

Παράλληλα και κατανεμημένα συστήματα

Εργασία 3

Μπεκιάρης Θεοφάνης ΑΕΜ:8200



Στόχος Εργασίας

Εκτέλεση του αλγόριθμου Non Local Mean με χρήση της Cuda για αποθορυβοποίηση εικόνων.

Παρατηρήσεις σχετικά με τον αλγόριθμο

- Ο αλγόριθμος απαιτεί την σωστή επιλογή μαθηματικής πράξης για το συμβολισμό $\|f(i)-f(j)\|$ που περιγράφει το κατά πόσο είναι όμοιες οι γειτονιές 2 pixels.
- Ο αλγόριθμος για τον υπολογισμό της νέας τιμής ενός pixel απαιτεί πρώτα τον υπολογισμό της τιμής $Z(i)$ που αντιστοιχεί σε ένα pixel και ύστερα το άθροισμα όλων των τιμών $w(i,j)*f(j)$, κάτι το οποίο καθιστά στην παραλληλοποίηση του αλγόριθμου περιπλοκότητα και επιπλέον καθυστερήσεις.

Κώδικες που περιέχονται στην εργασία

Για την διεκπεραίωση της εργασίας έχουν κατασκευαστεί:

- Ένα σειριακό πρόγραμμα το οποίο έγραψα και εκτελεί την υλοποίηση του αλγόριθμου Non Local Mean.
- Ένα πρόγραμμα που περιέχει παραλληλοποιημένη έκδοση του

προηγούμενου προγράμματος με χρήση της Cuda και το οποίο διαχειρίζεται την εικόνα σαν να αποτελείται από πολλές μικρές υποεικόνες -Ένα δεύτερο πρόγραμμα που περιέχει παραλληλοποιημένη έκδοση του προγράμματος με χρήση της Cuda το οποίο όμως λαμβάνει υπόψη του όλα τα pixel της εικόνας για την νέα τιμή των αποθυβοποιημένων pixel

Οι παραπάνω κώδικες που έχουν υλοποιηθεί σε γλώσσα C χρησιμοποιούνται σε συνδιασμό με το πρόγραμμα MATLAB.Εχω χρησιμοποιήσει τα αρχικά Scripts MATLAB που δόθηκαν μαζί με την εκφώνηση της εργασίας,τα οποία με τις απαραίτητες αλλαγές καλούν τις υλοποιήσεις του αλγόριθμου Non Local Mean.

Σειριακό Πρόγραμμα

Συγκεκριμένες λεπτομέριες για τον τρόπο υλοποίησης του αλγορίθμου δίνονται με σχόλια μέσα στο πρόγραμμα.Ένα σημείο που θα επιλέξω να αναλύσω είναι το κομμάτι για τον υπολογισμό της διαφοράς $\|f(i)-f(j)\|$.Για την συγκεκριμένη πράξη έχω επιλέξει να εκτελέσω την καρτεσιανή διαφορά των δυο pixel με βάρη.Συγκεκριμένα την ερμηνεύω ως το άθροισμα των τετράγωνων της διαφοράς (x_i-x_j) των τιμών των pixel γειτόνων πολλαπλασιαζόμενη με ένα βάρος w_i ,το βάρος που έχω επιλέξει είναι η απόσταση των pixel-γειτόνων από το κεντρικό pixel το οποίο επεξεργαζόμαστε.

$$\|f(i) - f(j)\| = \sum \{ (x_i - x_j)^2 \cdot w_k \} / \sum (w_k)$$

Το πρόγραμμα υλοποιεί την πράξη αυτή μέσω της συνάρτησης gaussianDistance().

Η συνάρτηση normFactor() υπολογίζει τις τιμές Z_i των pixel.

Η συνάρτηση weightingFunct() υπολογίζει τα βάρη $w(i,j)$.

Τέλος η συνάρτηση nonLocalMeans() χρησιμοποιεί τα βάρη $w(i,j)$ για να υπολογίσει το άθροισμα

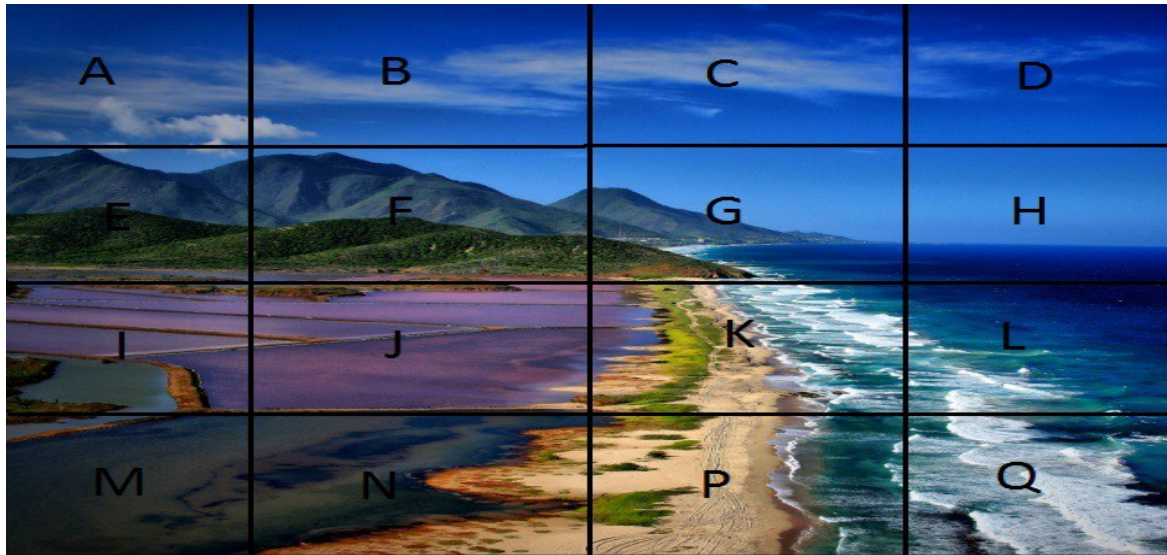
$$\hat{f}(\mathbf{x}) = \sum_{\mathbf{y} \in \Omega} w(\mathbf{x}, \mathbf{y}) f(\mathbf{y}), \quad \forall \mathbf{x} \in \Omega,$$

αρα και την νέα τιμή του pixel στο οποίο κάνουμε αποθυβοποίηση.

Πρώτο παραλληλοποιημένο πρόγραμμα σε Cuda

Η λογική με την οποία παραλληλοποίησα το αρχικό σειριακό πρόγραμμα ήταν η εξής:

Θεωρούμε ότι η αρχική εικόνα χωρίζεται σε μικρότερα κομμάτια όπως φαίνεται για παράδειγμα στην ακόλουθη εικόνα.



