



## ΣΥΣΤΗΜΑΤΑ ΠΑΡΑΛΛΗΛΗΣ ΕΠΕΞΕΡΓΑΣΙΑΣ

### Άσκηση 2: Παράλληλη επίλυση εξίσωσης θερμότητας

Χειμερινό εξάμηνο 2019-20 - Ροή Υ

Αντωνιάδης, Παναγιώτης  
el15009@central.ntua.gr

Μπαζώτης, Νικόλαος  
el15739@central.ntua.gr

---

*"Everybody who learns concurrency thinks they understand it, ends up finding mysterious races they thought weren't possible, and discovers that they didn't actually understand it yet after all."*

– Herb Sutter, *chair of the ISO C++ standards committee*

# 1 Περιγραφή προβλήματος

Η εξίσωση Laplace, γνωστή και ως εξίσωση θερμότητας, περιγράφει την κατανομή της θερμότητας σε μια δοσμένη περιοχή καθώς μεταβάλλεται ο χρόνος. Ορίζεται από την παρακάτω διαφορική εξίσωση:

$$\left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) u(x, y) = 0$$

Για την επίλυση του προβλήματος της διάδοσης θερμότητας σε δύο διαστάσεις, χρησιμοποιούνται τρεις υπολογιστικοί πυρήνες, οι οποίοι αποτελούν ευρέως διαδεδομένη δομική μονάδα για την επίλυση μερικών διαφορικών εξισώσεων: η μέθοδος Jacobi, η μέθοδος Gauss-Seidel με Successive OverRelaxation και η μέθοδος Red-Black SOR, που πραγματοποιεί Red-Black ordering στα στοιχεία του υπολογιστικού χωρίου και συνδυάζει τις δύο προηγούμενες μεθόδους.

## 1.1 Μέθοδος Jacobi

Στην μέθοδο αυτή, η τιμή ενός στοιχείου μια χρονική στιγμή ισούται με την μέση τιμή των γειτόνων του την προηγούμενη χρονική στιγμή, δηλαδή για το στοιχείο  $(x, y)$  την χρονική στιγμή  $t + 1$  έχουμε:

$$u_{x,y}^{t+1} = \frac{u_{x-1,y}^t + u_{x,y-1}^t + u_{x+1,y}^t + u_{x,y+1}^t}{4}$$

Από άποψη υλοποίησης, απαιτούνται τουλάχιστον δύο πίνακες, ένας για την τρέχουσα  $(t+1)$  κι ένας για την προηγούμενη  $(t)$  χρονική στιγμή. Στο τέλος κάθε βήματος, πραγματοποιείται έλεγχος σύγκλισης.

```
1 for (t = 0; t < T && !converged; t++) {  
2   for (i = 1; i < X - 1; i++)  
3     for (j = 1; j < Y - 1; j++)  
4       U[t+1][i][j] = (1/4) * (U[t][i-1][j] + U[t][i+1][j] +  
5         U[t][i][j-1] + U[t][i][j+1]);  
6   converged = check_convergence(U[t+1], U[t]);  
7 }  
8
```

## 1.2 Μέθοδος Gauss-Seidel SOR

Στην μέθοδο αυτή, προσπαθούμε να αντιμετωπίσουμε το γεγονός ότι η μέθοδος Jacobi έχει πολύ αργό ρυθμό σύγκλισης. Συγκεκριμένα, όταν υπολογίζεται το σημείο  $U[t+1][i][j]$  έχουν ήδη υπολογιστεί τα σημεία  $U[t+1][i-1][j]$  και  $U[t+1][i][j-1]$ . Συνεπώς, στην ανανέωση μπορούν να χρησιμοποιηθούν οι νέες τιμές των δύο αυτών γειτόνων, το οποίο θα επιταχύνει την όλη διαδικασία. Έτσι, προκύπτει ο παρακάτω τύπος για την μέθοδο Gauss-Seidel:

$$u_{x,y}^{t+1} = \frac{u_{x-1,y}^{t+1} + u_{x,y-1}^{t+1} + u_{x+1,y}^t + u_{x,y+1}^t}{4}$$

Για ακόμη μεγαλύτερη σύγκλιση, χρησιμοποιείται η τεχνική Successive Over-Relaxation (SOR), η οποία χρησιμοποιεί και την προηγούμενη τιμή του ίδιου του στοιχείου και καταλήγει στην εξής σχέση:

$$u_{x,y}^{t+1} = u_{x,y}^t + \omega \frac{u_{x-1,y}^{t+1} + u_{x,y-1}^{t+1} + u_{x+1,y}^t + u_{x,y+1}^t - 4u_{x,y}^t}{4}, \quad \omega \in (0, 2)$$

```

1 for (t = 0; t < T && !converged; t++) {
2   for (i = 1; i < X - 1; i++)
3     for (j = 1; j < Y - 1; j++)
4       U[t+1][i][j] = U[t][i][j] + (omega / 4) * (U[t+1][i-1][j] + U[t+1][i][j-1]
5         + U[t][i+1][j] + U[t][i][j+1] - 4*U[t][i][j]);
6   converged = check_convergence(U[t+1], U[t]);
7 }
8

```

### 1.3 Μέθοδος Red-Black SOR

Η μέθοδος αυτή είναι παρόμοια με την μέθοδο Gauss-Seidel SOR, με την διαφορά ότι χωρίζει τα στοιχεία σε δύο ομάδες, στα κόκκινα και στα μαύρα. Η ανανέωση γίνεται σε 2 φάσεις: στη κόκκινη φάση υπολογίζονται τα κόκκινα στοιχεία από τα μαύρα, ενώ στη μαύρη φάση υπολογίζονται τα μαύρα στοιχεία από τα κόκκινα. Έτσι, πετυχαίνουμε μεγαλύτερη ταχύτητα σύγκλισης.

