



Αριστοτέλειο Πανεπιστήμιο Θεσσαλονίκης (Α.Π.Θ.)

HY3603 Παράλληλα και Διανεμημένα Συστήματα

Εργασία 4: Εύρεση πλήθους τριγώνων μη κατευθυνόμενου απλού γράφου.

Αντωνιάδης Δημήτριος (8462): akdimitri@auth.gr

7 Αυγούστου 2019

Περιεχόμενα

1	Εισαγωγή.	2
2	Αλγόριθμοι εύρεσης πλήθους τριγώνων n_T μη κατευθυνόμενου απλού γράφου $G(V, E)$.	2
3	Περιεχόμενα φακέλου εργασίας και Μεταγλώττιση.	10
4	Αποτελέσματα.	11
5	Επίλογος.	13

1 Εισαγωγή.

Το παρόν έγγραφο αποτελεί την αναφορά της τέταρτης (4ης) εργασίας του μαθήματος *Παράλληλα και Διανεμημένα Συστήματα*. Στο πλαίσιο της εργασίας αυτής υλοποιήθηκε παράλληλος αλγόριθμος σε CUDA, ο οποίος διαβάζει σε αραιή μορφή τον πίνακα γειτνίασης A ενός απλού, μη κατευθυνόμενου γράφου $G(V, E)$ και υπολογίζει τον αριθμό των τριγώνων που σχηματίζονται σε αυτόν.

Ο πηγαίος κώδικας της εργασίας και λοιπά αρχεία περιλαμβάνονται στον παρακάτω σύνδεσμο:

<https://github.com/akdimitri/CUDA-Parallel-Triangle-Counting>

Παραπάνω έγινε μία σύντομη περιγραφή της εργασίας, στην επόμενη (2η) ενότητα αναλύονται οι αλγόριθμοι που υλοποιήθηκαν. Στην τρίτη (3η) ενότητα παρουσιάζονται τα αρχεία του repository που περιλαμβάνει η εργασία και ο τρόπος μεταγλώττισης του προγράμματος. Στην τέταρτη (4η) ενότητα παρουσιάζονται τα αποτελέσματα των αλγορίθμων από τις δοκιμές. Τέλος, στην πέμπτη (5η) ενότητα συνοψίζονται τα συμπεράσματα της εργασίας.

2 Αλγόριθμοι εύρεσης πλήθους τριγώνων n_T μη κατευθυνόμενου απλού γράφου $G(V, E)$.

Σύμφωνα με την εκφώνηση, η εύρεση του πλήθους των τριγώνων ενός μη κατευθυνόμενου απλού γράφου μπορεί να γίνει με τον παρακάτω αλγόριθμο.

Algorithm 1: TRIANGLES counting.

Input: A adjacency matrix of graph $G(V, E)$

Output: n_T graph's number of triangles

1 $C = ((A * A) \odot A)$ /* \odot is Hadamard product */

2 $n_T = \frac{1}{6} \sum_{ij} C_{ij}$

Αρχικά, έγινε προσπάθεια να υλοποιηθεί ο παραπάνω αλγόριθμος σε CUDA με τη χρήση της βιβλιοθήκης CUSPARSE. Πιο συγκεκριμένα, το πρόγραμμα διάβαζε τον πίνακα A σε αραιή μορφή COO (Coordinate list), τον μετέτρεπε σε μορφή CSR (Compressed sparse row) με τη χρήση της συνάρτησης `cusparsXcoo2csr(...)` και πραγματοποιούσε τον πολλαπλασιασμό $A * A$ με τη χρήση της συνάρτησης `cusparsScsrsgemm2(...)`. Το αποτέλεσμα του πολλαπλασιασμού ήταν ο πίνακας B σε μορφή CSR. Στη συνέχεια ο πίνακας αυτός μετατρεπόταν σε μορφή COO και διατηρούνταν μόνο τα στοιχεία τα οποία είχαν κοινές συντεταγμένες με τον πίνακα A .

Ωστόσο, ο αλγόριθμος αυτός απορρίφθηκε γρήγορα διότι η διαδικασία του πολλαπλασιασμού ήταν ιδιαίτερα χρονοβόρα και δεν παρατηρούνταν σημαντική επιτάχυνση έναντι του αντίστοιχου χρόνου που απαιτούσε ο αλγόριθμος στο MATLAB.