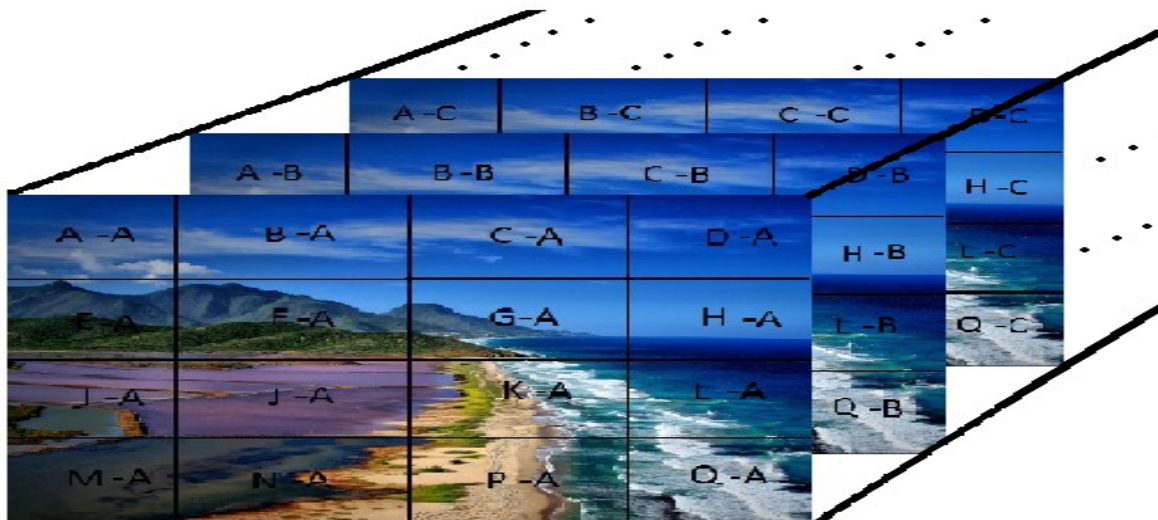
Μέσα από το Sript του MATLAB όταν θα καλέσουμε το πρόγραμμα kernel το οποίο υλοποιεί τον NonLocalMeans αλγόριθμο, δημιουργούμε ένα grid το οποίο θα περιέχει έναν αριθμό από blocks το καθένα από τα οποία θα αναλάβει να διαχειριστεί και να εκτελέσει τον αλγόριθμο σε ένα από τα κομμάτια στα οποία έχουμε χωρίσει την αρχική εικόνα. Τα blocks είναι δύο διαστάσεων και αποτελούνται από $N \times N$ threads. Ο αριθμός των κομματιών στα οποία θα χωρίσουμε την εικόνα καθώς και των blocks εξαρτάται από τον αριθμό των διαστάσεων που θα ορίσουμε να περιέχει το κάθε block. Ο αριθμός που έχω ορίσει είναι το 16 (δηλαδή blocks 16×16 ή 256 threads/block) καθώς διαιρεί ομοιόμορφα την εικόνα δεδομένου των διαστάσεων εικόνων που πρέπει να χρησιμοποιήσουμε και δεν χρειάζεται να εισάγουμε στον κώδικα επιπλέον συνθήκες ελέγχου για το αν έχουμε βγει εκτός ορίων, είναι αρκετά μεγάλος για να προκύψουν καλά αποτελέσματα αποθορυβοποίησης και δεν ξεπερνάει τα όρια των threads που μπορούμε να χρησιμοποιήσουμε ανά block. Το κάθε thread αναλαμβάνει να υπολογίσει την νέα τιμή ενός από τα pixel της εικόνας εκτελώντας τον αλγόριθμο μέσα στο block που ανήκει. Για παράδειγμα το block με συντεταγμένες (0,0) δηλαδή $blockIdx.x=0$ και $blockIdx.y=0$ φορτώνει από την global στην shared μνήμη τις τιμές των pixel που αντιστοιχούν στο χώρο A και έπειτα εκτελεί τον αλγόριθμο λαμβάνοντας υπόψη μόνο τα pixel της περιοχής A. Ομοίως η ίδια λογική ακολουθείται

και για όλες τις υπόλοιπες περιοχές B,C,D,..... κλπ. Με την παραπάνω υλοποίηση υπολογίζονται ταυτόχρονα όλες οι περιοχές τις εικόνας και η χρήση της shared memory βοηθάει ακόμα περαιτέρω τον χρόνο εκτέλεσης καθώς κάθε thread των block εκτελεί μόνο μία προσπέλαση στην global μνήμη μεταφέροντας μία τιμή ενός pixel στην shared μνήμη η οποία μπορεί να χρησιμοποιηθεί και απο τα υπόλοιπα thread του block στο οποίο ανήκει.

Δεύτερο παραλληλοποιημένο πρόγραμμα σε Cuda

Η προηγούμενη υλοποίηση αν και είναι αρκετά αποδοτική και σε χρόνο και σε ποιότητα της αποθορυβοποίησης των εικόνων διαχειρίζεται την αρχική εικόνα σαν σύνολο υπο-εικόνων αλλοιώνονται τον αλγόριθμο NonLocalMeans ο οποίος υποθέτει οτι για την νέα τιμή ενός pixel πρέπει να λάβουμε υπόψην όλα τα pixel της εικόνας. Επιπλέον με την προηγούμενη υλοποίηση παρατηρούμε οτι για πολύ μεγάλες εικόνες εμφανίζονται κάποιες γραμμές στην αποθορυβοποιημένη εικόνα. Για αυτόν τον λόγο και επειδή δεν ήξερα κατα πόσο σωστή είναι η παραπάνω λογική παραλληλοποίησης δημιούργησα και μία επιπλέον υλοποίηση η οποία συνυπολογίζει όλα τα pixel της εικόνας για να υπολογίσει την νέα τιμή ενός pixel.

Κατα την εκτέλεση του προγράμματος kernel το πρόγραμμα μας δημιουργεί ένα 3D Grid απο blocks. Οι 2 διαστάσεις του Grid διαμορφώνονται με τον ίδιο τρόπο όπως είδαμε και πριν, δηλαδή κάθε block θα υπολογίσει την νέα τιμή των pixel που του ανήκουν με την διαφορά όμως πως τα blocks πλέον δεν θα λαμβάνουν υπόψην τους μόνο την περιοχή που τους ανήκει αλλά ανάλογα με το βάθος που βρίσκονται μέσα στην τρίτη διάσταση θα φορτώνουν μία άλλη περιοχή της εικόνας και θα την συγκρίνουν με αυτήν που τους ανήκει. Το μέγεθος της τρίτης διάστασης Z είναι ίσο με το πλήθος όλων των δυνατών περιπτώσεων σύγκρισης μεταξύ περιοχών και αυτός ο αριθμός είναι ίσος με το γινόμενο των άλλων δύο διαστάσεων X,Y. Για παράδειγμα για εικόνα 256x256 δημιουργείται Grid μεγέθους (16,16,256). Η παρακάτω εικόνα θα βοηθήσει περισσότερο στην κατανόηση της παραπάνω λογικής.



Ας πάρουμε για παράδειγμα τα blocks με συντεταγμένες $(0,0)$, το κάθε ένα απο τα block φορτώνει στην shared μνήμη που του αντιστοιχεί τις τιμές των pixel της περιοχής A, θα φορτώσει όμως επιπλέον και έναν άλλο πίνακα στην shared μνήμη με τιμές απο κάποια άλλη περιοχή της εικόνας ανάλογα το βάθος στο οποίο βρίσκεται στην τρίτη διάσταση Z ($blockIdx.z$) και θα συγκρίνει πλέον τις νέες τιμές της περιοχής A με την αντίστοιχη(“δεύτερη”) περιοχή που φόρτωσε στην shared μνήμη. Δηλαδή στην ουσία πλέον όλα τα block με συντεταγμένες $blockIdx.x=0$ και $blockIdx.y=0$ θα είναι υπεύθυνα για τον υπολογισμό των τιμών της περιοχής A. Για παράδειγμα το block με συντεταγμένες $(0,0,0)$ θα φορτώσει το τμήμα A και θα το συγκρίνει με τον ίδιο το A. Το block με συντεταγμένες $(0,0,1)$ θα φορτώσει το τμήμα A και B και θα συγκρίνει το A με το B. Ομοίως το block με συντεταγμένες $(0,0,2)$ θα φορτώσει το τμήμα A και C και θα συγκρίνει το A με το C και ούτο καθεξής. Όταν μιλάμε για σύγκριση και υπολογισμούς του ενός σε σχέση με το άλλο στην ουσία αναφερόμαστε στον υπολογισμό των μερικών αθροισμάτων για τον υπολογισμό του Z_i και για τα μερικά αθροίσματα $w(i,j)*f(j)$. Δηλαδή ένα thread ενός block για ένα pixel i θα υπολογίσει την τιμή Z_i υπολογίζοντας τα μερικά αθροίσματα που προκύπτουν απο την σύγκριση των 2 περιοχών που έχουν φορτωθεί στη μνήμη shared. Στο τέλος κάθε σύγκρισης, κάθε block, για παράδειγμα με συντεταγμένες $(0,0)$ δηλαδή $blockIdx.x=0$ και $blockIdx.y=0$ της περιοχής A θα έχει υπολογίσει ένα απο τα μερικά αθροίσματα που αντιστοιχούν στον τύπο για το υπολογισμό του Z_i και έχουν να κάνουν με τα pixel μεταξύ των δύο περιοχών, τα οποία θα πρέπει να προστεθούν όλα μαζί για να προκύψει η τελική τιμή του Z_i για το pixel i κάτι το οποίο επιτυγχάνεται με την εντολή `atomicAdd`. Με την ίδια

