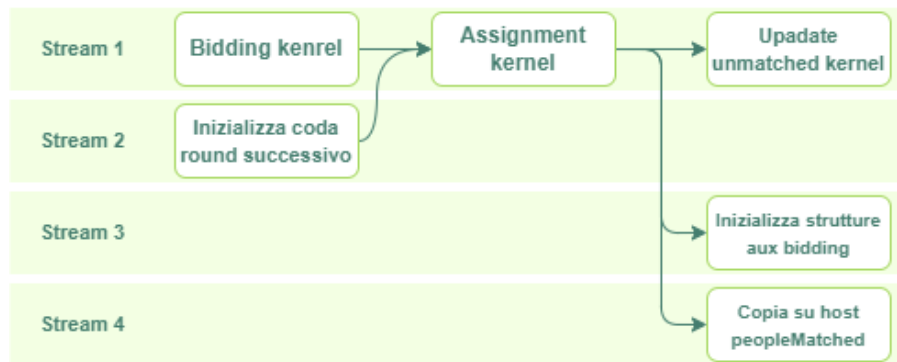


Figura 2: Suddivisione dell'esecuzione di un round tra stream diversi per l'implementazione parallela.

Tabella 1: Specifiche del sistema utilizzato per eseguire i test.

O.S.	Ubuntu 22.04.1 LTS
CPU	i5 4460 @ 3.2GHz
RAM	2x4GB DDR3 1866 MHz
GPU	GeForce GTX 960 @ 1.25GHz
Architettura	Maxwell
C.C.	5.2
SM	8
VRAM	1994 MB
versione CUDA	11.6.112

3 Risultati sperimentali

Per testare il comportamento della soluzione proposta si sono eseguiti diversi test, confrontando la soluzione seriale con quella parallela in più varianti. In particolare, sono state considerate le seguenti versioni:

- non deterministica, che presenta le caratteristiche descritte nella sezione precedente. Di conseguenza, nel caso in cui più offerenti propongano il medesimo rialzo d'offerta per uno stesso oggetto o quando più oggetti presentino un uguale profitto, l'accoppiamento e la scelta dell'obiettivo del rialzo avviene in maniera non deterministica a causa della variabilità insita nella schedulazione dei blocchi;
- una variante deterministica del precedente codice che si comporta come l'algoritmo seriale rispetto al caso di pareggi tra più offerenti, ossia assegnando un oggetto alla persona con indice minore tra quelle che propongono il rialzo maggiore in un certo round;
- una versione del primo codice che sfrutta i *CUDA graph* per l'esecuzione.

In tutti i casi si sono generate in maniera pseudo casuale le matrici dei pesi (positivi) per grafi bipartiti, assicurandosi di utilizzare un *seed* di generazione noto e garantendo per costruzione la presenza di un assegnamento tra persone ed oggetti. Per fare ciò si è generata per ciascuna istanza una permutazione casuale degli oggetti, utilizzandola per definire un assegnamento iniziale e successivamente si sono riempiti i valori dei pesi mancanti in base alla densità desiderata. Tali istanze, salvate in binario su file, sono state fornite in input agli algoritmi d'asta seriali e paralleli, misurando i tempi di esecuzione attraverso gli eventi CUDA.

I test sono stati eseguiti su una macchina con le specifiche riportate in tabella 1.

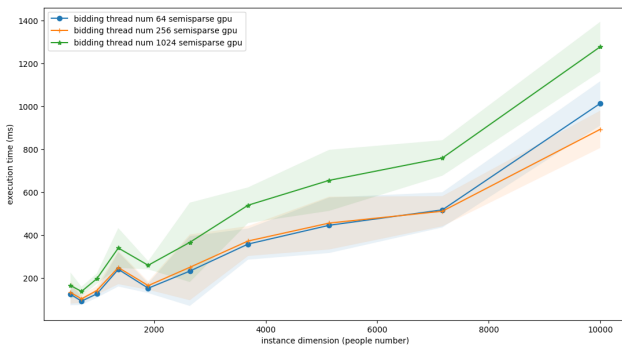


Figura 3: Tempi di esecuzione della versione non deterministica per diversi valori del numero di thread per il bidding kernel

