

# ΠΑΡΑΛΛΗΛΟΠΟΙΗΣΗ ΑΛΓΟΡΙΘΜΟΥ MEANSHIFT

*Παράλληλα και διανεμημένα συστήματα – 7<sup>ο</sup> Εξάμηνο*

## Ανάλυση προβλήματος

Στην τρίτη εργασία ασχολούμαστε με την παραλληλοποίηση του αλγορίθμου Mean shift στην GPU με CUDA. Ο αλγόριθμος παίρνει σαν είσοδο ένα σύνολο διανυσμάτων (σημείων)  $\mathbf{x}$  διάστασης  $d$  και στόχος είναι η εύρεση των διανυσμάτων  $\mathbf{y}$  από κάθε σημείο προς το μέγιστο για την ομαδοποίηση των δεδομένων. Λειτουργεί επαναληπτικά και σε κάθε βήμα υπολογίζονται τα νέα, πιο ακριβή  $\mathbf{y}$ , έως ότου το σφάλμα (μέτρο της διαφοράς των  $\mathbf{y}$  του τρέχοντος βήματος από το προηγούμενο) γίνει μικρότερο από μία μικρή σταθερά  $\epsilon$ .

## Πρώτες σκέψεις για παραλληλοποίηση

Ο Mean shift ενδείκνυται για παραλληλοποίηση σε δύο επίπεδα:

- 1) Ο αριθμητής και ο παρονομαστής του τύπου υπολογισμού των νέων διανυσμάτων  $\mathbf{y}$  μπορούν να υπολογιστούν ανεξάρτητα ο ένας από τον άλλον

$$\mathbf{y}_i^{[k+1]} = \frac{\sum_{j=1}^N k \left( \left\| \mathbf{y}_i^{[k]} - \mathbf{x}_j \right\|_2^2 \right) \mathbf{x}_j}{\sum_{j=1}^N k \left( \left\| \mathbf{y}_i^{[k]} - \mathbf{x}_j \right\|_2^2 \right)}$$

- 2) Αφού υπολογιστούν όλοι οι αριθμητές και παρονομαστές, τότε μπορούμε να συνεχίσουμε με τον υπολογισμό του νέου  $\mathbf{y}$  οπότε και του διανύσματος mean shift

$$\mathbf{m}_i^{[k]} = \mathbf{y}_i^{[k+1]} - \mathbf{y}_i^{[k]}.$$

Τα παραπάνω γίνονται ακόμα πιο κατανοητά αν δούμε όλα τα διανύσματα σαν πίνακες διάστασης  $N \times d$ , δηλαδή  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{N \times d}$ . Έτσι, ο παραπάνω τύπος απλοποιείται στη μορφή:

$$\mathbf{Y}^{[k+1]} = k(\mathbf{Z} \cdot \wedge 2) \cdot \mathbf{X} ./ k(\mathbf{Z} \cdot \wedge 2) \cdot \vec{1}$$

$$\text{Όπου } \mathbf{Z} \in \mathbb{R}^{N \times N}, Z_{ij} = \begin{cases} \left\| \mathbf{y}_i^{[k]} - \mathbf{x}_j \right\|_2, & \text{αν } \left\| \mathbf{y}_i^{[k]} - \mathbf{x}_j \right\|_2 > \sigma^2 \\ 0, & \text{αλλιώς} \end{cases}$$

και  $\vec{1} = [1 \quad \dots \quad 1]^T \in \mathbb{R}^d$ . Οι τελεστές  $\cdot \wedge$  και  $./$  είναι αυτοί του MATLAB, οι οποίοι πραγματοποιούν ύψωση σε δύναμη και διαίρεση αντίστοιχα ανά στοιχείο.

## Παραδοχές & Σημειώσεις

Καθώς η ανάγνωση δεδομένων από αρχείο mat είναι πολύπλοκη και απαιτεί βιβλιοθήκες που ενδεχομένως να μην είναι διαθέσιμες σε κάποια συστήματα, τα δεδομένα των

δύο αρχείων της άσκησης αποθηκεύτηκαν σε μορφή απλού κειμένου txt · ένα για κάθε μεταβλητή MATLAB στην πρώτη σειρά των οποίων αναγράφονται οι διαστάσεις του πίνακα που περιέχουν. Ακόμη οι σταθερές  $\epsilon$  και  $\sigma$  ορίζονται στην αρχή του προγράμματος, οπότε για να αλλαχθούν χρειάζεται εκ νέου compile.