λογική υπολογίζονται τα Z_i των pixel και για τις υπόλοιπες περιοχές. Όταν τελικά υπολογιστούν οι τιμές Z_i τότε τα block εκτελούν με την ίδια λογική και τους υπολογισμούς για τα μερικά αθροίσματα $w(i,j)*f(j)$. Στην παρακάτω εικόνα γίνεται προσπάθεια για την απεικόνιση της παραπάνω λογικής για τα pixel της περιοχής A τα οποία τα διαχειρίζονται τα blocks με συντεταγμένες blockIdx.x=0 και blockIdx.y=0.

Κάποιο pixel της περιοχής A

$$Z(i) = \sum_j e^{-\frac{\|f(N_i) - f(N_j)\|_{G(a)}^2}{\sigma^2}} + \sum_j e^{-\frac{\|f(N_i) - f(N_j)\|_{G(a)}^2}{\sigma^2}} + \sum_j e^{-\frac{\|f(N_i) - f(N_j)\|_{G(a)}^2}{\sigma^2}} + \dots$$

Block (0,0,0)

A-A

Block (0,0,1)

A-B

Block (0,0,2)

A-C

Παρατήρη για μια λεπτομέρεια στην υλοποίηση του κώδικα

Ο αλγόριθμος απαιτεί πρώτα τον συνολικό υπολογισμό της τιμής Z_i ενός pixel και αυτό δεδομένου της παραπάνω υλοποίησης δημιουργεί πρόβλημα που οφείλεται στην αδυναμία για συγχρονισμό μεταξύ των block. Όταν ένα block για παράδειγμα γράψει τα μερικά αθροίσματα της τιμής Z_i ή των αθροισμάτων $w(i,j)*f(j)$ θα πρέπει να περιμένει τα άλλα block να τελειώσουν τις εγγραφές στην τιμή Z_i ώστε να χρησιμοποιηθεί η συνολική τιμή της μεταβλητής Z_i . Για να πετύχουμε λοιπόν τον παραπάνω συγχρονισμό μέσα απο το Script Matlab υπολογίζουμε πρώτα όλες τις τιμές Z_i , έπειτα μεταφέρουμε τον έλεγχο στην cpu (για τον συγχρονισμό των blocks) και στην συνέχεια ξανά καλούμε την gpu να χρησιμοποιήσει τις τιμές Z_i ώστε να υπολογίσει εν τέλη τις νέες τιμές των pixel. Για να μην δημιουργήσουμε δύο συναρτήσεις kernel το παραπάνω επιτυγχάνεται με την κλίση του ίδιου προγράμματος kernel μόνο που ανάλογα με το τι απο τα δύο θέλουμε, δηλαδή υπολογισμό του Z_i ή των $w(i,j)*f(j)$, εισάγουμε μια επιπλέον παράμετρο flag που θα καθορίσει την επιλογή μας. Έτσι για flag=0 το πρόγραμμα kernel υπολογίζει τα Z_i ενώ για flag=1 υπολογίζει τα $w(i,j)*f(j)$.

Σχεδίαση και περιγραφή τεχνητής εισόδου, για έλεγχο ορθότητας.

---Το σειριακό πρόγραμμα καλείται απο το script του matlab ως MEX file,το script εκτελεί αυτόματα την εντολή mex και επομένως δεν χρειάζεται να εκτελέσουμε κάποια λειτουργία περαιτέρω.

---Τα παραλληλοποιημένα προγράμματα σε cuda χρειάζονται πριν την εκτέλεση του script να μεταγλωττιστούν με τον compiler nvcc της nvidia και γιαυτό τον σκοπό πρέπει πρώτα να εκτελεστεί το Makefile που βρίσκεται στο φάκελο cuda.

Το script είναι μία συνάρτηση που παίρνει ως ορίσματα :

- Το μέγεθος των patch με τιμές 3,5,7,9,... κλπ

- Την επιθυμητή τιμή σ

- Το όνομα της εικόνας που θέλουμε να αποθρουβοποιήσουμε και βρίσκεται στον φάκελο data.Το όνομα το εισάγουμε ως string δηλαδή είναι της μορφής 'ονομαεικονας'.

Επομένως για να μπορέσουμε να τρέξουμε τα προγράμματα και να ελέγξουμε την ορθότητα τους καλούμε το script-συνάρτηση με όνομα pipeline_nonLocalMeans.m με τα αντίστοιχα ορίσματα.Ένα παράδειγμα είναι της μορφής : pipeline_nonLocalMeans(3,0.04,'window64')

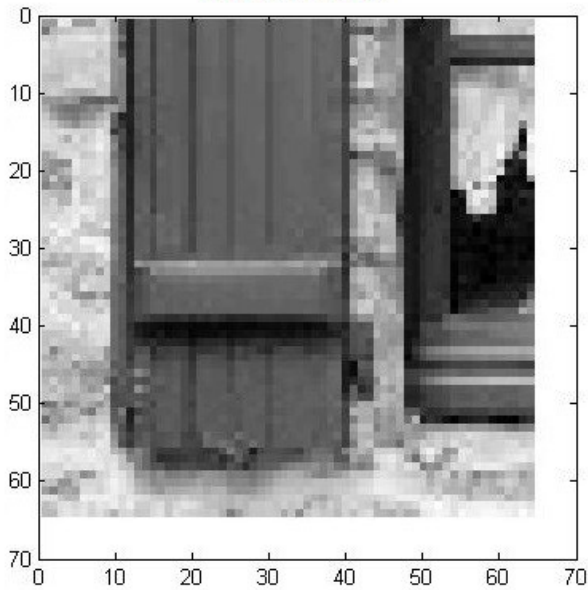
Τα προγράμματα κατά την εκτέλεση του αποθηκεύει τις εικόνες που προκύπτουν στον φάκελο images,και επομένως σε περίπτωση που η εκτέλεση του κώδικα γίνεται στο Διάδη μπορούν να κατέβουν και να απεικονιστούν στο Matlab σε κάποιο προσωπικό υπολογιστή.

Αποτελέσματα εικόνων απο το σειριακό πρόγραμμα

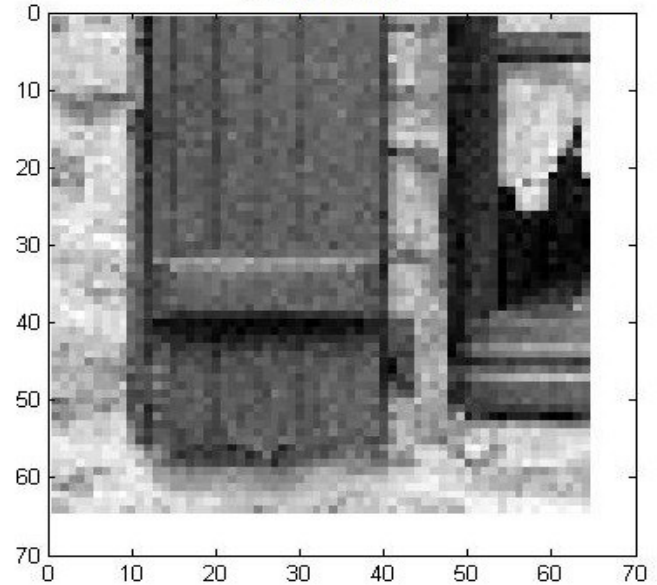
Τα παρακάτω αποτελέσματα προκύπτουν απο το σειριακό πρόγραμμα που έγραψα και όχι απο το προτεινόμενο.