Για την βελτίωση της ταχύτητας εκτέλεσης του προγράμματος δοκιμάστηκαν και παραλλαγές του αλγορίθμου 1.

Πιο αναλυτικά δοκιμάστηκαν οι δύο παρακάτω αλγόριθμοι.

Algorithm 2: TRIANGLES counting.

Input: A adjacency matrix of graph $G(V, E)$, L Lower triangular part of matrix A **Output:** n_T graph's number of triangles

```
1  $C = ((A * A) \odot L)$  /*  $\odot$  is Hadamard product */  
2  $n_T = \frac{1}{3} \sum_{ij} C_{ij}$ 
```

Algorithm 3: TRIANGLES counting.

Input: L Lower triangular part of graph's $G(V, E)$ adjacency matrix A **Output:** n_T graph's number of triangles

```
1  $C = ((L * L) \odot L)$  /*  $\odot$  is Hadamard product */  
2  $n_T = \frac{1}{2} \sum_{ij} C_{ij}$ 
```

Στους παραπάνω αλγορίθμους έγινε χρήση του κάτω(άνω) τριγωνικού πίνακα του A . Παρόλο αυτά, ο αλγόριθμος με τη χρήση της βιβλιοθήκης CUSPARSE δεν εμφάνισε αξιόλογη απόδοση και απορρίφθηκε. Η αργή εκτέλεση του οφείλονταν κυρίως στο γεγονός ότι ο πολλαπλασιασμός $A * A$ κατανάλωνε άσκοπο χρόνο για περιττές πράξεις, δηλαδή για τον υπολογισμό στοιχείων τα οποία θα απορρίπτονταν αργότερα από το γινόμενο Hadamard. Ένας ακόμη λόγος που απορρίφθηκε ο αλγόριθμος που έκανε χρήση της βιβλιοθήκης CUSPARSE, ήταν υπερβολική μνήμη που καταλάμβαναν τα στοιχεία. Για τη χρήση των συναρτήσεων απαιτούνταν πέρα από τους πίνακες με τις συντεταγμένες των στοιχείων και πίνακες με τις τιμές τους. Οι πίνακες με τις τιμές ήταν άχρηστοι εφόσον ήταν γνωστό ότι όλες οι τιμές των μη μηδενικών στοιχείων των πινάκων ισούταν με 1.

Εφόσον παρατηρήθηκε ότι ο πολλαπλασιασμός είναι μία χρονοβόρα διαδικασία επιχειρήθηκε να υπολογισθούν μόνο τα γινόμενα των *χρήσιμων* στοιχείων. Χρήσιμα στοιχεία είναι τα στοιχεία που έχουν τιμή 1 στον πίνακα A .

Έστω ο πίνακας γειτνίασης A μεγέθους N του απλού μη κατευθυνόμενου γράφου $G(V, E)$.

$$A = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & \\ \hline 1 & 0 & 1 & 1 & 1 & \\ \hline 0 & 1 & 0 & 0 & 1 & \\ \hline 1 & 1 & 0 & 0 & 1 & \\ \hline 0 & 1 & 1 & 1 & 0 & \\ \hline \end{array}$$

Εφόσον ο πίνακας έχει μόνο τιμές 0 και 1 αρκεί να αποθηκεύσουμε μόνο τις συντεταγμένες των στοιχείων με τιμή 1. Ο A έχει $NNZ = 14$ μη μηδενικά στοιχεία.

Ο πίνακας σε μορφή CSR αποθηκεύεται ως εξής (η αρίθμηση των πινάκων ξεκινάει από το 0) [1]:

$$Rows = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 2 & 6 & 8 & 11 & 14 \\ \hline \end{array}$$
$$Cols = \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 1 & 3 & 0 & 2 & 3 & 4 & 1 & 4 & 0 & 1 & 4 & 1 & 2 & 3 \\ \hline \end{array}$$

Ο πίνακας $Rows$ περιέχει στη θέση i τη θέση του πρώτου στοιχείου της γραμμής i στον πίνακα $Cols$. Ο πίνακας $Cols$ περιέχει τις αντίστοιχες στήλες των στοιχείων κάθε γραμμής με τιμή 1.