$$u_{x,y}^{t+1} = u_{x,y}^t + \omega \frac{u_{x-1,y}^t + u_{x,y-1}^t + u_{x+1,y}^t + u_{x,y+1}^t - 4u_{x,y}^t}{4}, \text{ when } (x+y)\%2 = 0$$

$$u_{x,y}^{t+1} = u_{x,y}^t + \omega \frac{u_{x-1,y}^{t+1} + u_{x,y-1}^{t+1} + u_{x+1,y}^{t+1} + u_{x,y+1}^{t+1} - 4u_{x,y}^t}{4}, \text{ when } (x+y)\%2 = 1$$

```

1 for (t = 0; t < T && !converged; t++) {
2   //Red phase
3   for (i = 1; i < X - 1; i++)
4     for (j = 1; j < Y - 1; j++)
5       if ((i+j)%2==0)
6         U[t+1][i][j] = U[t][i][j] + (omega / 4) * (U[t][i-1][j] + U[t][i][j-1]
7           + U[t][i+1][j] + U[t][i][j+1] - 4*U[t][i][j]);
8   //Black phase
9   for (i = 1; i < X - 1; i++)
10    for (j = 1; j < Y - 1; j++)
11      if ((i+j)%2==1)
12        U[t+1][i][j] = U[t][i][j] + (omega / 4) * (U[t+1][i-1][j] + U[t+1][i][j-1]
13          + U[t+1][i+1][j] + U[t+1][i][j+1] - 4*U[t][i][j]);
14   converged = check_convergence(U[t+1], U[t]);
15 }
16

```

## 2 Σχεδιασμός Παραλληλοποίησης

Στην γενική περίπτωση, ο πίνακας  $U$  είναι διάστασης  $N \times M$ . Τον χωρίζουμε σε ορθογώνια block μεγέθους  $n \times m$  και υποθέτουμε ότι  $N \bmod n = 0$ ,  $M \bmod m = 0$  (αν όχι ορίζεται κατάλληλο padding). Για τον σχεδιασμό των task graph, θεωρούμε  $\frac{N}{n} = \frac{M}{m} = 4$ . Άρα ο πίνακας χωρίζεται στα παρακάτω blocks:

$A_{00}$	$A_{01}$	$A_{02}$	$A_{03}$
$A_{10}$	$A_{11}$	$A_{12}$	$A_{13}$
$A_{20}$	$A_{21}$	$A_{22}$	$A_{23}$
$A_{30}$	$A_{31}$	$A_{32}$	$A_{33}$

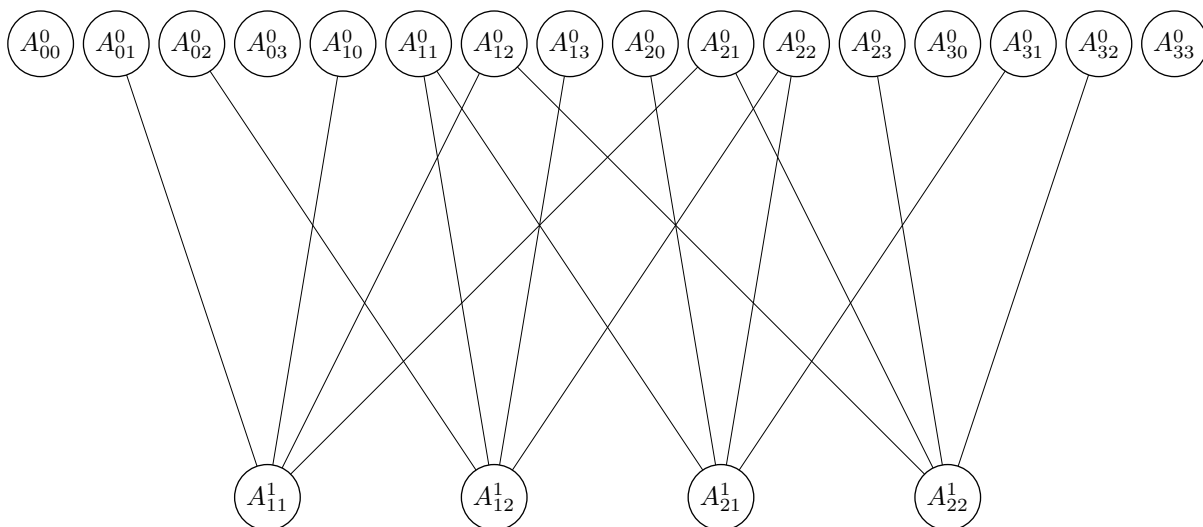
### 2.1 Μέθοδος Jacobi

#### 1. Κατανομή Υπολογισμών

Θεωρούμε, ως task τον υπολογισμό της θερμότητας για ένα block του διδιάστατου πλέγματος για μία χρονική στιγμή.

#### 2. Ορισμός ορθής σειράς εκτέλεσης μέσω task graph

Οι εξαρτήσεις που υπάρχουν μεταξύ των διαφορετικών tasks παρουσιάζονται για το συγκεκριμένο παράδειγμα στο παρακάτω task graph, όπου ορίζονται οι δύο πρώτες χρονικές στιγμές μόνο. Ουσιαστικά, κάθε χρονική στιγμή για να γίνει ο υπολογισμός ενός block χρειάζεται τις προηγούμενες τιμές από τα γειτονικά του μόνο blocks.