Για εικόνα 64X64 pixel, patchSize=5 και sigma=0.04

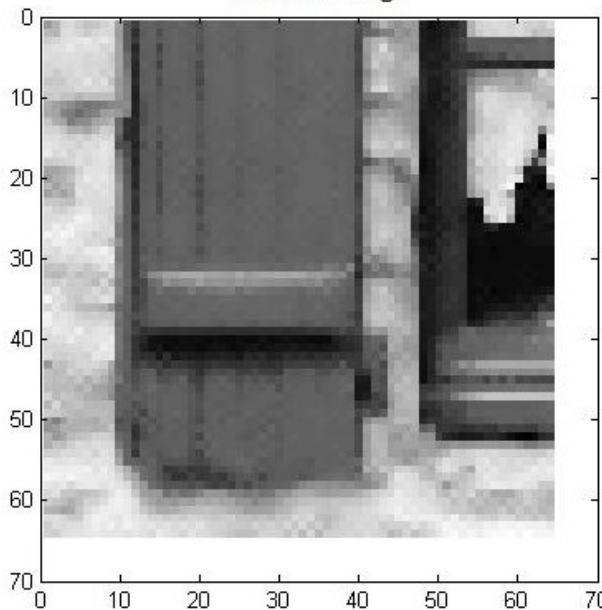
Original Image



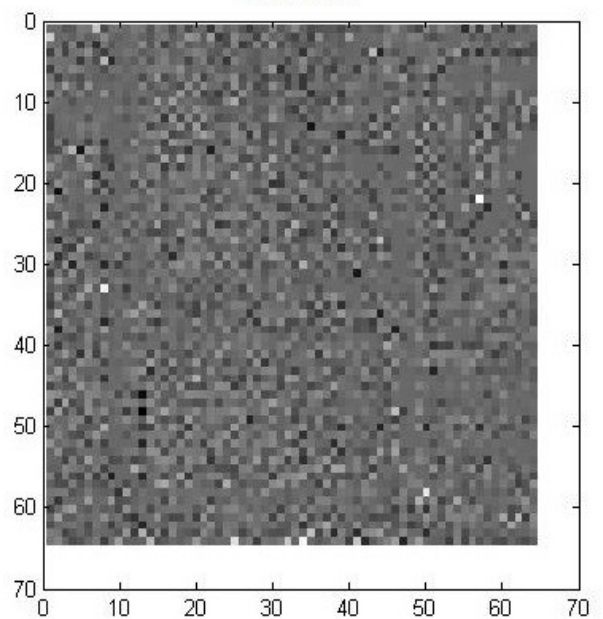
Noisy Image



Filtered Image

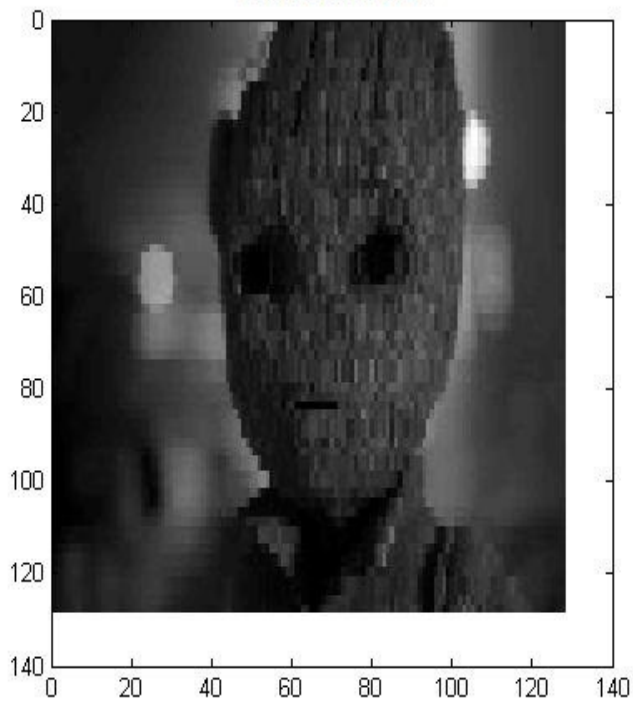


Residual

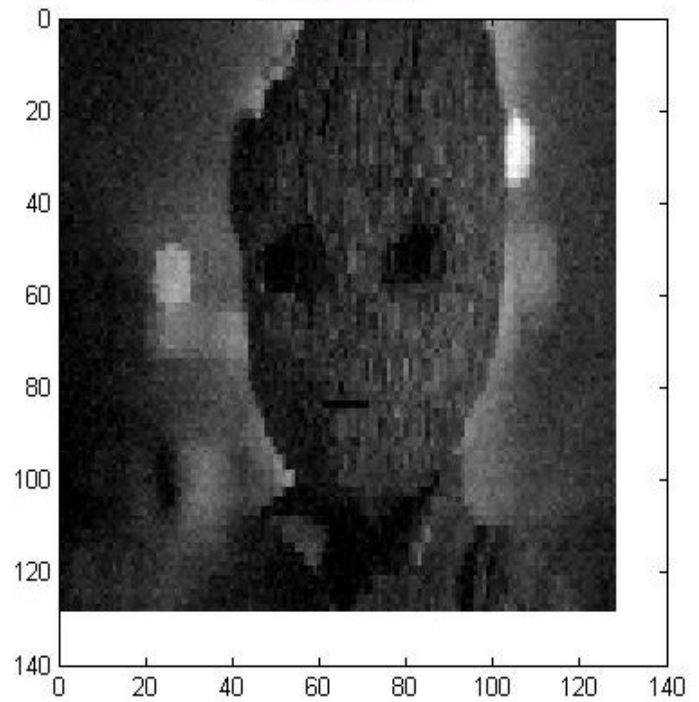


Για εικόνα 128X128 pixel, patchSize=5 και sigma=0.04

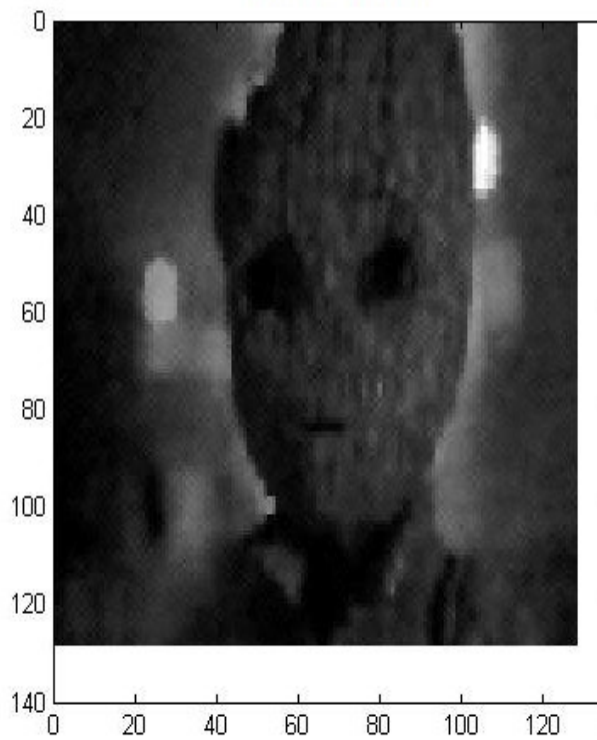
Original Image



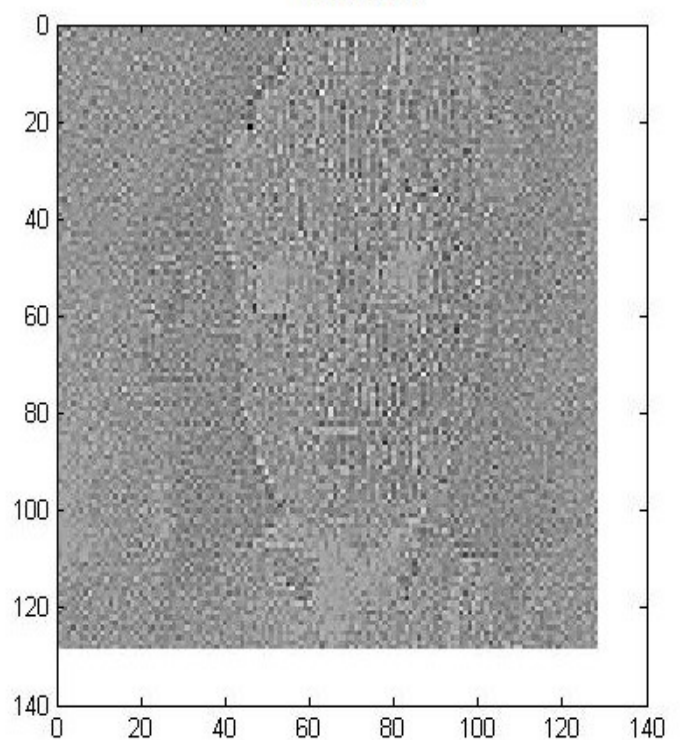