## Παραλληλοποίηση με CUDA

Στην αρχή εισάγονται τα δεδομένα από το αρχείο που ορίσθηκε κατά την εκτέλεσή του προγράμματος. Αφού γίνουν οι απαραίτητες κατανομές μνήμης για κάθε πίνακα που χρειάζεται στο Device (cudaMalloc), ξεκινάει η μεταφορά του πίνακα  $X$  στην GPU και αμέσως μετά αντιγράφονται τα δεδομένα του στον  $Y$ . Εδώ είναι σημαντικό να τονισθεί ότι όλοι οι πίνακες αποθηκεύονται σε «ξετυλιγμένη» (unrolled) μορφή στη μνήμη της κάρτας γραφικών, δηλαδή σαν διανύσματα με διάσταση το γινόμενο των διαστάσεων του πίνακα από τον οποίο προέκυψαν, και ανά σειρά.

Έπειτα ξεκινάει η επαναληπτική διαδικασία με την αρχικοποίηση του σφάλματος, του αριθμητή και του παρονομαστή στο μηδέν. Εκτελείται το πρώτο kernel (cudaMeanShift1) με blocks διαστάσεων  $BLOCK\_SIZE \times BLOCK\_SIZE^1$ . Κάθε thread αντιστοιχεί σε ένα συγκεκριμένο  $i$  (rowY) και  $j$  (rowX), τα οποία δίνουν και τους δείκτες για τον πίνακα  $Z$  (colZ & rowZ). Έτσι, αφού σιγουρευτούμε ότι οι δείκτες που προέκυψαν είναι έγκυροι για τους πίνακες που έχουμε<sup>2</sup>, υπολογίζουμε το τετράγωνο της νόρμας των διανυσμάτων  $y_i$  και  $x_j$  (diffYiXj). Αν λοιπόν η ρίζα του diffYiXj είναι μικρότερη του  $\sigma^2$ , τότε υπολογίζεται η γκαουσιανή της και προστίθεται<sup>3</sup> στη σειρά  $i$  του παρονομαστή, ενώ για τον αριθμητή υπολογίζονται τα ανά διάσταση γινόμενα της γκαουσιανής με την τιμή του  $x_j$  στην αντίστοιχη διάσταση και προστίθενται στη σειρά  $i$  και αντίστοιχη διάσταση του αριθμητή. Με άλλα λόγια, πραγματοποιούνται οι πράξεις του τύπου που σημειώθηκε παραπάνω και συμπληρώνονται οι πίνακες denom και numer διαστάσεων  $N \times 1$  και  $N \times d$  αντίστοιχα. Στην περίπτωση που η απόσταση των δύο σημείων  $y_i$  και  $x_j$  είναι μεγαλύτερη από  $\sigma^2$ , τότε απλά δε συμμετέχουν στους υπολογισμούς.

Στο σημείο αυτό και προτού συνεχίσουμε στον υπολογισμό των νέων τιμών του  $Y$ , είναι απαραίτητο να έχουν τελειώσει όλα τα threads του προηγούμενου kernel, ώστε οι πίνακες του αριθμητή και του παρονομαστή να έχουν πάρει τις τελικές τιμές τους. Αυτό εξασφαλίζεται με την εισαγωγή ενός barrier στον Host μέσω της συνάρτησης cudaDeviceSynchronize.

Πλέον είμαστε έτοιμοι να εκτελέσουμε το δεύτερο και τελευταίο kernel του αλγορίθμου. Αυτή τη φορά εκκινούμε block πλήθους που εξαρτάται από τις διαστάσεις του  $Y$  και όχι του  $Z$  όπως προηγουμένως. Σε κάθε thread βρίσκουμε τις τιμές των row και col που υποδεικνύουν ποιο στοιχείο του πίνακα των νέων  $Y$  θα υπολογιστεί στο τρέχον thread. Στη συνέχεια ελέγχουμε την εγκυρότητα των δεικτών σύμφωνα με τις διαστάσεις του  $Y$ , όπως και προηγουμένως, και αν είναι εντάξει αποθηκεύουμε στην αντίστοιχη με το thread θέση του  $Y$  το νέο πηλίκο των στοιχείων του αριθμητή και του παρονομαστή