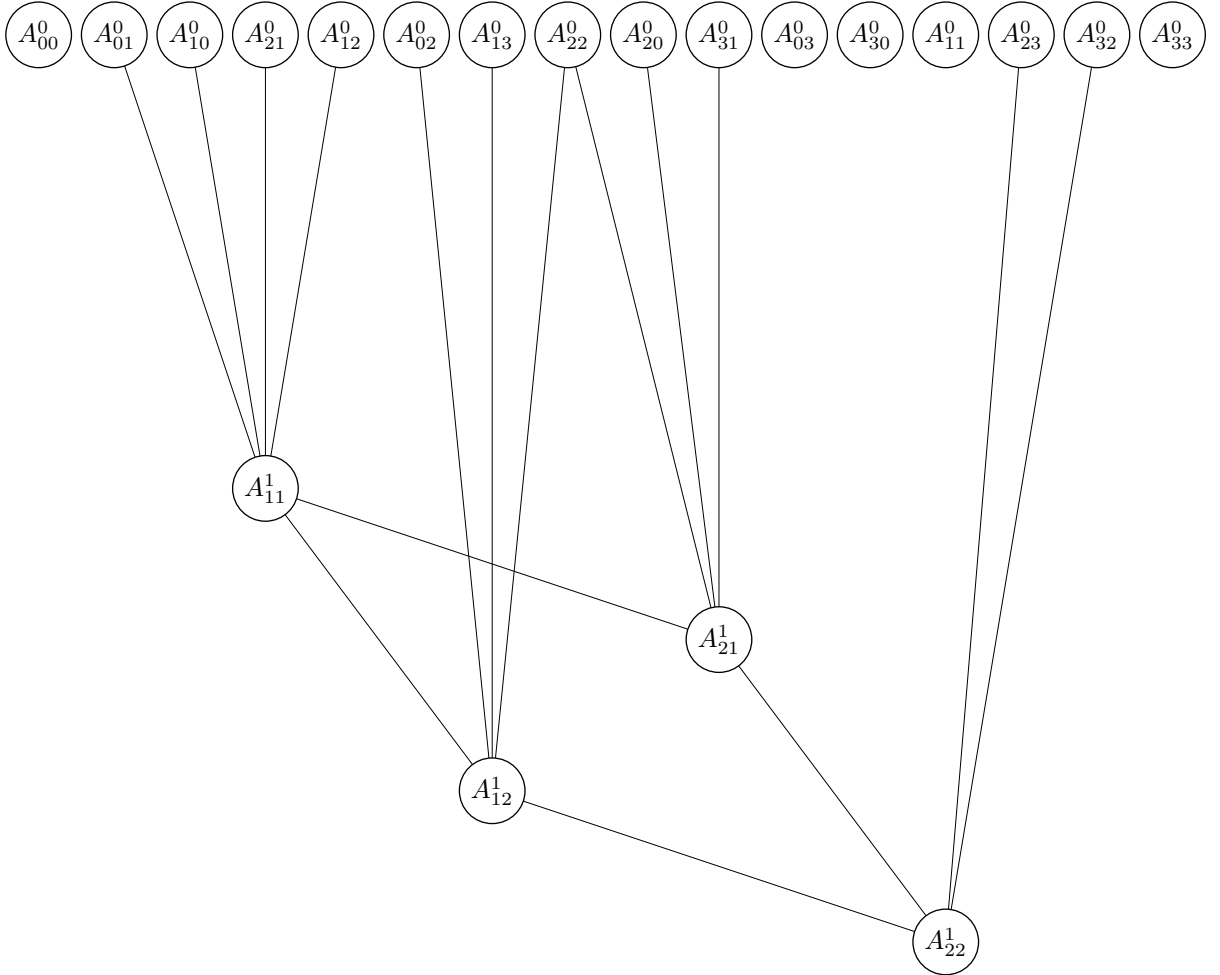
### 2.2 Μέθοδος Gauss-Seidel SOR

#### 1. Κατανομή Υπολογισμών

Αντίστοιχα με παραπάνω, θεωρούμε ως task τον υπολογισμό της θερμότητας για ένα block του διδιάστατου πλέγματος για μία χρονική στιγμή.

## 2. Ορισμός ορθής σειράς εκτέλεσης μέσω task graph

Οι εξαρτήσεις εδώ διαφέρουν από αυτές του Jacobi. Κάθε χρονική στιγμή για να γίνει ο υπολογισμός ενός block χρειάζεται τις προηγούμενες τιμές από το δεξιά, το από κάτω και το ίδιο block και τις ανανεωμένες τιμές του από πάνω και αριστερά block. Έτσι, με βάση το παράδειγμα παραπάνω έχουμε το εξής task graph:



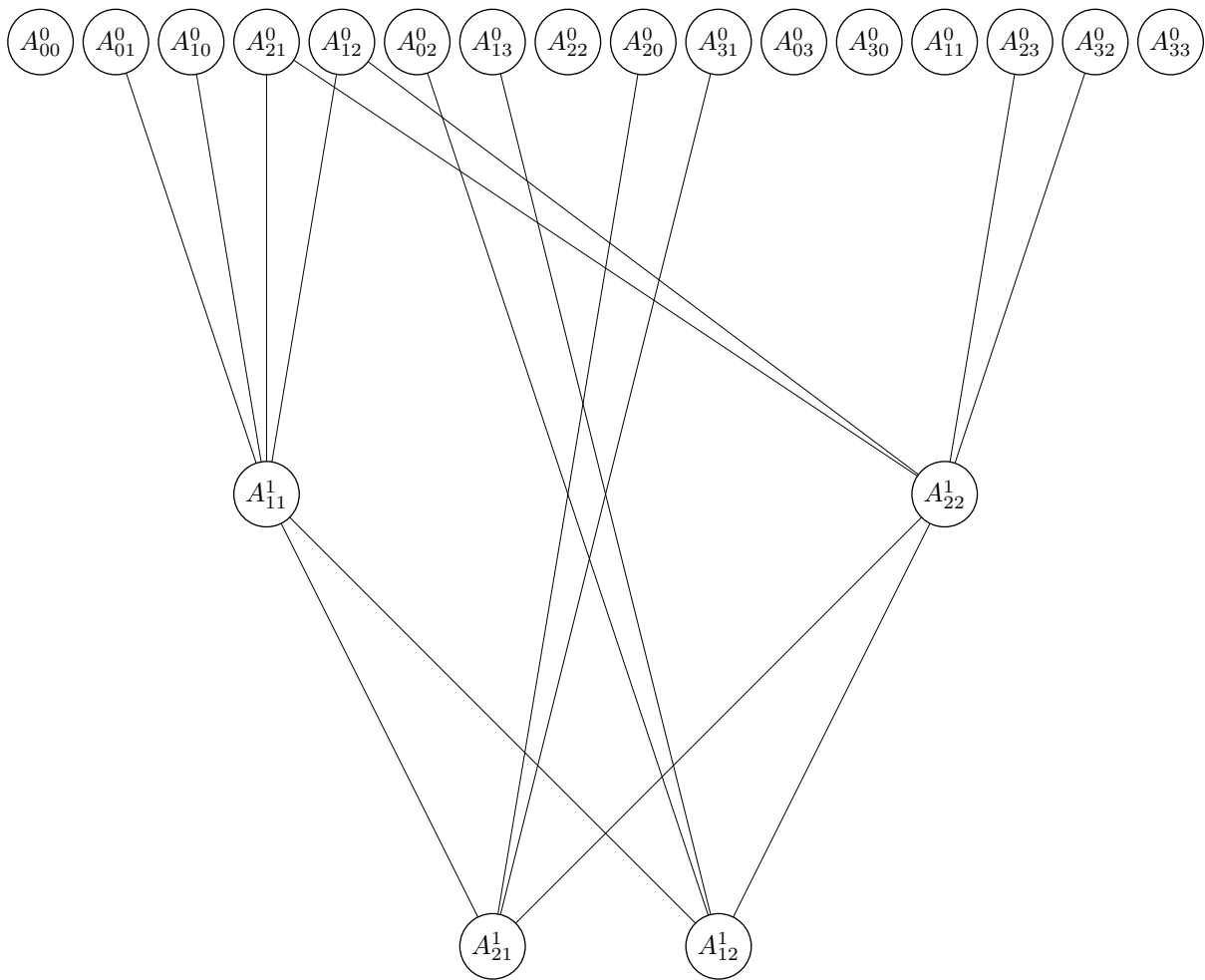
## 2.3 Μέθοδος Red-Black SOR

### 1. Κατανομή Υπολογισμών

Αντίστοιχα με παραπάνω, θεωρούμε ως task τον υπολογισμό της θερμότητας για ένα block του δισδιάστατου πλέγματος για μία χρονική στιγμή.