Noisy Image



Filtered Image



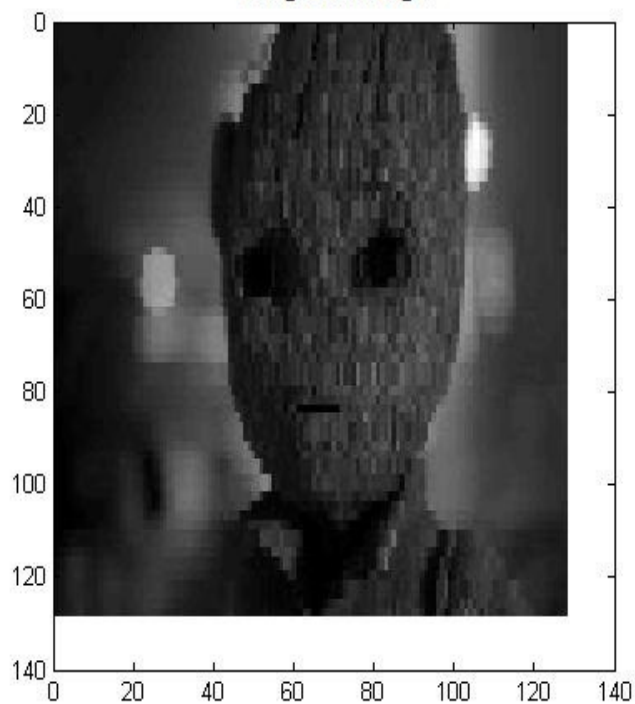
Residual



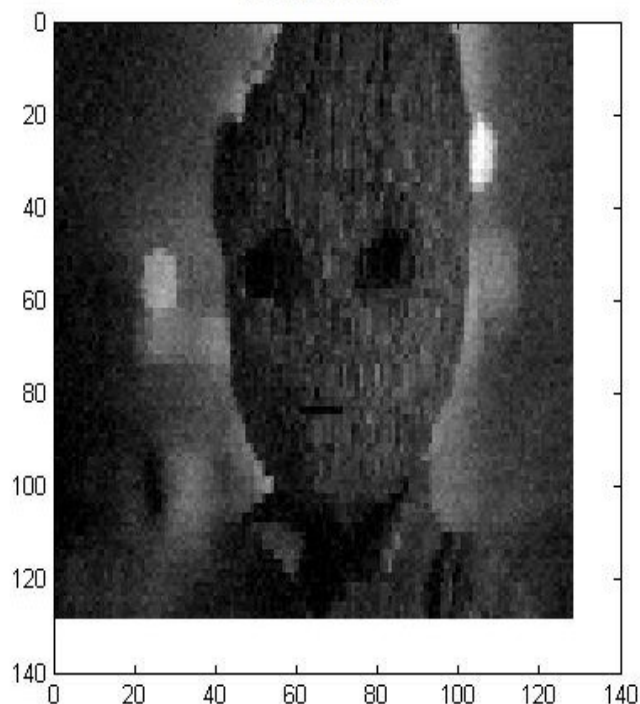
Πρώτο Παραλληλοποιημένο Πρόγραμμα

Για εικόνα 128X128 pixel, patchSize=5 και sigma=0.04

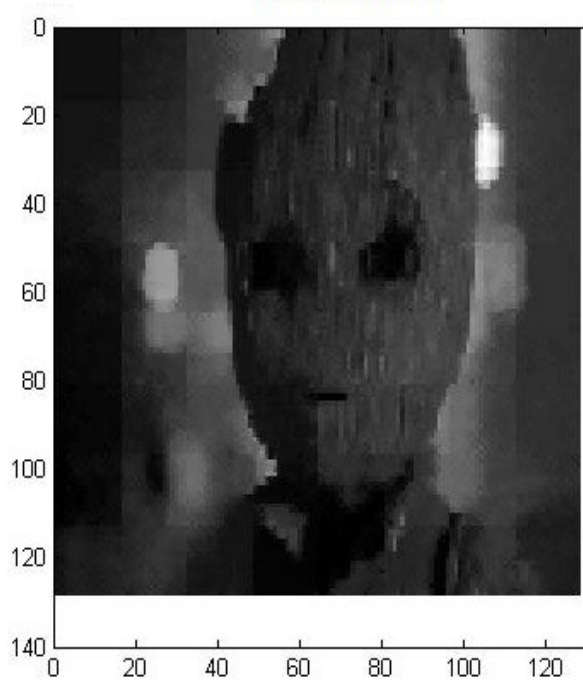
Original Image



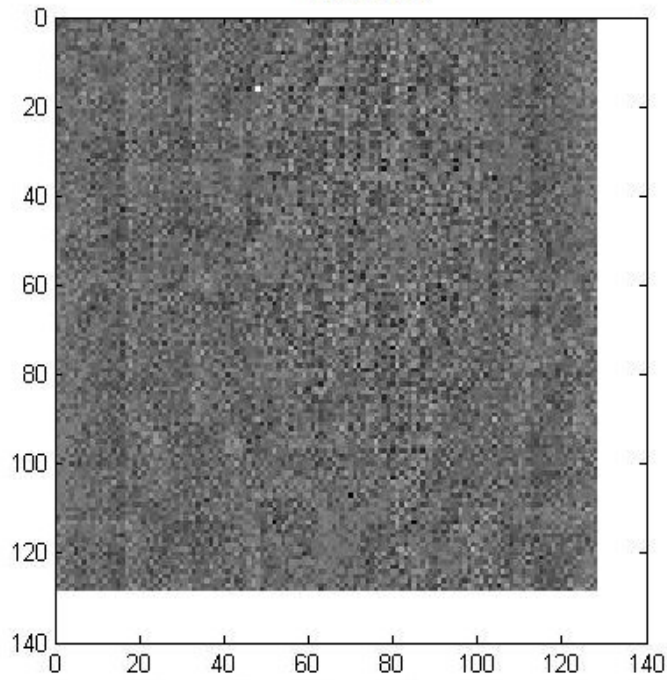
Noisy Image



Filtered Image

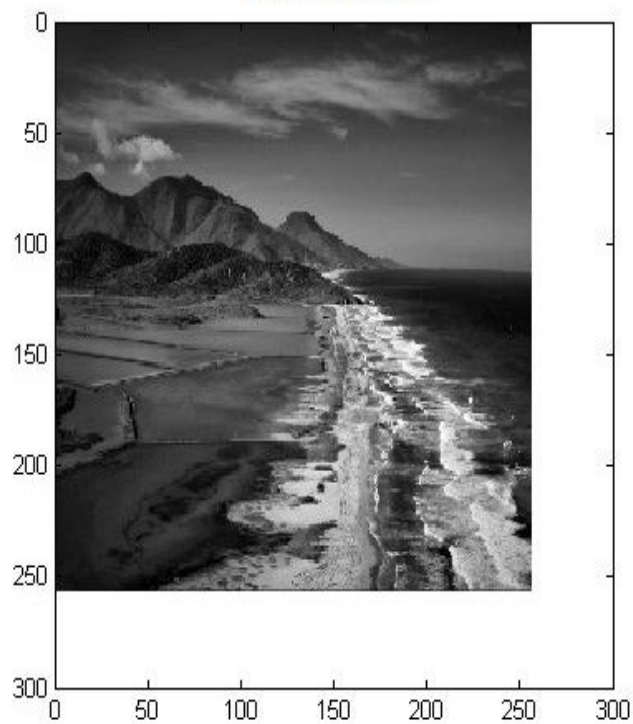


Residual