Δηλαδή, το πρώτο στοιχείο της γραμμής 0 βρίσκεται στη θέση 0 του πίνακα Cols. Η τιμή στη θέση 0 του πίνακα Cols είναι 1. Συνεπώς αυτό σημαίνει ότι στη θέση (0,1) του πίνακα, υπάρχει στοιχείο με τιμή 1.

ή για παράδειγμα,

Το στοιχείο στη θέση 2 του πίνακα Rows είναι το 6. Αυτό σημαίνει ότι η γραμμή 2 αποθηκεύει τα στοιχεία της στον πίνακα Cols από τη θέση 6 και μετά. Στη θέση 6 του Cols η τιμή είναι ένα, άρα το στοιχείο (2,1) έχει τιμή 1.

Εφόσον είναι γνωστές οι συντεταγμένες των μη μηδενικών τιμών, είναι δυνατό να υπολογίσουμε το γινόμενο $A * A$ μόνο για τις συντεταγμένες αυτές.

Τα στοιχεία με τιμή 1 έχουν τις παρακάτω συντεταγμένες:

(0,1) (0,2)

(1,0) (1,2) (1,3) (1,4)

(2,1) (2,4)

(3,0) (3,1) (3,4)

(4,1) (4,2) (4,3)

Ας δούμε για παράδειγμα το στοιχείο στη θέση (3,4). Για να υπολογισθεί ο πολλαπλασιασμός $A * A$ στη θέση (3,4) αρκεί να πολλαπλασιασθεί η γραμμή 3 με τη στήλη 4. Εφόσον ο πίνακας είναι συμμετρικός αρκεί να υπολογισθεί το γινόμενο της γραμμής 3 με τη γραμμή 4. Δεδομένου ότι τα στοιχεία είναι αποθηκευμένα σε μορφή CSR, αρκεί να βρεθεί το σύνολο των κοινών στοιχείων της γραμμής 3 και 4 στον πίνακα Cols.

Δηλαδή:

$$\begin{array}{|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 1 \\ \hline \end{array} * \begin{array}{|c|} \hline 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ \hline \end{array} = 1 * 0 + 1 * 1 + 0 * 1 + 0 * 1 + 1 * 0 = 1$$

ή

Στήλες με μη μηδενικά στοιχεία της γραμμής 3:

0	1	4
---	---	---

Στήλες με μη μηδενικά στοιχεία της γραμμής 4:

1	2	3
---	---	---

Όπως φαίνεται παραπάνω έχουν ένα κοινό στοιχείο το 1. Άρα το αποτέλεσμα στη θέση (3,4) του γινομένου $A * A$ θα είναι ίσο με 1.

Επομένως προκειμένου να υπολογισθεί πιο γρήγορα η έκφραση $(A * A) \odot A$ αρκεί για κάθε μη μηδενικό στοιχείο (i, j) του πίνακα A να βρεθούν τα κοινά στοιχεία των γραμμών i, j στον πίνακα Cols.

Input: A adjacency matrix of graph $G(V, E)$ in CSR format.

- $Rows[N + 1]$

- $Cols[NNZ]$

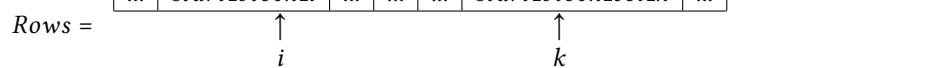
```

sum = 0
for i ← 0 to N - 1 do
    sum ← sum + B[i]

```

10 $n_T = sum/6$

Η συνάρτηση *compare_blocks* δέχεται ως ορίσματα τον πίνακα Cols, τη θέση του πρώτου στοιχείου της γραμμής *i* στον πίνακα Cols η οποία περιγράφεται στον παραπάνω αλγόριθμο από τη μεταβλητή *start_block_i*. Ακόμη, δέχεται ως όρισμα τη θέση του τελευταίου στοιχείου της γραμμής *i* στον πίνακα Cols η οποία περιγράφεται από τη μεταβλητή *stop_block_i*. Η μεταβλητή *k* περιγράφει την τεταγμένη ενός στοιχείου στη γραμμή *i* στο οποίο υπάρχει μη μηδενική τιμή και επειδή *A* συμμετρικός η στήλη *k* θα ισούται με τη γραμμή *k*. Αντίστοιχα ορίζονται οι μεταβλητές *start_block_col_k* και *stop_block_col_k*.