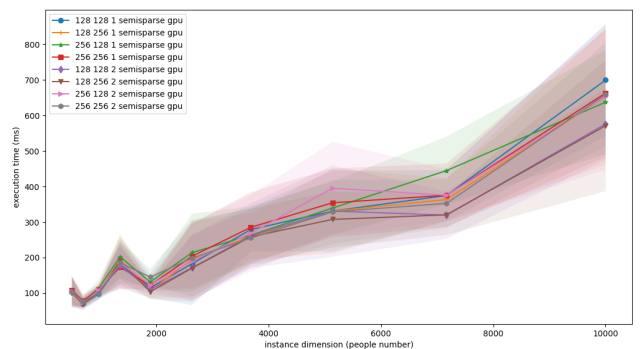


Figura 4: Tempi di esecuzione della terza versione parallela con diversi numero di thread e di blocchi

Durante lo sviluppo si sono considerate diversi valori per il numero dei thread e le dimensioni dei blocchi lanciati. Ad esempio nella figura 3 si sono valutati i valori 64, 256 e 1024 per il

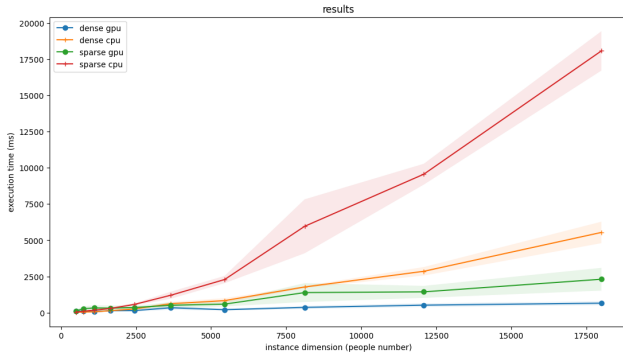


Figura 5: Tempi di esecuzione della versione non deterministica contro quella seriale

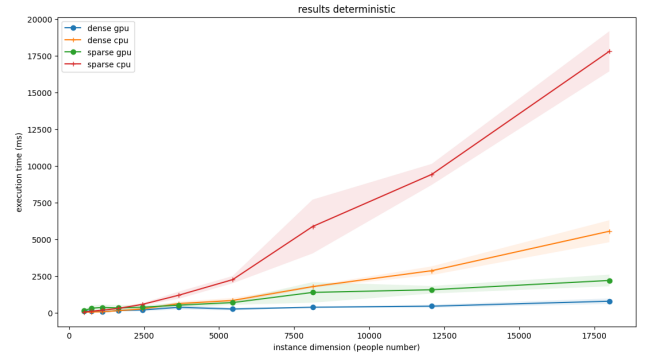


Figura 6: Tempi di esecuzione della versione parallela deterministica contro quella seriale

numero di thread da utilizzare nel lancio del *bidding kernel*, mentre la figura 4 rappresenta i risultati ottenuti utilizzando i *CUDA graph* e variando il numero di thread di entrambe le fasi dell'*auction algorithm* così come la dimensione y della griglia. In particolare, per quest'ultimo parametro si sono considerati valori pari al numero di SM della scheda oppure al doppio di esso.

In entrambi i casi i test hanno considerato delle istanze con un numero variabile di persone, equivalente al numero di oggetti, ed una densità al 50%. Per ciascuna dimensione si sono generate casualmente 3 istanze, considerandone i tempi medi d'esecuzione e utilizzando un valore di ϵ pari a 1.

Si è dunque scelta la configurazione migliore, che utilizza 128 thread per il kernel d'offerta e 256 per quello d'assegnamento, impostando un valore doppio rispetto al numero di SM per la dimensione y delle relative griglie. Tale configurazione viene quindi sottintesa nei risultati seguenti.

Per valutare la differenza di performance tra le versioni parallele si sono eseguiti dei test al variare del numero di persone, utilizzando lo stesso numero di oggetti e due diversi valori di densità: 10% per ottenere un grafo sparso e 100% per uno denso. In particolare si sono valutati 10 valori per m a partire da 500 e fino ad arrivare a 18000 seguendo una progressione geometrica, utilizzando 3 istanze per dimensione e un valore di ϵ pari a 1. Inoltre, per limitare il numero di passi necessari al raggiungimento della convergenza, aspetto che come verrà evidenziato in seguito risulta di particolare importanza, si è deciso di considerare solamente pesi nell'intervallo $[1, 10000]$.