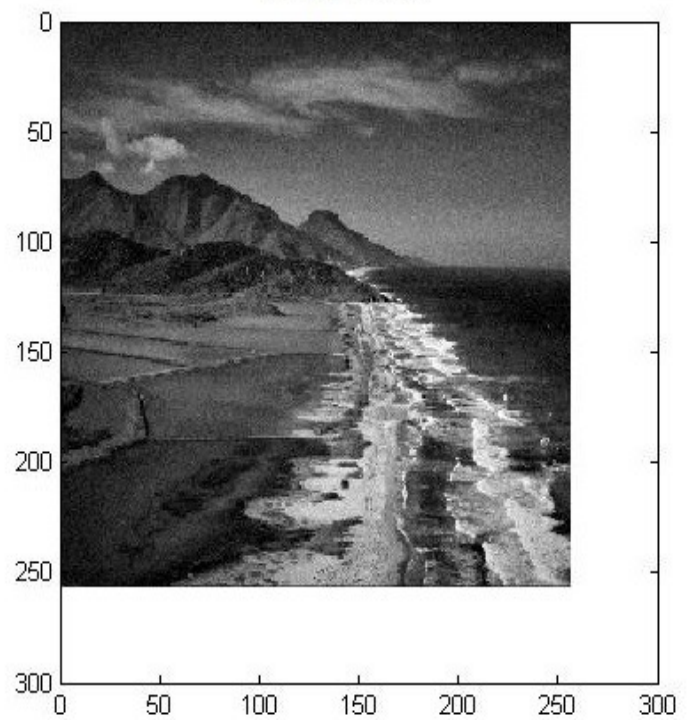


Για εικόνα 256X256 pixel, patchSize=5 και sigma=0.04

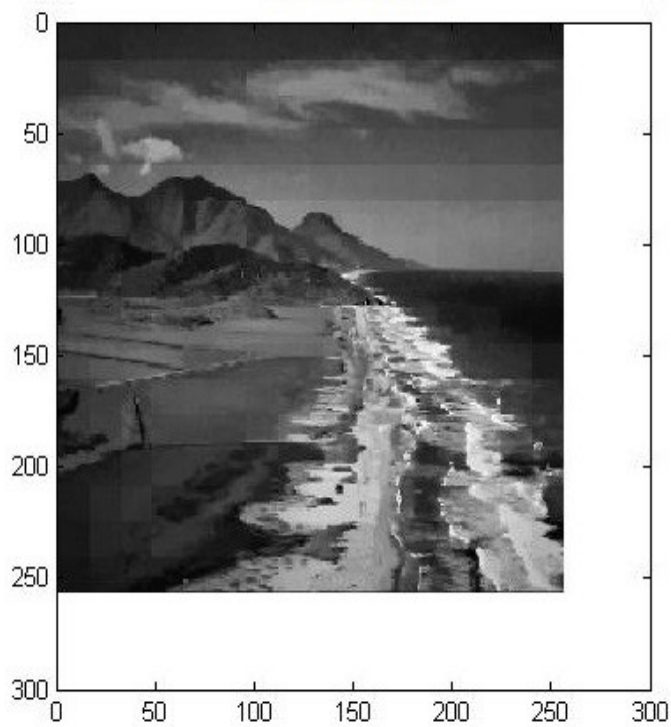
Original Image



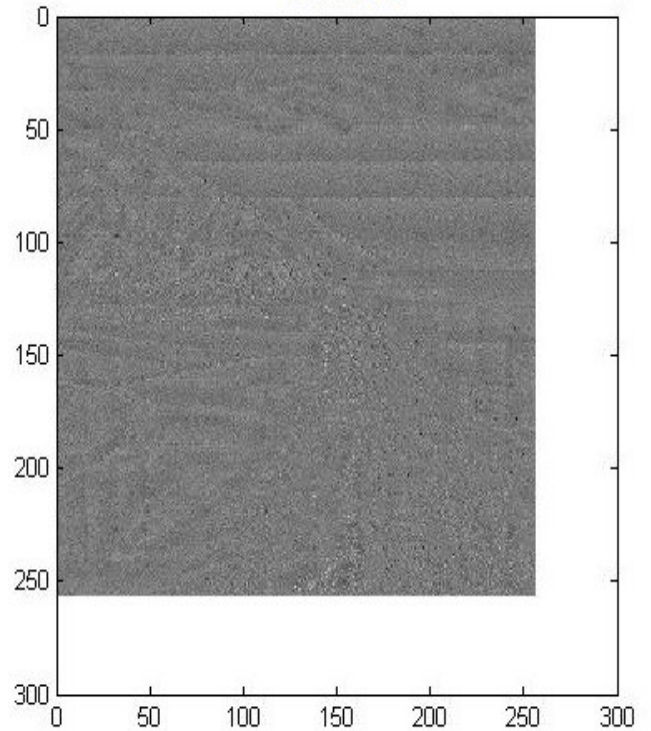
Noisy Image



Filtered Image

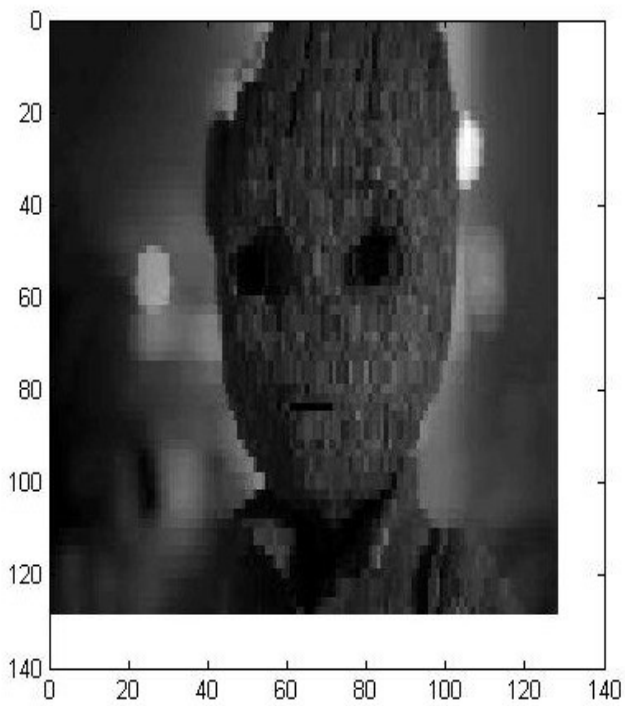


Residual

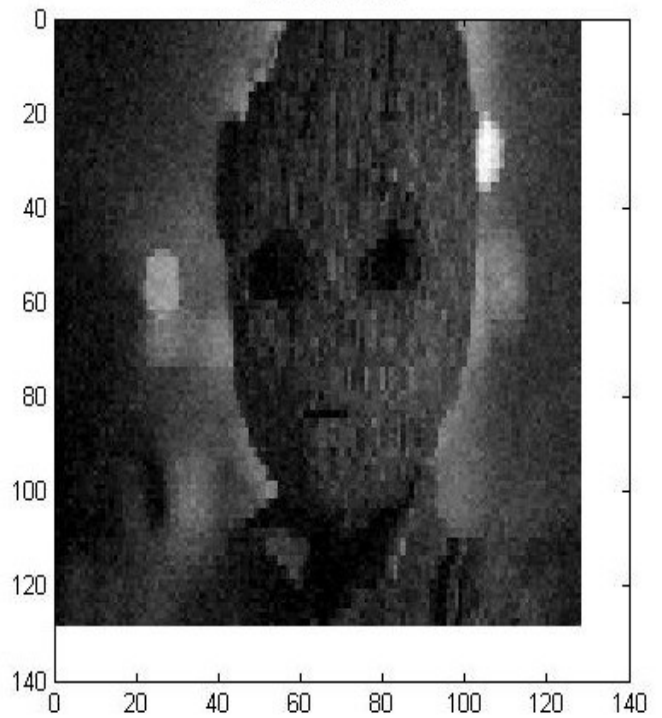


Δεύτερο Παραλληλοποιημένο Πρόγραμμα
Για εικόνα 128X128 pixel, patchSize=5 και sigma=0.04