Algorithm 5: COMPARE blocks.

Input:

- $Cols[NNZ]$
- $start_block_i$
- $stop_block_i$
- $start_block_col_k$
- $stop_block_col_k$

Output: $matches_found$

```
1  $i = start\_block\_i$ 
2  $k = start\_block\_col\_k$ 
3  $matches\_found = 0$ 
4 while ( $i \leq stop\_block\_i$ ) && ( $k \leq stop\_block\_col\_k$ ) do
5   if  $Cols[i] == Cols[k]$  then
6      $i++$ 
7      $k++$ 
8      $matches\_found++$ 
9   else if  $Cols[i] < Cols[k]$  then
10     $i++$ 
11   else
12     $k++$ 
13 return  $matches\_found$ 
```

Ο παραπάνω αλγόριθμος μπορεί να γίνει ακόμη πιο αποδοτικός αν σκεφτεί κανείς τη σχέση

$$(A * A) \odot L$$

ή

$$(A * A) \odot U$$

Οι παραπάνω σχέσεις μπορούν να υλοποιηθούν με τη χρήση μίας εντολής *If* η οποία θα επιτρέπει την εκτέλεση του αλγορίθμου μόνο για $k > i$, δηλαδή μόνο για τα στοιχεία που βρίσκονται στον κάτω(άνω) τριγωνικό. Με τον τρόπο αυτό θα εκτελεσθούν ακριβώς οι μισές συγκρίσεις εφόσον ο A είναι συμμετρικός. Έτσι καταλήγουμε στον τελικό αλγόριθμο:

Algorithm 6: TRIANGLES counting - Serial Algorithm .

Input: A adjacency matrix of graph $G(V, E)$ in CSR format.

- $Rows[N + 1]$
- $Cols[NNZ]$

Output: n_T graph's number of triangles

```
1  $sum = 0$  for  $i \leftarrow 0$  to  $N - 1$  do
2    $start\_block\_i = Rows(i)$ 
3    $stop\_block\_i = Rows(i + 1) - 1$ 
4   for  $j \leftarrow start\_block\_i$  to  $stop\_block\_i$  do
5     if  $i < k$  then
6        $k = Cols[j]$ 
7        $start\_block\_col\_k = Rows(k)$ 
8        $stop\_block\_col\_k = Rows(k + 1) - 1$ 
9        $matches\_found =$ 
10         $compare\_blocks(Cols, start\_block\_i, stop\_block\_i, start\_block\_col\_k, stop\_block\_col\_k)$ 
11    $sum = sum + matches\_found$ 
12  $n_T = sum/3$ 
```

Δεδομένου ότι έχει υλοποιηθεί ο σειριακός αλγόριθμος τώρα πρέπει ο αλγόριθμος αυτός να μετατραπεί σε παράλληλο. Αυτό που μπορεί να παραλληλοποιηθεί είναι το εξωτερικό For loop του σειριακού αλγορίθμου και η διαδικασία σύγκρισης.

Οι κάρτες γραφικών CUDA διαθέτουν έναν αριθμό Streaming Multiprocessors. Κάθε Streaming Multiprocessor μπορεί να διαχειρίζεται threads από μόνο ένα block. Ακόμη ο μέγιστος αριθμός threads που μπορεί να διαχειρίζεται ο Streaming Multiprocessor είναι συγκεκριμένος.

Επομένως αποφασίστηκε να καλείται ένα block για κάθε γραμμή του πίνακα. Κάθε block αποφασίστηκε να έχει το μέγιστο δυνατό αριθμό threads που μπορεί να διαχειριστεί κάθε φορά ένας Streaming Multiprocessor. Αυτός ο αριθμός ονομάζεται wrap size.

Τώρα κάθε thread ενός block αποτελεί ένα στοιχείο της γραμμής $blockIdx.x$. Άρα κάθε thread ελέγχει το σύνολο των κοινών στοιχείων που έχει η γραμμή $blockIdx.x$ στον πίνακα $Cols$ με τη στήλη(γραμμή) που του αντιστοιχεί στον πίνακα $Cols$.