που του αναλογούν. Στην περίπτωση όπου ο παρονομαστής έχει κάποιο μηδενικό, τότε και το αντίστοιχο διάνυσμα του αριθμητή θα είναι το μηδενικό, καθώς ο αριθμητής εμπεριέχει μέσα του πολλαπλασιασμό με ίδιο όρο με τον παρονομαστή. Έτσι για να αποφευχθεί διαίρεση με το μηδέν, αντικαθιστούμε τον συγκεκριμένο παρονομαστή με τη μονάδα και στη θέση αυτή του νέου  $Y$  προκύπτει μηδέν χωρίς να δημιουργηθούν προβλήματα.

Στο εσωτερικό του ίδιου kernel πραγματοποιείται και ο υπολογισμός του σφάλματος της τρέχουσας επανάληψης. Αντί λοιπόν να αποθηκεύσουμε ολόκληρο τον πίνακα  $m$ , εφ' όσον δε χρειάζεται κάπου αλλού απλά προσθέτουμε το τετράγωνο της διαφοράς του νέου με το παλιό στοιχείο του  $Y$  που αναλογεί στο τρέχον thread. Το άθροισμα αυτό είναι το εσωτερικό της νόρμας (και όχι η νόρμα). Συνεχίζοντας, αντικαθίσταται η παλιά τιμή του  $Y$  για το συγκεκριμένο thread με την νέα της.

Ολοκληρώνοντας το βήμα της επανάληψης, και πάλι περιμένουμε να συγχρονιστούν όλα τα thread της GPU, μιας και πρέπει να ολοκληρωθεί ο υπολογισμός του mean shift πριν πάρουμε τη ρίζα του η οποία και αντιστοιχεί στη νόρμα του πίνακα  $m$ .

#### ➤ Υλοποίηση με shared memory:

Η υλοποίηση με τη χρήση της κοινής μνήμης στο εσωτερικό κάθε block εισάγει μια επιπρόσθετη πολυπλοκότητα στον αλγόριθμο με σκοπό τη μείωση χρήσης της global memory. Πιο αναλυτικά, πρέπει κάθε thread στην αρχή της εκτέλεσής του να μεταφέρει τα δεδομένα που χρειάζεται από την global στη shared memory, να κάνει όποιες πράξεις στα δεδομένα αυτά (με τις φθηνότερες κλήσεις read/write στη shared memory) και αν αλλάξει κάτι που αφορά και threads εκτός του block του πρέπει να τοποθετήσει τις νέες τιμές στην global.

Η χρησιμότητα της shared memory προκύπτει (πέραν της αλληλεπίδρασης των threads ενός block μεταξύ τους που εδώ δεν συμβαίνει γιατί το κάθε ένα χρησιμοποιεί ξεχωριστό κομμάτι από την κοινή μνήμη) όταν κάθε thread χρησιμοποιεί πολύ συχνά κάποια θέση μνήμης, καθώς αν τη χρησιμοποιεί ελάχιστα χειροτερεύει ο χρόνος εκτέλεσης αντί να καλυτερεύει αφού πραγματοποιείται τουλάχιστον μία επιπλέον κλήση της global μνήμης (π.χ.  $t_{rg} + t_{ws} + t_{rs} + t_{wg} > t_{rg} + t_{wg}$ ). Έτσι στο δεύτερο kernel είναι προφανές πως η χρήση της shared μνήμης είναι απαγορευτική για τους πίνακες  $denom$  και  $numer$ , ενώ για τον  $Y$  δεν κερδίζουμε κάτι ουσιαστικό αφού πραγματοποιούνται μονάχα δύο αναγνώσεις ανά thread.

Με παρόμοιο σκεπτικό, στο πρώτο kernel χρησιμοποιούμε shared memory μόνο για τον πίνακα  $X^4$ , γεγονός που ωφελεί κατά ένα μικρό βαθμό λόγω μικρότερου κόστους πράξεων. Στην πράξη όμως (για μικρά  $d$ ) ο χρόνος χειροτερεύει, αφού γίνονται περισσότερες πράξεις για να βρούμε τους σχετικούς δείκτες για τους shared πίνακες και απαιτείται χρόνος για την καταχώρηση και μεταφορά της μνήμης. Για μεγάλα  $d$  τα πράγματα είναι καλύτερα, διότι αντικαθίστανται περισσότερες προσπελάσεις της global μνήμης.