Come si può osservare dai risultati riportati nelle figure 6, 5 e 8, riportati in forma numerica

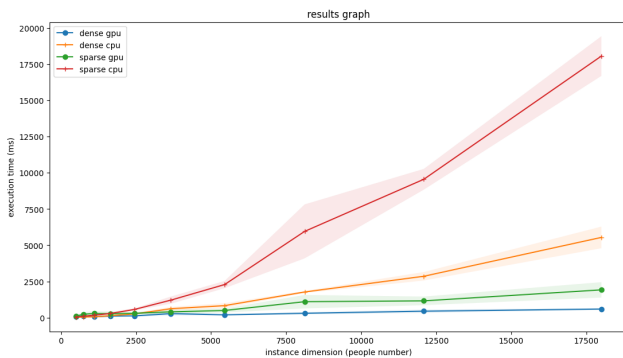


Figura 7: Tempi di esecuzione della versione con i CUDA graph contro quella seriale

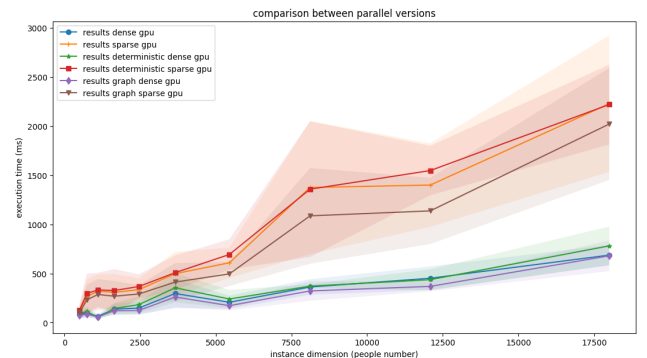


Figura 8: Confronto sui tempi di esecuzione delle varianti parallele

nella sezione B, i tempi di esecuzione della versione seriale crescono più rapidamente dei quelli delle versioni parallele che, al contrario, scalano meglio all'aumentare della dimensione dell'istanza. In particolare, si può notare che in entrambi i casi le performance risultano migliori per i grafi densi, con una differenza notevole per la versione seriale (peggioramento di 7 volte circa sull'istanza maggiore).

Tale comportamento può essere spiegato guardando al numero di round svolti in media al variare della dimensione dell'istanza. Come si evince dalla figura 9, infatti, le istanze con una densità minore richiedono sempre un numero di round inferiore per il raggiungimento di una configurazione finale, riducendo il tempo di esecuzione dell'algoritmo.

Il numero di iterazioni in generale costituisce un fattore determinante per la bontà dei tempi di esecuzione. Ciò può essere constatato guardando ai risultati riportati nella figura 10, riguardanti nuovamente il confronto dell'algoritmo seriale con quello parallelo, sulle medesime istanze. In questo caso però si è utilizzato un valore di ϵ pari a $\frac{1}{m+1}$ con m numero delle persone, che garantisce il raggiungimento di un assegnamento ottimale.

Tale modifica ha portato ad un significativo aumento dei tempi di esecuzione per entrambi gli algoritmi, risultando in una interruzione prematura dei test a causa del raggiungimento di un timeout impostato preventivamente. Inoltre, rispetto ai casi precedenti si può osservare un inversione di tendenza rispetto alla densità delle istanze. Ancora una volta la motivazione di tale comportamento può essere osservata guardando al numero di round richiesti dall'esecuzione. Infatti in questo caso la crescita di questi ultimi diviene esponenziale (si noti la scala logaritmica per l'asse delle ascisse), causando un degrado significativo dei tempi d'esecuzione.

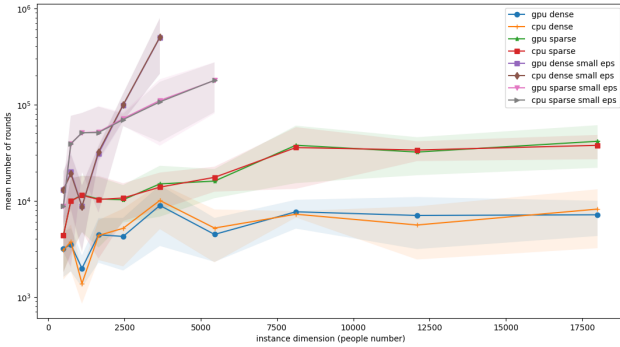


Figura 9: Numero di round eseguiti in media al variare della dimensione delle istanze per gli algoritmi seriale e parallelo non deterministico con diversi valori di ϵ