Algorithm 7: MAIN function Triangles counting - Parallel Algorithm .

Input: A adjacency matrix of graph $G(V, E)$ in CSR format.

- $Rows[N + 1]$
- $Cols[NNZ]$

Output: n_T graph's number of triangles

```
1 ...
2 total_sum_d = 0 /* assuming variable total_sum at device has been set to 0 */
3 kernel«< N, wrap_size»»( Rows, Cols, total_sum_d)
4 ...
5  $n_T = total\_sum/3$ 
```

Η συνάρτηση **kernel** υλοποιεί τον παράλληλο kernel σε CUDA όπως φαίνεται στον αλγόριθμο παρακάτω.

Algorithm 8: KERNEL function Triangles counting - Parallel Algorithm .

Input: A adjacency matrix of graph $G(V, E)$ in CSR format.

- $Rows[N + 1]$
- $Cols[NNZ]$

Output: $total_sum_d$: 3 times graph's number of triangles

```
1  $i = blockIdx.x$ 
2  $start\_block\_i = Rows[i]$ 
3  $stop\_block\_i = Rows[i + 1] - 1$ 
4  $j = start\_block\_i + threadIdx.x$ 
5 while  $j \leq stop\_block\_i$  do
6    $temp = 0$ 
7    $k = Cols[j]$ 
8   if  $k > blockIdx.x$  then
9      $start\_block\_col\_k = Rows[k]$ 
10     $stop\_block\_col\_k = Rows[k + 1] - 1$ 
11     $tepm =$ 
12       $compare\_blocks(Cols, start\_block\_i, stop\_block\_i, start\_block\_col\_k, stop\_block\_col\_k)$ 
13     $atomicAdd(total\_sum\_d, temp)$ 
14   $j = j + blockDim.x$ 
```

Ουσιαστικά αυτό που πραγματοποιεί ο παραπάνω αλγόριθμος είναι να εκτελεί παράλληλα σε πρώτο επίπεδο τη διαδικασία για κάθε γραμμή σε ξεχωριστό Streaming Multiprocessor και σε επόμενο επίπεδο να εκτελεί παράλληλα κάθε thread τη διαδικασία σύγκρισης του block της γραμμής i με αυτό της στήλης k των μη μηδενικών στοιχείων.

Δεδομένου ότι εντός ενός block με threads διαβάζεται σε κάθε σύγκριση το block με τα στοιχεία της γραμμής, το block με τα στοιχεία της γραμμής του πίνακα $Cols$ μπορεί να φορτωθεί στη shared memory.

Στην περίπτωση αυτή ο αλγόριθμος είναι ίδιος με τον παραπάνω με μία μικρή τροποποίηση στο Loop και στα ορίσματα της `compare_blocks`.

Algorithm 9: KERNEL function Triangles counting - Parallel Algorithm - Shared Memory.

Input: A adjacency matrix of graph $G(V, E)$ in CSR format.

- `Rows[N + 1]`
- `Cols[NNZ]`

Output: `total_sum_d`: 3 times graph's number of triangles

```

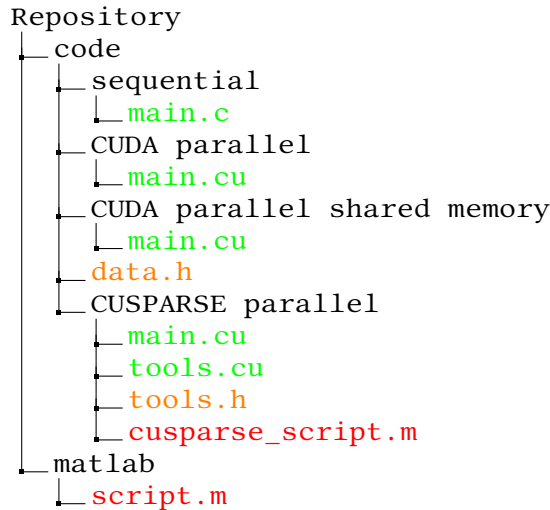
1 i = blockIdx.x
2 start_block_i = Rows[i]
3 stop_block_i = Rows[i + 1] - 1
4 nnz_at_row = stop_block_i - start_block_i
5 j = threadIdx.x
6 while j < nnz_at_row do
7   shared[j] = Cols[start_block_i + j]
8   j = j + blockDim.x
9 __syncthreads()
10 j = start_block_i + threadIdx.x
11 while j <= stop_block_i do
12   temp = 0
13   k = Cols[j]
14   if k > blockIdx.x then
15     start_block_col_k = Rows[k]
16     stop_block_col_k = Rows[k + 1] - 1
17     temp =
18       compare_blocks(Cols, shared, nnz_at_row, start_block_col_k, stop_block_col_k)
19     atomicAdd(total_sum_d, temp)
20   j = j + blockDim.x
21   __syncthreads()