Τέλος, στην shared υλοποίηση ο μέγιστος αριθμός των threads ανά block μειώνεται δραστικά, καθώς η shared μνήμη για κάθε block είναι μικρή (~10KB) γεγονός που περιορίζει ακόμα περισσότερο το μέγεθος των εισαγόμενων δεδομένων.

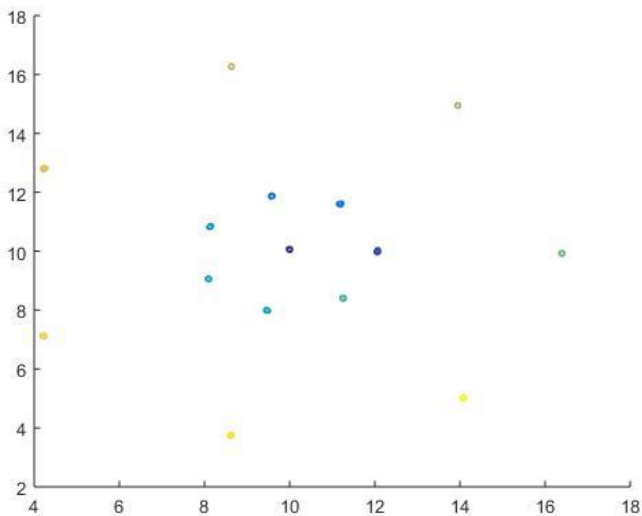
### Μετρήσεις χρόνων εκτέλεσης και αποτελέσματα (Διάδης)

	MATLAB	Σειριακός -O0	Σειριακός -O3	CUDA (global)	CUDA (shared)
a Χρόνος	0,62 secs	0,25 secs	0.069 secs	0,0028 secs	0,0031 secs
Επιτάχυνση	1x	2.5x	8x	<b>214x</b>	194x
b Χρόνος	0,27 secs	3,4 secs	0,72 secs	0,076 secs	0,098 secs
Επιτάχυνση	13x	1x	5x	<b>45x</b>	35x

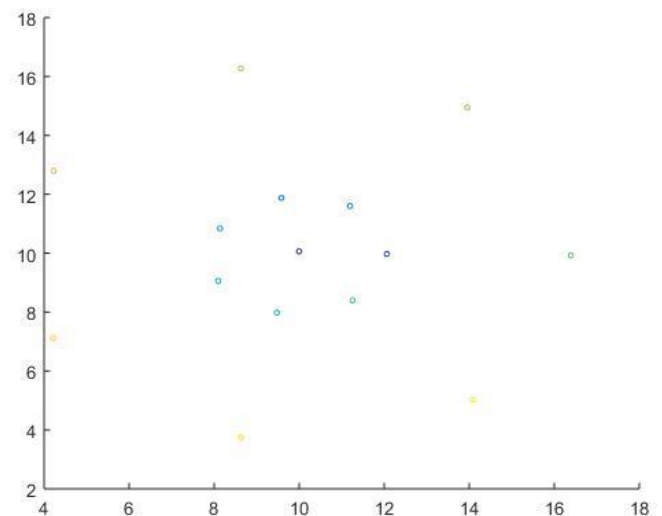
Οι παραπάνω χρόνοι αφορούν μόνο τη διάρκεια εκτέλεσης του αλγορίθμου στην GPU χωρίς τις καθυστερήσεις από μεταφορά δεδομένων από και προς την Device. Για τις μετρήσεις χρησιμοποιήθηκαν το data set διαστάσεων 600 x 2 που δόθηκε με την εκφώνηση της εργασίας (a) και ένα ακόμα διαστάσεων 1024 x 32 από την κατηγορία “DIM-sets (high)” της σελίδας <https://cs.joensuu.fi/sipu/datasets> (b). Σημειώνεται πως για την εκτέλεση με το (a) χρησιμοποιήθηκαν  $\epsilon = 0.0001$  και BLOCK\_SIZE = 20 και με το (b)  $\epsilon = 0.00015$  και BLOCK\_SIZE = 10.