### 2. Ορισμός ορθής σειράς εκτέλεσης μέσω task graph

Ακολουθεί το task graph με βάση το αρχικό παράδειγμα που έχουμε:



### 3 Ανάπτυξη παράλληλων προγραμμάτων

Σε κάθε περίπτωση, ο αρχικός πίνακας  $U$  χωρίζεται σε block και κάθε διεργασία αναλαμβάνει από ένα block. Για κάθε επανάληψη, οι διεργασίες υπολογίζουν τις νέες τιμές τους αφού επικοινωνήσουν μεταξύ τους για να ανταλλάξουν δεδομένα. Έχουμε σύγκλιση όταν όλα τα block έχουν συγκλίνει στις τελικές τιμές τους. Στο τέλος, οι τιμές των επιμέρους block συνενώνονται στον αρχικό πίνακα. Ακολουθούν ορισμένες διαδικασίες που είναι κοινές και για τις 3 υλοποιήσεις.

1. Δημιουργία 2D καρτεσιανού communicator για την ευκολότερη διαχείριση των διεργασιών

```
1 MPI_Comm CART_COMM; //CART_COMM: the new 2D-cartesian communicator
2 int periods[2]={0,0}; //periods={0,0}: the 2D-grid is non-periodic
3 int rank_grid[2]; //rank_grid: the position of each process on the new communicator
4
5 MPI_Cart_create(MPI_COMM_WORLD,2,grid,periods,0,&CART_COMM); //communicator creation
6 MPI_Cart_coords(CART_COMM,rank,2,rank_grid); //rank mapping on the new communicator
```

2. Υπολογισμός διαστάσεων του επιμέρους πίνακα κάθε block. Σε περίπτωση που ο διαμοιρασμός των block δεν χωράει ακριβώς κάνουμε padding στον αρχικό πίνακα.

```
1 for (i=0;i<2;i++) {
2     if (global[i]%grid[i]==0) {
3         local[i]=global[i]/grid[i];
4         global_padded[i]=global[i];
5     }
6     else {
7         local[i]=(global[i]/grid[i])+1;
8         global_padded[i]=local[i]*grid[i];
9     }
10 }
```

3. Διαμοιρασμός του αρχικού global πίνακα  $U$  στον local πίνακα κάθε διεργασίας

```
1 //-----Datatype definition for the 2D-subdomain on the global matrix-----//
2 MPI_Datatype global_block;
3 MPI_Type_vector(local[0],local[1],global_padded[1],MPI_DOUBLE,&dummy);
4 MPI_Type_create_resized(dummy,0,sizeof(double),&global_block);
5 MPI_Type_commit(&global_block);
6
7 //-----Datatype definition for the 2D-subdomain on the local matrix-----//
8 MPI_Datatype local_block;
9 MPI_Type_vector(local[0],local[1],local[1]+2,MPI_DOUBLE,&dummy);
10 MPI_Type_create_resized(dummy,0,sizeof(double),&local_block);
11 MPI_Type_commit(&local_block);
12
13 //-----Rank 0 defines positions and counts of local blocks (2D-subdomains) on global
14 //matrix-----//
15 int * scatteroffset, * scattercounts;
16 if (rank==0) {
17     U_start = &(U[0][0]);
18     scatteroffset=(int*)malloc(size*sizeof(int));
19     scattercounts=(int*)malloc(size*sizeof(int));
20     for (i=0;i<grid[0];i++)
21         for (j=0;j<grid[1];j++) {
22             scattercounts[i*grid[1]+j]=1;
23             scatteroffset[i*grid[1]+j]=(local[0]*local[1]*grid[1]*i+local[1]*j);
24         }
```

```

23     }
24 }
25
26 //-----Rank 0 scatters the global matrix-----//
27 MPI_Scatterv(U_start, scattercounts, scatteroffset, global_block, &(u_previous[1][1]),
28             1, local_block, 0, MPI_COMM_WORLD);
29 MPI_Scatterv(U_start, scattercounts, scatteroffset, global_block, &(u_current[1][1]),
30             1, local_block, 0, MPI_COMM_WORLD);

```

4. Ορισμός μιας δομής column, που θα διευκολύνει την μεταφορά μεταξύ διεργασιών μιας στήλης ενός πίνακα.

```

1  /* Define a datatype that corresponds to a column of a local block */
2  MPI_Datatype column;
3  MPI_Type_vector(local[0], 1, local[1]+2, MPI_DOUBLE, &dummy);
4  MPI_Type_create_resized(dummy, 0, sizeof(double), &column);
5  MPI_Type_commit(&column);

```

5. Υπολογισμός για κάθε διεργασία του rank των γειτόνων του ώστε να μπορεί να επικοινωνήσει μαζί τους.

```

1  int north, south, east, west;
2  MPI_Cart_shift(CART_COMM, 0, 1, &north, &south);
3  MPI_Cart_shift(CART_COMM, 1, 1, &west, &east);