```

Η συνάρτηση `compare_blocks(..)` είναι σχεδόν πανομοιότυπη. Μπορεί να βρεθεί στο αρχείο `./code/CUDA parallel shared memory/main.cu`.

3 Περιεχόμενα φακέλου εργασίας και Μεταγλώττιση.

Στο φάκελο της εργασίας περιλαμβάνονται δύο φάκελοι. Ο φάκελος code και ο φάκελος matlab.



Στο φάκελο matlab περιλαμβάνεται το αρχείο script.m. Το αρχείο αυτό υπολογίζει τον αριθμό των τριγώνων για έναν πίνακα A και υπολογίζει και το χρόνο εκτέλεσης του αλγορίθμου στο MATLAB. Επιπλέον, τυπώνει σε αρχείο τον πίνακα A σε μορφή CSR. Πιο συγκεκριμένα, δημιουργεί δύο αρχεία, ένα αρχείο που περιέχει τον πίνακα Cols και ένα αρχείο που περιέχει τον πίνακα Rows.

Στο φάκελο code περιλαμβάνεται το αρχείο data.h. Το αρχείο αυτό περιλαμβάνει το path για τα αρχεία με τους πίνακες Cols και Rows καθώς και για τις μεταβλητές N (μέγεθος πλευράς πίνακα) και NNZ (πλήθος μη μηδενικών στοιχείων πίνακα).

Στο φάκελο sequential περιλαμβάνεται ο σειριακός αλγόριθμος ενώ στους υπόλοιπους περιλαμβάνεται ο παράλληλος και ο παράλληλος με τη χρήση της shared memory αντίστοιχα.

Στο φάκελο CUSPARSE parallel περιλαμβάνεται η υλοποίηση σε CUSPARSE όπως αυτή περιγράφτηκε στην αρχή της ενότητας 2. Το αρχείο cusparse_script.m χρησιμοποιείται για την εξαγωγή των πινάκων σε μορφή COO. Το πρόγραμμα αυτό δεν κάνει χρήση του αρχείου data.h. Η εκτέλεση του προγράμματος σε CUSPARSE απαιτεί και δύο ορίσματα, τον αριθμό των γραμμών και τον αριθμό των μη μηδενικών στοιχείων.

Compilation για το σειριακό αλγόριθμο μπορεί να κάνει κανείς με τη χρήση της εντολής:

- gcc main.c -o main

Compilation για τους παράλληλους αλγορίθμους μπορεί να κάνει κανείς με τη χρήση της εντολής:

- nvcc main.cu -o main

Compilation για τον αλγόριθμο CUSPARSE μπορεί να κάνει κανείς με τη χρήση της εντολής:

- nvcc main.cu tools.cu -lcusparse -lcudart -o main

4 Αποτελέσματα.

Για τη μελέτη των αλγορίθμων χρησιμοποιήθηκαν τα παρακάτω δεδομένα:

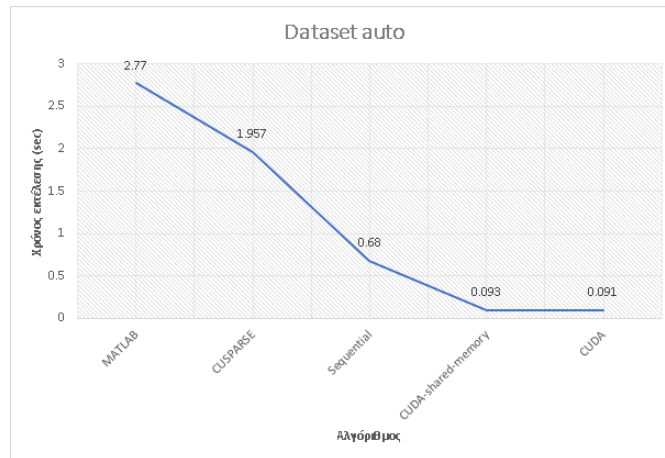
- auto : <https://sparse.tamu.edu/DIMACS10/auto>
- grat-britain_osm : https://sparse.tamu.edu/DIMACS10/great-britain_osm
- delaunay_n22 : https://sparse.tamu.edu/DIMACS10/delaunay_n22
- delaunay_n23 : https://sparse.tamu.edu/DIMACS10/delaunay_n23

Οι δοκιμές έγιναν σε υπολογιστή με επεξεργαστή **Intel® Core™ i5-4690K** ο οποίος έχει ταχύτητα ρολογιού 3.50GHz με 4 πυρήνες. Η κάρτα γραφικών που χρησιμοποιήθηκε ήταν η **GeForce GTX 650 Ti**.

Στοιχεία των datasets.

Dataset	Number of Rows	Non Zero Elements	Density	Triangles
auto	448695	6629222	$3.29 * 10^{-5}$	6245184
great-britain_osm	7733822	16313034	$2.72 * 10^{-7}$	10908
delaunay_n22	4194304	25165738	$1.43 * 10^{-6}$	8436672
delaunay_n23	8388608	50331568	$7.15 * 10^{-7}$	16873359

Αρχικά παρουσιάζονται οι χρόνοι εκτέλεσης των αλγορίθμων που παρουσιάστηκαν παραπάνω για το dataset auto.

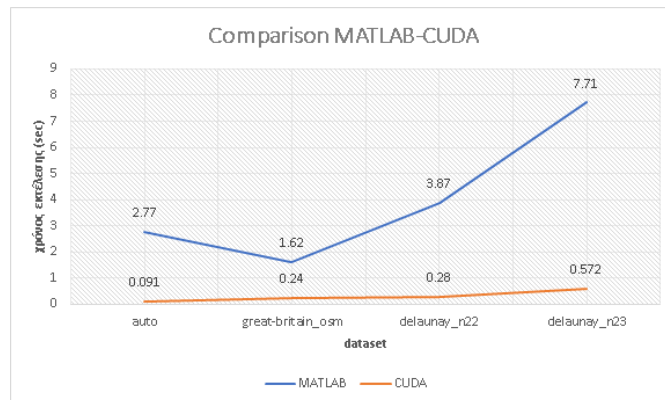


Σχήμα 1: Χρόνοι εκτέλεσης αλγορίθμων στο dataset auto.

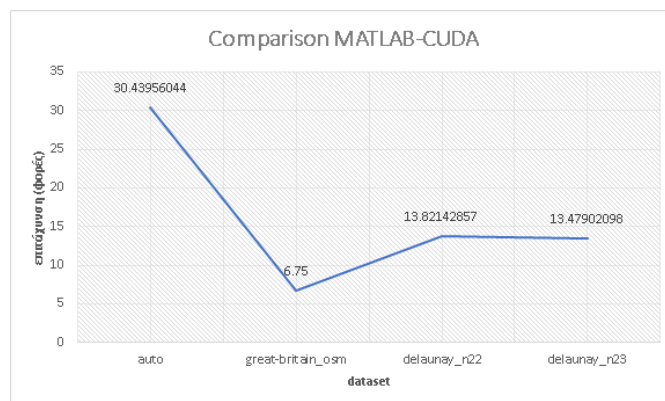
Όπως φαίνεται παραπάνω ο αλγόριθμος στο MATLAB είναι ο πιο αργός, ενώ ο αλγόριθμος σε CUSPARSE είναι μόλις **1.41** φορές πιο γρήγορος από αυτόν σε MATLAB. Ο σειριακός αλγόριθμος που υλοποιήθηκε είναι **4.07** φορές πιο γρήγορος από αυτόν σε MATLAB. Οι παράλληλοι αλγόριθμοι σε CUDA τόσο με τη χρήση Shared Memory όσο και χωρίς παρουσιάζουν ίδιο περίπου χρόνο εκτέλεσης.