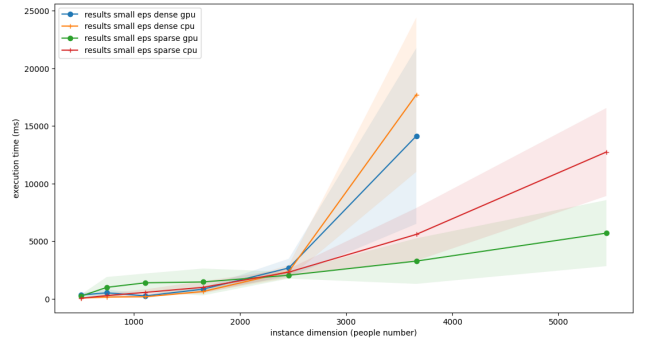


Figura 10: Confronto sui tempi di esecuzione delle varianti parallele non deterministiche con valori di ϵ sufficientemente piccolo per garantire l'ottimalità

4 Osservazioni finali

Il confronto con tra le diverse versioni ha permesso di osservare come l'implementazione parallela risulti essere in generale una scelta adeguata per l'*auction algorithm*, riuscendo a superare sempre le prestazioni della versione seriale, fatta eccezione per le istanze di dimensione minore. Tuttavia tale approccio presenta anche alcuni svantaggi come la maggior difficoltà implementativa e la necessita dell'impiego di un dispositivo di computazione ausiliario.

Riguardo alle diverse implementazioni proposte per l'algoritmo parallelo, le differenze pratiche sono risultate minime, facendo preferire le versioni non deterministiche per la maggiore semplicità del codice associata, unita alla minor probabilità di ricadere in casi "sfortunati".

Nonostante l'*auction algorithm* garantisca di raggiungere una soluzione ottimale nel caso di un ϵ sufficientemente piccolo, si è osservato come la riduzione di tale parametro causi un significativo aumento dei tempi di computazione a causa dell'elevato numero di iterazioni necessarie alla convergenza. Per tale ragione nei test si è preferito l'utilizzo di un valore unitario per ϵ che può rappresentare un buon compromesso qualora ci si possa accontentare di una soluzione sub-ottimale. In generale a seconda dello specifico scenario si potrebbe pensare di agire su tale parametro partendo da un valore anche più elevati, riducendolo fino al raggiungimento di un risultato soddisfacente.

L'approccio sviluppato presenta comunque alcune limitazioni ed inefficienze, come la serializzazione dell'esecuzione di round diversi e l'assenza di sovrapposizione nell'esecuzione dei kernel. Sebbene tali questioni siano in parte dovute alla struttura dell'algoritmo, potrebbe essere vantaggioso indagare l'utilizzo di strutture dati alternative che permettano di limitare o sovrapporre almeno parzialmente il lavoro svolto nella fase di offerta, osservando che i profitti riguardanti gli oggetti non modificati durante la fase di assegnamento rimangono costanti rispetto al round precedente.

Di seguito viene infine fornita una breve guida all'utilizzo dei principali eseguibili utilizzati.

A Miniguia all'utilizzo

Per compilare i sorgenti viene fornito un *Makefile* che permette di compilare singolarmente o cumulativamente i file per l'esecuzione dell'*auction algorithm*. In particolare l'opzione di default permette di ottenere gli eseguibili `gpu_auction`, `gpu_auction_deterministic`, `gpu_auction_graph` e `cpu_auction`, corrispondenti con le tre versioni parallele e con quella seriale. Interfaccia per l'utilizzo di questi è la medesima, permettendone l'utilizzo in modo *standalone* secondo le seguenti modalità:

- fornendo all'eseguibile il nome di un file binario contenete m , n e una matrice intera $m \times n$, rappresentante i pesi del grafo, e un valore per ϵ si possono ottenere in output l'assegnamento trovato insieme ai prezzi finali di ciascun offerente e al costo totale della soluzione;
- utilizzando gli argomenti `-generate m n d eps [seed]` è possibile specificare gli omonimi parametri, dove d indica la densità degli archi, per generare in maniera pseudo casuale un'istanza con pesi interi limitati a 10000, salvata in formato sia testuale che binario, su cui viene eseguito l'algoritmo ritornando i risultati analogamente al caso precedente.