```

6. Ορισμός του range πάνω στο οποίο η διεργασία θα ανανεώνει τον local πίνακά της. Εδώ, πρέπει να λάβουμε υπόψην τόσο τα ghost cell που έχουν τοποθετηθεί για την διευκόλυνση της επικοινωνίας όσο και τα πιθανά padding cells.

```

1  int i_min, i_max, j_min, j_max;
2
3  /* internal process (ghost cell only) */
4  i_min = 1;
5  i_max = local[0] + 1;
6
7  /* boundary process - no possible padding */
8  if (north == MPI_PROC_NULL) {
9      i_min = 2; // ghost cell + boundary
10 }
11
12 /* boundary process and padded global array */
13 if (south == MPI_PROC_NULL) {
14     i_max = (global_padded[0] - global[0]) + 1;
15 }
16
17 /* internal process (ghost cell only) */
18 j_min = 1;
19 j_max = local[1] + 1;
20
21 /* boundary process - no possible padding */
22 if (west == MPI_PROC_NULL) {
23     j_min = 2; //ghost cell + boundary
24 }
25
26 /* boundary process and padded global array */
27 if (east == MPI_PROC_NULL) {

```



```

28 j_max == (global_padded[1] - global[1]) + 1;
29 }

```

7. Έλεγχος σύγκλισης, το οποίο συμβαίνει όταν όλα τα local block έχουνε συγχλίνει

```

1 if (t%C==0) {
2     /*Test convergence*/
3     gettimeofday(&tcvs, NULL);
4     converged=converge(u_previous,u_current,local[0],local[1]);
5     MPI_Allreduce(&converged, &global_converged, 1, MPI_INT, MPI_MIN, MPI_COMM_WORLD);
6     gettimeofday(&tcvf, NULL);
7     tconv += (tcvf.tv_sec-tcvf.tv_sec)+(tcvf.tv_usec-tcvf.tv_usec)*0.000001;
8 }

```

8. Συνένωση των επιμέρους local πινάκων στον αρχικό global πίνακα U, τον οποίο τον χειρίζεται η διεργασία με rank = 0

```

1 //-----Rank 0 gathers local matrices back to the global matrix-----//
2 if (rank==0) {
3     U=allocate2d(global_padded[0],global_padded[1]);
4     U_start = &(U[0][0]);
5 }
6 MPI_Gatherv(&u_current[1][1], 1, local_block, U_start, scattercounts, scatteroffset,
    global_block, 0, MPI_COMM_WORLD);

```

### 3.1 Μέθοδος Jacobi

Στην μέθοδο αυτή, κάθε σημείο για να ανανεωθεί χρειάζεται τις τιμές των γειτόνων του. Συνεπώς, κάθε φορά πριν το υπολογιστικό μέρος, κάθε block χρειάζεται να ανταλλάξει με τους γείτονές του τα σύνορά τους. Οι τιμές που έρχονται από τους γείτονες αποθηκεύονται στα ghost cells που έχουμε δεσμεύσει στην αρχή, με αποτέλεσμα να μην χρειάζεται να αλλάξει καθόλου ο βρόγχος ανανέωσης του σειριακού προγράμματος. Μόλις ολοκληρωθεί η επικοινωνία κάθε κελί θα έχει τις σωστές τιμές των γειτόνων του.

```

1 //-----Computational core-----//
2 gettimeofday(&tts, NULL);
3 #ifdef TEST_CONV
4 for (t=0;t<T && !global_converged;t++) {
5     #endif
6     #ifndef TEST_CONV
7     #undef T
8     #define T 256
9     for (t=0;t<T;t++) {
10        #endif
11        swap=u_previous;
12        u_previous=u_current;
13        u_current=swap;
14        /*Compute and Communicate*/
15        // Communicate with north
16        if (north != MPI_PROC_NULL){
17            MPI_Sendrecv(&u_previous[1][1], local[1], MPI_DOUBLE, north, 0, &u_previous[0][1],
18                local[1], MPI_DOUBLE, north, 0, MPI_COMM_WORLD, &status );
19        }
20        // Communicate with south
21        if (south != MPI_PROC_NULL){

```

```

22     MPI_Sendrecv(&u_previous[local[0]][1], local[1], MPI_DOUBLE, south, 0, &u_previous[
23     local[0]+1][1],
24     local[1], MPI_DOUBLE, south, 0, MPI_COMM_WORLD, &status );
25 }
26 // Communicate with east
27 if (east != MPI_PROC_NULL){
28     MPI_Sendrecv(&u_previous[1][local[1]], 1, column, east, 0, &u_previous[1][local
29     [1]+1],
30     1, column, east, 0, MPI_COMM_WORLD, &status );
31 }
32 // Communicate with west
33 if (west != MPI_PROC_NULL){
34     MPI_Sendrecv(&u_previous[1][1], 1, column, west, 0, &u_previous[1][0],
35     1, column, west, 0, MPI_COMM_WORLD, &status );
36 }
37 /*Add appropriate timers for computation*/
38 gettimeofday(&tcs, NULL);
39
40 for (i=i_min; i<i_max; i++)
41     for (j=j_min; j<j_max; j++)
42         u_current[i][j]=(u_previous[i-1][j]+u_previous[i+1][j]+u_previous[i][j-1]+
43         u_previous[i][j+1])/4.0;
44
45 gettimeofday(&tcf, NULL);
46 tcomp += (tcf.tv_sec-tcs.tv_sec)+(tcf.tv_usec-tcs.tv_usec)*0.000001;
47 #ifndef TEST_CONV
48 if (t%C==0) {
49     /*Test convergence*/
50     gettimeofday(&tcvs, NULL);
51     converged=converge(u_previous,u_current,local[0],local[1]);
52     MPI_Allreduce(&converged, &global_converged, 1, MPI_INT, MPI_MIN, MPI_COMM_WORLD);
53     gettimeofday(&tcvf, NULL);
54     tconv += (tcvf.tv_sec-tcvs.tv_sec)+(tcvf.tv_usec-tcvs.tv_usec)*0.000001;
55 }
56 #endif
57 }
58 gettimeofday(&tts, NULL);
59 tttotal=(tts.tv_sec-ttf.tv_sec)+(tts.tv_usec-ttf.tv_usec)*0.000001;
60
61 MPI_Reduce(&tttotal,&total_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);
62 MPI_Reduce(&tcomp,&comp_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);
63 MPI_Reduce(&tconv,&conv_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);

```

### 3.2 Μέθοδος Gauss-Seidel SOR

Στην μέθοδο αυτή, η υλοποίηση της επικοινωνίας δεν είναι τόσο απλή. Κάθε block πρέπει:

- Πριν τον υπολογισμό, να στείλει στον πάνω και στον αριστερά τα δεδομένα της προηγούμενης χρονικής στιγμής.
- Πριν τον υπολογισμό, να λάβει από τον πάνω και τον αριστερά τα δεδομένα της τωρινής χρονικής στιγμής.
- Πριν τον υπολογισμό, να λάβει από τον δεξιά και τον κάτω τα δεδομένα της προηγούμενης χρονικής στιγμής.