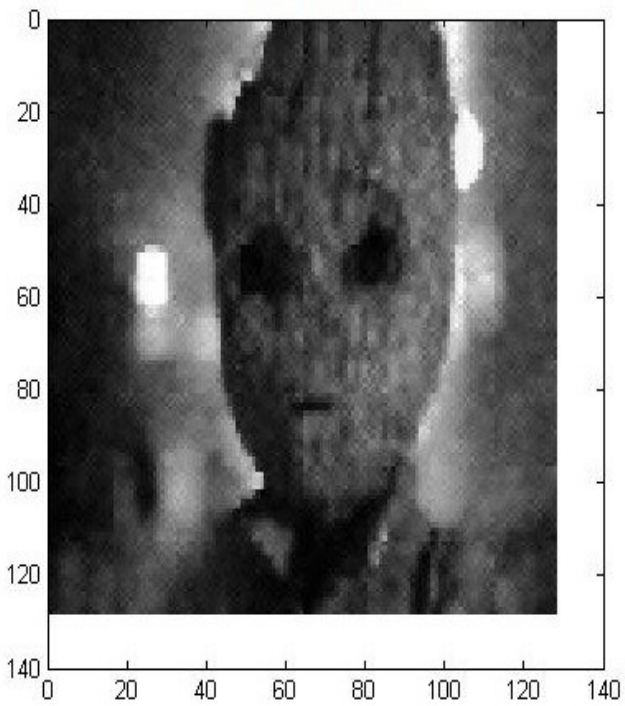
Original Image



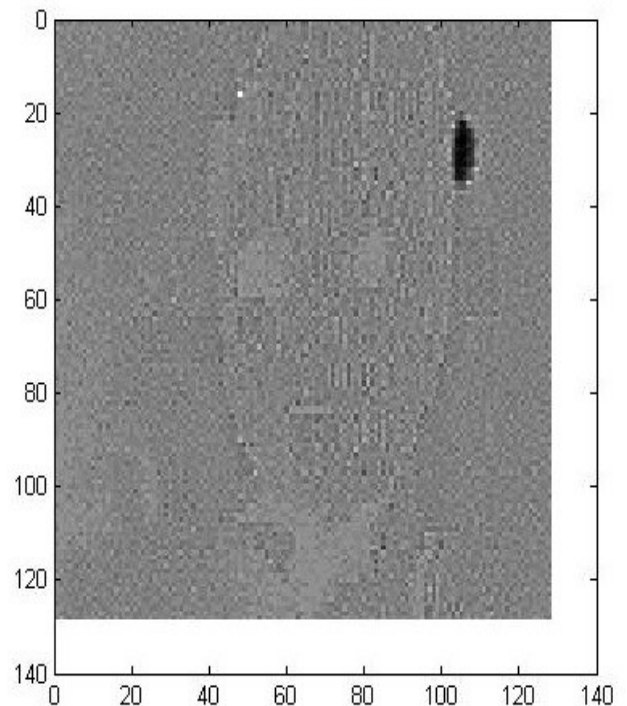
Noisy Image



Filtered Image

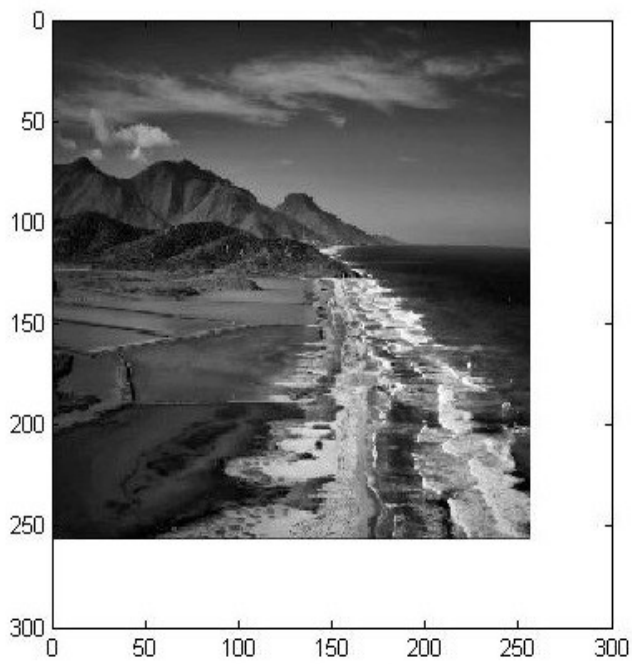


Residual

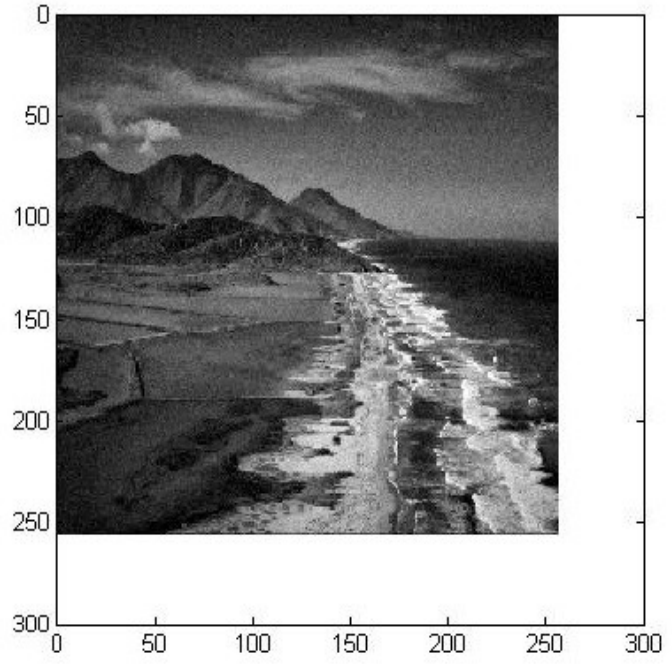


Για εικόνα 256X256 pixel, patchSize=5 και sigma=0.02

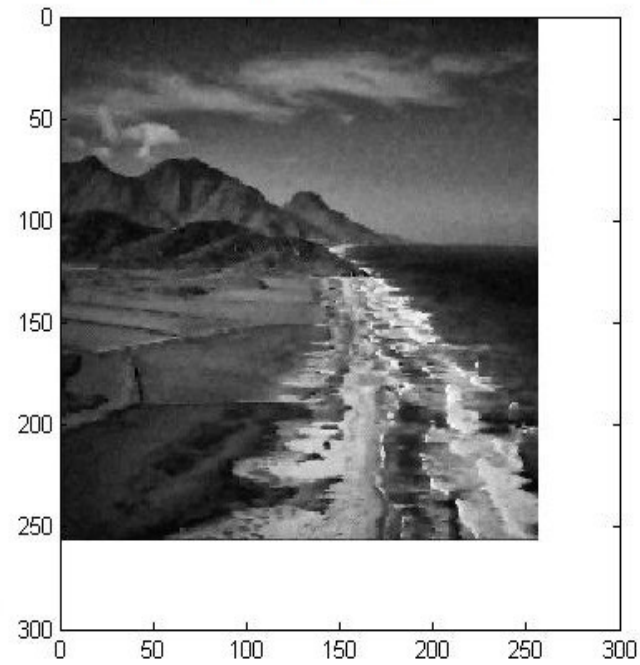
Original Image



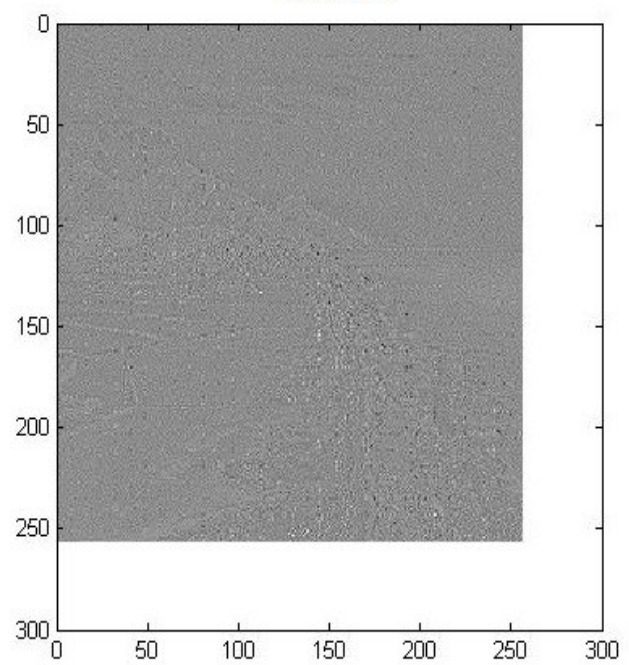
Noisy Image



Filtered Image



Residual



Οι χρόνοι εκτέλεσης των προγραμμάτων

Εικόνα 64x64	Σειριακό	Πρώτο Παραλληλοποιημένο	Δεύτερο Παραλληλοποιημένο
3x3 patch	5.97 sec	0.0044 sec	0.031 sec
5x5 patch	11.53 sec	0.0098 sec	0.074 sec
7x7 patch	20.17 sec	0.017sec	0.137sec

Εικόνα 128x128	Σειριακό	Πρώτο Παραλληλοποιημένο	Δεύτερο Παραλληλοποιημένο
3x3 patch	95 sec	0.011 sec	0.45 sec
5x5 patch	183 sec	0.024 sec	1.12 sec
7x7 patch	311 sec	0.044 sec	2.08 sec

Εικόνα 256x256	Σειριακό	Πρώτο Παραλληλοποιημένο	Δεύτερο Παραλληλοποιημένο
3x3 patch	– sec	0.033 sec	7.67 sec
5x5 patch	– sec	0.079 sec	18.21 sec
7x7 patch	– sec	0.14 sec	33.52 sec

Σχόλια και παρατηρήσεις

Στους παραπάνω πίνακες φαίνονται τα αποτελέσματα απο την εκτέλεση των προγραμμάτων για μεγέθοι εικόνων 64x64,128x128 και 256x256 και μεγέθοι των patches 3x3 5x5 και 7x7. Αρχικά αυτό που παρατηρούμε είναι η πολύ μεγάλη άυξηση στον χρόνο εκτέλεσης του σειριακού προγράμματος καθώς το μέγεθος εικόνας και patch αυξάνονται. Για μέγεθος εικόνων 256x256 οι χρόνοι ήταν αρκετά μεγάλη για να μετρηθούν, παρόλα αυτά η σύγκριση στις αποδόσεις των προγραμμάτων είναι εμφανής. Για την αποφυγή κάποιου μπερδέματος θα υπενθυμίσω ότι οι χρόνοι προκύπτουν απο το σειριακό πρόγραμμα που δημιούργησα εγώ και όχι απο το προτεινόμενο που δόθηκε μαζί με την εκφώνηση της εργασία. Οι χρόνοι και απο τα 2 παραλληλοποιημένα προγράμματα είναι με πολύ μεγάλη διαφορά καλύτεροι απο το σειριακό πρόγραμμα. Το πρώτο πρόγραμμα το οποίο διαχειρίζεται την εικόνα σαν σύνολο υπο-εικόνων