### Παρατηρήσεις και συμπεράσματα

Παρατηρούμε πως και στις δύο περιπτώσεις η υλοποίηση CUDA (global) είναι η καλύτερη με αρκετά μεγάλη διαφορά. Η shared είναι πιο αργή σύμφωνα με τα όσα περιγράφηκαν παραπάνω, ενώ με την αύξηση του d βλέπουμε να πλησιάζει την global, όπως αναμέναμε.



**Σχήμα 1.** Αποτελέσματα MATLAB  
(600 x 2)



**Σχήμα 2.** Αποτελέσματα σειριακού  
και παράλληλου αλγορίθμου  
(600 x 2)

Τέλος, αν συγκρίνουμε προσεκτικά τα αποτελέσματα του κώδικα στο MATLAB με αυτά του σειριακού και του παράλληλου αλγορίθμου, θα δούμε ότι αυτά στο σχήμα 2 είναι ίδια ανά περιοχή. Αυτό συμβαίνει λόγω του ότι ο αλγόριθμος στο MATLAB είναι πιο αναλυτικός σε σχέση με τον απλό κλασματικό τύπο της εκφώνησης, οπότε και έχει πολύ μεγαλύτερη ακρίβεια. Το σφάλμα βέβαια είναι πολύ μικρό ανά σημείο (δεν έχει σχέση με το  $\epsilon$  και το πλήθος των επαναλήψεων, αλλά με την πορεία εύρεσης του αποτελέσματος και γι' αυτό θα δούμε ότι οι τελευταίοι αλγόριθμοι ολοκληρώνονται με ελαφρώς διαφορετικά σφάλματα επαναλήψεων ανά εκτέλεση σε αντίθεση με αυτόν του MATLAB).

### Υποσημειώσεις

<sup>1</sup> Η σταθερά BLOCK\_SIZE μπορεί να πάρει οποιαδήποτε ακέραια τιμή, αρκεί το BLOCK\_SIZE<sup>2</sup> να είναι μικρότερο από το μέγιστο πλήθος threads ανά block που επιτρέπει η κάρτα γραφικών που χρησιμοποιείται. Ένας επιπλέον περιορισμός αφορά το πλήθος registers ανά thread, οπότε και στις παράλληλες υλοποιήσεις χρησιμοποιήθηκαν οι μικρότεροι (από άποψη μνήμης) δυνατοί τύποι μεταβλητών (π.χ. float, unsigned short κτλ) για μέγιστη απόδοση.

<sup>2</sup> Ο έλεγχος της εγκυρότητας των δεικτών σε κάθε thread είναι απαραίτητος, καθώς ενδέχεται να μην ταιριάζει πάντα ακριβώς το πλήθος των threads με τις διαστάσεις των δεδομένων. Έτσι, μερικά threads αναπόφευκτα θα παραμείνουν απραγή. Γενικά η χρήση εντολών υπό συνθήκη (if, for κα) είναι κακή πρακτική στην CUDA. Στη συγκεκριμένη περίπτωση όμως δε δημιουργείται πρόβλημα, καθώς i) οι χρόνοι εκτέλεσης των threads είναι παραπλήσιοι (για οποιαδήποτε συνθήκη) και ii) τα for εκτελούνται για τον ίδιο αριθμό επαναλήψεων για όλα τα threads.

<sup>3</sup> Για την αποφυγή κάποιου πιθανού race condition λόγω της ταυτόχρονης εγγραφής σε ίδιες θέσεις μνήμης από πολλά threads, χρησιμοποιείται η συνάρτηση atomicAdd, η οποία με αριθμητή κινητής υποδιαστολής υποστηρίζεται για computing capability της GPU  $\geq 2.0$ .

<sup>4</sup> Δοκιμάσθηκε και υλοποίηση με τους πίνακες Y, denom και numer του πρώτου kernel στη shared memory (επισυνάπτεται με σχόλια στον κώδικα του cudaMeanshift1 της shared υλοποίησης), αλλά απορρίφθηκε λόγω της επιπλέον καθυστέρησης που εισάγει μιας και χρειάζεται σε ελάχιστες κλήσεις ανά thread.