- Μετά τον υπολογισμό, να στείλει στον δεξιά και στον κάτω τα δεδομένα της τωρινής χρονικής στιγμής.

Ουσιαστικά, αυτό που θα συμβεί είναι ότι αρχικά όλες οι διεργασίες, εκτός από αυτές που βρίσκονται πάνω και αριστερά στο ταμπλό, θα μπλοκάρουν περιμένοντας τα τωρινά δεδομένα από τον πάνω και τον αριστερά τους. Οι διεργασίες που δεν έχουν γείτονα πάνω και αριστερά θα υπολογίσουν τις νέες τιμές και θα τις στείλουν στον κάτω και δεξιά γείτονά τους. Η διαδικασία αυτή θα επαναληφθεί μέχρι όλα τα block να ανανεωθούν κ.ο.κ.. Για να υλοποιήσουμε αυτή την διαδικασία, χρησιμοποιήσαμε non-blocking send και receive συναρτήσεις και θέσαμε στα σημεία όπου χρειαζόμαστε τις γειτονικές τιμές το κατάλληλο wait, όπως φαίνεται και στον κώδικα παρακάτω.

```

1 MPI_Request before_req[6];
2 MPI_Status before_status[6];
3 MPI_Request after_req[2];
4 MPI_Status after_status[2];
5 int before_req_len = 0;
6 int after_req_len = 0;
7
8 //-----Computational core-----//
9 gettimeofday(&tts, NULL);
10 #ifndef TEST_CONV
11 for (t=0;t<T && !global_converged;t++) {
12     #endif
13     #ifndef TEST_CONV
14     #undef T
15     #define T 256
16     for (t=0;t<T;t++) {
17         #endif
18
19         before_req_len = 0;
20         after_req_len = 0;
21
22         swap=u_previous;
23         u_previous=u_current;
24         u_current=swap;
25
26         /* Send data to north */
27         if (north != MPI_PROC_NULL){
28             MPI_Isend(&u_previous[1][1], local[1], MPI_DOUBLE, north, 0, MPI_COMM_WORLD, &
29             before_req[before_req_len]);
30             MPI_Irecv(&u_current[0][1], local[1], MPI_DOUBLE, north, 0, MPI_COMM_WORLD, &
31             before_req[before_req_len + 1]);
32             before_req_len += 2;
33         }
34
35         /* Send data to west */
36         if (west != MPI_PROC_NULL){
37             MPI_Isend(&u_previous[1][1], 1, column, west, 0, MPI_COMM_WORLD, &before_req[
38             before_req_len]);
39             MPI_Irecv(&u_current[1][0], 1, column, west, 0, MPI_COMM_WORLD, &before_req[
40             before_req_len + 1]);
41             before_req_len += 2;
42         }
43
44         /* Get data from south */
45         if (south != MPI_PROC_NULL){

```

```

42     MPI_Irecv(&u_previous[local[0]+1][1], local[1], MPI_DOUBLE, south, 0, MPI_COMM_WORLD
, &before_req[before_req_len]);
43     before_req_len++;
44 }
45 /* Get data from east */
46 if (east != MPI_PROC_NULL){
47     MPI_Irecv(&u_previous[1][local[1]+1], 1, column, east, 0, MPI_COMM_WORLD, &
before_req[before_req_len]);
48     before_req_len++;
49 }
50 MPI_Waitall(before_req_len, before_req, before_status);
51
52 /*Compute and Communicate*/
53
54 /*Add appropriate timers for computation*/
55 gettimeofday(&tcs, NULL);
56 for (i=i_min; i<i_max; i++)
57     for (j=j_min; j<j_max; j++)
58         u_current[i][j]=u_previous[i][j] + (omega/4.0)*(u_current[i-1][j]+u_previous[i
+1][j]+u_current[i][j-1]+u_previous[i][j+1]-4*u_previous[i][j]);
59
60 gettimeofday(&tcf, NULL);
61 tcomp += (tcf.tv_sec-tcs.tv_sec)+(tcf.tv_usec-tcs.tv_usec)*0.000001;
62
63 /* Send data to south */
64 if (south != MPI_PROC_NULL){
65     MPI_Isend(&u_current[local[0]][1], local[1], MPI_DOUBLE, south, 0, MPI_COMM_WORLD, &
after_req[after_req_len]);
66     after_req_len++;
67 }
68
69 /* Send data to east */
70 if (east != MPI_PROC_NULL){
71     MPI_Isend(&u_current[1][local[1]], 1, column, east, 0, MPI_COMM_WORLD, &after_req[
after_req_len]);
72     after_req_len++;
73 }
74 MPI_Waitall(after_req_len, after_req, after_status);
75
76 #ifndef TEST_CONV
77 if (t%C==0) {
78     /*Test convergence*/
79     gettimeofday(&tcvs, NULL);
80     converged=converge(u_previous, u_current, local[0], local[1]);
81     MPI_Allreduce(&converged, &global_converged, 1, MPI_INT, MPI_MIN, MPI_COMM_WORLD);
82     gettimeofday(&tcvf, NULL);
83     tconv += (tcvf.tv_sec-tcvf.tv_sec)+(tcvf.tv_usec-tcvf.tv_usec)*0.000001;
84 }
85 #endif
86 }
87 gettimeofday(&ttf, NULL);
88 tttotal=(ttf.tv_sec-tts.tv_sec)+(ttf.tv_usec-tts.tv_usec)*0.000001;
89
90 MPI_Reduce(&tttotal, &total_time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);
91 MPI_Reduce(&tcomp, &comp_time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);
92 MPI_Reduce(&tconv, &conv_time, 1, MPI_DOUBLE, MPI_MAX, 0, MPI_COMM_WORLD);