Per generare delle istanze in blocco a partire da un file di configurazione è possibile richiamare l'eseguibile `test_case_generator` che salva i risultati in file binari posti nella cartella corrente o nella destinazione indicata come secondo parametro. In questo caso i file di configurazione devono presentare la seguente struttura:

- `SEED [LIMIT]` sulla prima riga per determinare il seed di generazione e opzionalmente un limite per la dimensione massima dei pesi da utilizzare (altrimenti viene scelto il valore massimo possibile);
- una serie di righe, una per istanza, con formato `m n density` per indicare i corrispondenti parametri di inizializzazione.

Per generare i casi di test ed eseguire gli algoritmi così da ottenere i risultati riportati nella sezione B discussi precedentemente, è possibile richiamare in sequenza i comandi `make generate_testcases` (utilizzerà approssimativamente 14 GB di spazio su disco) e `make`

`run_tests` per avviare la computazione. I risultati, vengono salvati nella cartella corrente in dei file con formato *json*.

Per la compilazione degli eseguibili in formato *standalone* viene richiesta la presenza del compilatore cuda `nvcc` mentre per eseguire i test è necessario avere installato Python in versione 3 invocabile attraverso il comando `python3`.

B Tabelle dati

Tabella 2: Tempi di esecuzione medi e deviazione standard per i diversi algoritmi proposti

numero offerenti	tempi medi di esecuzione (ms)													
	dense gpu	dense cpu	sparse gpu	sparse cpu	det. gpu	dense gpu	det. sparse gpu	graph den- se gpu	graph spar- se gpu	small dense gpu	eps. dense cpu	small sparse gpu	eps. sparse cpu	
500	85.00 (37.54)	20.65 (5.34)	117.37 (58.07)	31.48 (6.26)	92.94 (42.71)	128.68 (67.40)	73.74 (32.33)	101.06 (51.33)	332.91 (170.34)	67.15 (25.50)	235.50 (137.29)	53.56 (15.53)		
744	97.97 (41.19)	40.67 (9.58)	265.52 (172.16)	93.99 (40.66)	112.45 (50.99)	296.72 (200.90)	86.62 (36.46)	231.96 (154.01)	526.87 (284.19)	155.06 (55.82)	1001.40 (896.04)	283.08 (204.67)		
1108	64.73 (19.74)	51.04 (11.59)	320.12 (171.10)	169.61 (57.40)	52.81 (15.88)	332.83 (174.42)	57.80 (14.00)	286.57 (155.68)	256.91 (150.48)	185.28 (89.25)	1395.63 (804.18)	565.46 (263.96)		
1650	137.25 (56.77)	132.21 (26.90)	308.50 (185.74)	288.21 (97.78)	143.84 (54.01)	327.17 (217.55)	120.71 (40.80)	269.00 (155.25)	842.73 (137.65)	635.01 (76.98)	1467.23 (1161.82)	1000.44 (528.70)		
2458	148.29 (65.22)	254.13 (60.48)	335.39 (114.21)	571.08 (78.00)	182.95 (87.58)	367.07 (123.70)	124.06 (42.32)	291.47 (89.15)	2681.10 (796.02)	2367.75 (561.92)	2049.85 (286.46)	2317.74 (280.92)		
3660	296.05 (148.00)	618.13 (124.18)	498.15 (220.84)	1208.28 (245.63)	354.84 (146.66)	508.78 (176.32)	260.63 (103.46)	412.66 (193.53)	14135.33 (7620.05)	17722.16 (6690.81)	3278.19 (1973.56)	5596.85 (2286.47)		
5451	207.40 (66.80)	838.33 (177.88)	609.23 (153.59)	2306.07 (250.25)	240.65 (91.84)	693.91 (153.86)	172.92 (47.32)	496.16 (122.28)			5707.31 (2863.20)	12747.94 (3816.37)		
8117	362.77 (78.52)	1783.74 (78.09)	1375.81 (676.95)	5988.42 (1855.90)	371.49 (18.88)	1358.39 (689.89)	322.60 (99.48)	1087.09 (486.54)						
12087	451.84 (119.83)	2885.70 (297.15)	1400.44 (424.43)	9545.33 (738.69)	437.07 (106.00)	1547.90 (251.25)	368.84 (48.86)	1137.45 (336.49)						
18000	687.58 (103.78)	5541.21 (741.21)	2228.96 (693.61)	18053.06 (1392.17)	781.46 (196.09)	2220.35 (405.54)	680.20 (155.95)	2021.77 (568.53)						