βλέπουμε ότι έχει πολύ μικρούς χρόνους εκτέλεσης αφού στην ουσία όσο και να αυξάνουμε τον αριθμό των pixel δημιουργούνται επιπλέον block τα οποία καλύπτουν τους υπολογισμούς χωρίς να επηρεάζονται απο το μέγεθος της εικόνας εφόσον κάθε block θα εκτελέσει τον αλγόριθμο σε μία περιοχή $16 \times 16 = 256$ pixels και κάθε thread του μπλόκ θα υπολογίσει την τιμή ενός απο τα pixel του block. Στην ουσία αυτο που γίνεται είναι ότι για κάθε επιπλέον pixel που προκύπτει απο μία μεγαλύτερη εικόνα παρέχουμε επιπλέον και περισσότερους υλικούς πόρους (επεξεργαστές) που αναλαμβάνουν να καλύψουν τους υπολόγισμούς, φυσικά υπάρχει ένα όριο για τους πόρους που μπορούν χρησιμοποιηθούν ατομικά απο ένα νήμα. Οι χρόνοι παρόλα αυτά έχουν μία μικρή (σχετικά) αύξηση επειδή όσο περισσότερα νήματα-block δημιουργούμε γίνεται πιο δύσκολος ο συγχρονισμός τους και επιπλέον υπάρχουν καθυστερήσεις απο την προσπάθεια για ταυτόχρονη προσπέλαση της global μνήμης απο πολλά νήματα. Το δεύτερο παραλληλοποιημένο πρόγραμμα έχει πολύ καλύτερους χρόνους απο το σειριακό πρόγραμμα, όπως βλέπουμε όμως είναι αρκετά μεγαλύτεροι σε σχέση με το προηγούμενο παραλληλοποιημένο πρόγραμμα. Η διαφορά αυτή οφείλεται στο ότι το δεύτερο πρόγραμμα αφού λαμβάνει υποψη όλες τις τιμές των pixel δημιουργεί πολύ περισσότερα thread-block τα οποία δεν είναι δυνατόν να είναι όλα ενεργά ταυτόχρονα, δημιουργούνται επιπλέον καθυστερήσεις απο τον τεράστιο αριθμό νημάτων που εκτελούν ταυτόχρονες προσπελάσεις στην μνήμη και τέλος έχουμε καθυστερήσεις και απο τον συγχρονισμό μεταξύ των block με την εντολή `atomicadd()` για την ατομική ανανέωση κοινών θέσεων μνήμης. Τα αποτελέσματα της αποθορυβοποίησης των εικόνων απο τα προγράμματα και στις δύο περιπτώσεις είναι πολύ καλά, με την παρατήρηση ότι στην περίπτωση του πρώτου παραλληλοποιημένου λόγω του μικρού δείγματος pixel για κάθε block βλέπουμε ότι στα όρια των block δημιουργούνται στην εικόνα κάποιες γραμμές. Αυτό αποφεύγεται στην δεύτερη υλοποίηση αφού πλέον λαμβάνουμε υπόψη όλα τα pixel της εικόνας. Τέλος βλέπουμε απο τις εικόνες Residual ότι η αποθορυβοποίηση δεν καταστρέφει την πληροφορία της αρχικής εικόνας αφού στις εικόνες Residual φαίνεται ο θόρυβος που έχει αφαιρεθεί απο την αρχική εικόνα ενώ η πληροφορία της αρχικής εικόνας φαίνεται αμυδρά έως καθόλου.

Τελος