Το γεγονός αυτό οφείλεται στο μικρό μέγεθος των στοιχείων που φορτώνεται στη Shared Memory σε κάθε νέο block. Ο χρόνος που κερδίζεται από τη χρήση της Shared Memoery χάνεται από τη διαδικασία φόρτωσης των δεδομένων σε αυτή. Η επιτάχυνση που παρουσιάζουν οι αλγόριθμοι σε CUDA είναι περίπου **30** φορές πιο γρήγοροι από αυτόν στο MATLAB και **7.4** φορές πιο γρήγοροι από τον σειριακό.

Στη συνέχεια παρουσιάζονται οι χρόνοι εκτέλεσης και των τεσσάρων datasets στο MATLAB και στην CUDA.



Σχήμα 2: Χρόνοι εκτέλεσης MATLAB/CUDA.



Σχήμα 3: Επιτάχυνση αλγορίθμου MATLAB/CUDA.

Όπως γίνεται αντιληπτό από τα παραπάνω διαγράμματα η επιτάχυνση κυμαίνεται από **6.75** φορές έως και **30.43** φορές.

Αν ταξινομήσουμε τα σετ δεδομένων βάση του μήκους της πλευράς τους τότε αυτά θα έχουν την εξής σειρά:

$$auto < delaunay_n22 < great - britain_osm < delaunay_n23$$

$$448695 < 4194304 < 7733822 < 8388608$$

Συνεπώς, η επιτάχυνση δεν εξαρτάται από το μέγεθος του πίνακα.

Αν ταξινομήσουμε τα σετ δεδομένων βάση των μη μηδενικών στοιχείων τότε αυτά θα έχουν την εξής σειρά:

$$\begin{aligned} auto &< great - britain_osm < delaunay_n22 << delaunay_n23 \\ 6629222 &< 16313034 < 25165738 < 50331568 \end{aligned}$$

Φαίνεται λοιπόν ότι το πλήθος των μη μηδενικών στοιχείων δεν είναι σημαντικός παράγοντας επιτάχυνσης.

Ο ουσιαστικός παράγοντας ο οποίος επηρεάζει την επιτάχυνση είναι η **πυκνότητα των δεδομένων**. Αν ταξινομήσουμε τα σετ δεδομένων βάση της πυκνότητας των μη μηδενικών στοιχείων τότε αυτά θα έχουν την εξής σειρά:

$$\begin{aligned} great - britain_osm &< delaunay_n23 < delaunay_n22 < auto \\ 2.72 * 10^{-7} &< 7.15 * 10^{-7} < 1.43 * 10^{-6} < 3.29 * 10^{-5} \end{aligned}$$

Το σετ δεδομένων *auto* επιτυγχάνει την καλύτερη επιτάχυνση διότι έχει την υψηλότερη πυκνότητα δεδομένων. Αντίθετα, το σετ δεδομένων *great-britain_osm* έχει την χαμηλότερη πυκνότητα δεδομένων. Το *great-britain_osm* έχει μέσο όρο 2 στοιχεία ανά γραμμή ενώ το *auto* έχει μέσο όρο 14 στοιχεία ανά γραμμή. Επομένως, ο αλγόριθμος επιτυγχάνει κακή παραλληλοποίηση για το πρώτο, ενώ επιτυγχάνει εμφανώς καλύτερη για το δεύτερο.

Προφανώς η **ορθότητα του αλγορίθμου** επιβεβαιώνεται από το σωστό υπολογισμό των τριγώνων ο οποίος συμπίπτει με τον υπολογισμό των τριγώνων στο MATLAB.

5 Επίλογος.

Στο έγγραφο αυτό παρουσιάστηκαν αλγόριθμοι εύρεσης του πλήθους των τριγώνων ενός απλού μη κατευθυνόμενου γράφου. Ο βέλτιστος αλγόριθμος που υλοποιήθηκε, υλοποιήθηκε σε CUDA και κατάφερε να επιτύχει βελτίωση του χρόνου εκτέλεσης από **6.75** έως **30.43** φορές, έναντι του αλγορίθμου 1 στο MATLAB.

Αναφορές

[1] Wikipedia Sparse Matrix. https://en.wikipedia.org/wiki/Sparse_matrix.