Tabella 3: Numero di round medi e deviazione standard per i diversi algoritmi proposti

numero offerenti	Numero medio di round per la convergenza													
	dense gpu	dense cpu	sparse gpu	sparse cpu	det. gpu	dense gpu	det. sparse gpu	graph den- se gpu	graph spar- se gpu	small dense gpu	eps. dense cpu	small sparse gpu	eps. sparse cpu	
500	3176.87 (1558.82)	3071.33 (1552.49)	4372.33 (2552.93)	4364.33 (2557.94)	3071.33 (1552.49)	4364.33 (2557.94)	3163.73 (1526.89)	4372.33 (2552.93)	12983.27 (7189.66)	12974.00 (7424.21)	8825.00 (5829.74)	8773.67 (5866.24)		
744	3507.33 (1662.54)	3737.00 (1873.35)	9935.13 (7039.40)	9940.67 (7189.57)	3737.00 (1873.35)	9940.67 (7189.57)	3571.80 (1682.33)	9922.07 (7091.79)	19938.93 (11491.31)	19108.33 (11107.58)	39266.60 (37265.98)	39119.33 (37047.20)		
1108	1979.93 (714.35)	1377.67 (526.72)	11520.53 (6794.83)	11397.33 (6637.62)	1377.67 (526.72)	11397.33 (6637.62)	2027.13 (593.52)	11580.00 (6869.16)	8743.67 (5668.91)	8699.33 (5798.44)	50844.73 (30931.16)	50865.00 (30892.46)		
1650	4425.67 (2154.68)	4353.00 (1932.48)	10452.33 (7188.58)	10308.00 (7801.87)	4353.00 (1932.48)	10308.00 (7801.87)	4487.33 (1763.11)	10528.87 (6859.08)	30853.53 (5291.65)	32172.33 (5270.22)	51902.67 (44364.14)	51202.67 (43619.23)		
2458	4256.67 (2373.14)	5170.33 (3090.82)	10241.33 (4299.81)	10676.33 (4504.33)	5170.33 (3090.82)	10676.33 (4504.33)	4168.73 (1843.62)	10541.07 (4061.93)	98341.80 (30987.04)	99424.00 (31197.76)	71565.07 (11296.83)	69787.33 (10297.77)		
3660	8935.07 (5543.74)	10066.67 (4992.63)	14965.73 (8134.88)	13846.33 (5763.02)	10066.67 (4992.63)	13846.33 (5763.02)	9234.07 (4502.62)	14501.40 (8401.92)	497157.07 (288563.89)	502452.33 (292669.05)	110249.27 (72795.18)	106652.00 (65676.08)		
5451	4465.53 (2148.99)	5199.33 (2910.33)	15995.80 (5355.98)	15995.80 (5076.97)	5199.33 (2910.33)	15995.80 (5076.97)	4374.47 (2003.38)	15702.53 (5367.41)			177905.27 (96852.25)	178896.00 (95227.85)		
8117	7676.47 (2541.58)	7244.33 (543.67)	37770.93 (22505.28)	35732.33 (22439.67)	7244.33 (543.67)	35732.33 (22439.67)	8457.73 (4305.18)	36962.40 (21007.62)						
12087	7030.67 (3883.93)	5592.33 (3142.12)	32082.73 (13655.47)	33686.33 (7832.03)	5592.33 (3142.12)	33686.33 (7832.03)	5680.53 (1954.80)	30132.13 (14260.18)						
18000	7145.13 (2837.96)	8186.00 (4972.83)	41501.73 (19470.77)	37676.67 (10612.75)	8186.00 (4972.83)	37676.67 (10612.75)	8258.53 (4945.32)	42904.80 (18815.52)						

Riferimenti bibliografici

Bertsekas, D. P. (1979). A distributed algorithm for the assignment problem. *Lab. for Information and Decision Systems Working Paper, MIT*.

Bertsekas, D. P. e Castanon, D. A. (1989). The auction algorithm for the transportation problem. *Annals of Operations Research*, 20(1):67–96.

Vasconcelos, C. N. e Rosenhahn, B. (2009). Bipartite graph matching computation on gpu. In *Energy Minimization Methods in Computer Vision and Pattern Recognition: 7th International Conference, EMMCVPR 2009, Bonn, Germany, August 24-27, 2009. Proceedings 7*, pages 42–55. Springer.