```

### 3.3 Μέθοδος Red-Black SOR

Τέλος, στην μέθοδο Red-Black SOR έχουμε δύο φάσεις υπολογισμού. Στην πρώτη φάση, τα κόκκινα στοιχεία πρέπει να γνωρίζουν τους γείτονές τους και στην δεύτερη φάση τα μαύρα στοιχεία αντίστοιχα. Από την στιγμή, που τα κόκκινα και τα μαύρα στοιχεία τοποθετούνται εναλλάξ στις άρτιες και περιττές θέσεις, και στις δύο φάσεις όλα τα block χρειάζεται να επικοινωνήσουν με τους γείτονές τους. Έτσι, καταλήγουμε σε ένα αντίστοιχο μοντέλο με αυτό του Jacobi, με την διαφορά ότι τώρα θα γίνει 2 φορές η επικοινωνία, μία για να ενημερωθεί ο πίνακας `u_previous` και να υπολογιστούν τα κόκκινα και μία για να ενημερωθεί ο πίνακας `u_current` και να υπολογιστούν τα μαύρα στοιχεία.

```
1  //-----Computational core-----//
2  gettimeofday(&tts, NULL);
3  #ifdef TEST_CONV
4  for (t=0;t<T && !global_converged;t++) {
5  #endif
6  #ifndef TEST_CONV
7  #undef T
8  #define T 256
9  for (t=0;t<T;t++) {
10 #endif
11
12     swap=u_previous;
13     u_previous=u_current;
14     u_current=swap;
15
16     // Communicate with north
17     if (north != MPI_PROC_NULL){
18         MPI_Sendrecv(&u_previous[1][1], local[1], MPI_DOUBLE, north, 0, &u_previous[0][1],
19                     local[1], MPI_DOUBLE, north, 0, MPI_COMM_WORLD, &status );
20     }
21
22     // Communicate with south
23     if (south != MPI_PROC_NULL){
24         MPI_Sendrecv(&u_previous[local[0]][1], local[1], MPI_DOUBLE, south, 0, &u_previous[
25         local[0]+1][1],
26                     local[1], MPI_DOUBLE, south, 0, MPI_COMM_WORLD, &status );
27     }
28
29     // Communicate with east
30     if (east != MPI_PROC_NULL){
31         MPI_Sendrecv(&u_previous[1][local[1]], 1, column, east, 0, &u_previous[1][local
32         [1]+1],
33                     1, column, east, 0, MPI_COMM_WORLD, &status );
34     }
35
36     // Communicate with west
37     if (west != MPI_PROC_NULL){
38         MPI_Sendrecv(&u_previous[1][1], 1, column, west, 0, &u_previous[1][0],
39                     1, column, west, 0, MPI_COMM_WORLD, &status );
40     }
41
42     /*Add appropriate timers for computation*/
43     gettimeofday(&tcs, NULL);
44     for (i=i_min;i<i_max;i++)
45         for (j=j_min;j<j_max;j++)
```

```

45         if ((i+j)%2==0)
46             u_current[i][j]=u_previous[i][j]+(omega/4.0)*(u_previous[i-1][j]+u_previous[
47 i+1][j]+u_previous[i][j-1]+u_previous[i][j+1]-4*u_previous[i][j]);
48 gettimeofday(&tcf, NULL);
49 tcomp += ( tcf.tv_sec-tcs.tv_sec)+(tcf.tv_usec-tcs.tv_usec)*0.000001;
50
51 // Communicate with north
52 if (north != MPI_PROC_NULL){
53     MPI_Sendrecv(&u_current[1][1], local[1], MPI_DOUBLE, north, 0, &u_current[0][1],
54 local[1], MPI_DOUBLE, north, 0, MPI_COMM_WORLD, &status );
55 }
56
57 // Communicate with south
58 if (south != MPI_PROC_NULL){
59     MPI_Sendrecv(&u_current[local[0]][1], local[1], MPI_DOUBLE, south, 0, &u_current[
60 local[0]+1][1],
61 local[1], MPI_DOUBLE, south, 0, MPI_COMM_WORLD, &status );
62 }
63
64 // Communicate with east
65 if (east != MPI_PROC_NULL){
66     MPI_Sendrecv(&u_current[1][local[1]], 1, column, east, 0, &u_current[1][local[1]+1],
67 1, column, east, 0, MPI_COMM_WORLD, &status );
68 }
69
70 // Communicate with west
71 if (west != MPI_PROC_NULL){
72     MPI_Sendrecv(&u_current[1][1], 1, column, west, 0, &u_current[1][0],
73 1, column, west, 0, MPI_COMM_WORLD, &status );
74 }
75
76 gettimeofday(&tcs, NULL);
77
78 for (i=i_min; i<i_max; i++)
79     for (j=j_min; j<j_max; j++)
80         if ((i+j)%2==1)
81             u_current[i][j]=u_previous[i][j]+(omega/4.0)*(u_current[i-1][j]+u_current[i
82 +1][j]+u_current[i][j-1]+u_current[i][j+1]-4*u_previous[i][j]);
83
84 gettimeofday(&tcf, NULL);
85 tcomp += ( tcf.tv_sec-tcs.tv_sec)+(tcf.tv_usec-tcs.tv_usec)*0.000001;
86
87 #ifdef TEST_CONV
88     if (t%C==0) {
89         /*Test convergence*/
90         gettimeofday(&tcvs, NULL);
91         converged=converge(u_previous,u_current,local[0],local[1]);
92         MPI_Allreduce(&converged, &global_converged, 1, MPI_INT, MPI_MIN, MPI_COMM_WORLD);
93         gettimeofday(&tcvf, NULL);
94         tconv += ( tcvf.tv_sec-tcvs.tv_sec)+(tcvf.tv_usec-tcvs.tv_usec)*0.000001;
95     }
96 #endif
97 }
98
99 gettimeofday(&ttf, NULL);
100 ttotal=(ttf.tv_sec-tts.tv_sec)+(ttf.tv_usec-tts.tv_usec)*0.000001;
101 MPI_Reduce(&ttotal,&total_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);
102 MPI_Reduce(&tcomp,&comp_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);
103 MPI_Reduce(&tconv,&conv_time,1,MPI_DOUBLE,MPI_MAX,0,MPI_COMM_WORLD);

```

## 4 Μετρήσεις επίδοσης

### 4.1 Μετρήσεις με έλεγχο σύγκλισης

Εδώ, πραγματοποιήσαμε μετρήσεις με έλεγχο σύγκλισης για τα παράλληλα προγράμματα που αναπτύχθηκαν σε MPI, για μέγεθος πίνακα 1024x1024. Γι' αυτό το μέγεθος πίνακα, λάβαμε μετρήσεις (συνολικός χρόνος, χρόνος υπολογισμών, χρόνος ελέγχου σύγκλισης) για 64 MPI διεργασίες. Οι μετρήσεις αυτές συνοψίζονται στον παρακάτω πίνακα:

Algorithm	Grid Size	Iterations	Compute	Converge	Communicate	Total
Jacobi	16x4	798201	36.06	23.57	536.5	599.13
GaussSeidelSOR	16x4	3201	0.47	0.23	2.58	3.28
RedBlackSOR	16x4	2501	0.29	0.14	3.39	3.82
Jacobi	4x16	798201	38.9	28.23	418.86	485.99
GaussSeidelSOR	4x16	3201	0.44	0.22	2.22	2.88
RedBlackSOR	4x16	2501	0.29	0.11	2.7	3.1
Jacobi	8x8	798201	40.28	9.49	217.01	266.78
GaussSeidelSOR	8x8	3201	0.48	0.13	1.33	1.94
RedBlackSOR	8x8	2501	0.31	0.08	1.38	1.77

Συνοψίζουμε τις παραπάνω μετρήσεις, στα διαγράμματα που ακολουθούν:

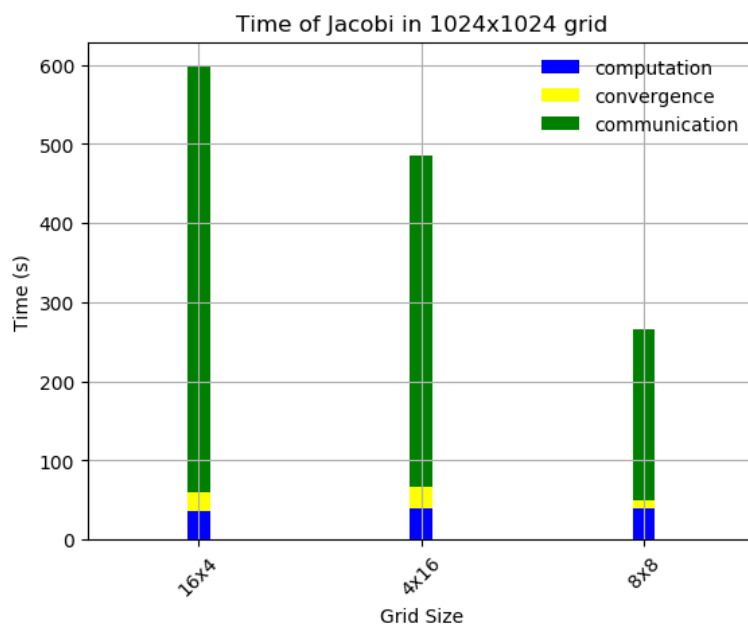
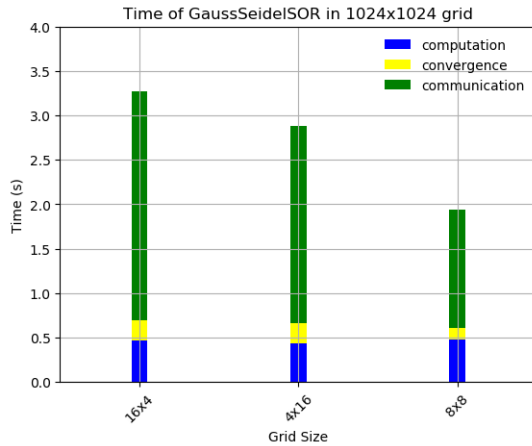
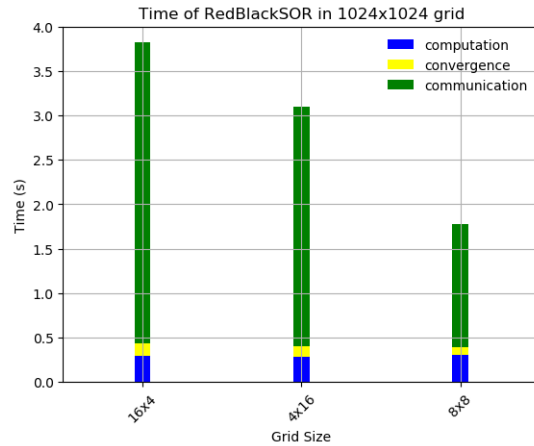


Figure 1: Barplot για Jacobi



(a) Barplot για GaussSeidelSOR



(b) Barplot για RedBlackSOR

Αρχικά, παρατηρούμε ότι σε όλες τις περιπτώσεις μας συμφέρει να ορίσουμε grid size 8x8, καθώς έτσι έχουμε μεταφορά ίσου μεγέθους δεδομένων μεταξύ των διεργασιών. Στις άλλες περιπτώσεις, η μεταφορά των δεδομένων στην μεγαλύτερη πλευρά ενός block καθυστερούσε όλη την διαδικασία. Ως προς το είδος του αλγορίθμου τα συμπεράσματα από τα πειράματα ήταν και αυτά που αναμέναμε θεωρητικά. Δηλαδή:

- ✓ Ο Jacobi είναι πολύ πιο αργός από τους άλλους δύο αλγορίθμους. Αυτό συμβαίνει διότι έχει πάρα πολύ αργό ρυθμό σύγκλισης. Αν δούμε τον παραπάνω πίνακα, συγκλίνει μετά από 798201 επαναλήψεις αριθμός 200-πλάσιος από τις επαναλήψεις των άλλων δύο. Συνεπώς, δεν έχει καμία σημασία ο χρόνος υπολογισμού, επικοινωνίας και σύγκλισης του Jacobi σε μία επανάληψη από την στιγμή που οι επαναλήψεις του είναι τόσο πολλές.
- ✓ Ανάμεσα στους άλλους δύο, είναι λίγο πιο γρήγορος ο RedBlackSOR. Ουσιαστικά, η διαφορά των δύο αλγορίθμων είναι ότι από την μία ο Gauss έχει μικρότερο χρόνο επικοινωνίας και υπολογισμού ανά επανάληψη αλλά από την άλλη ο RedBlack συγκλίνει πιο γρήγορα (700 επαναλήψεις λιγότερο). Συνεπώς, οι λιγότερες επαναλήψεις του δεύτερου υπερνικούν του καλύτερου χρόνου ανά επανάληψη που έχει ο πρώτος. Για τον λόγο αυτό, βλέπουμε τις πράσινες περιοχές σχεδόν ίσες αλλά τις κίτρινες και μπλε μικρότερες στον RedBlackSOR.
- ✓ Παρατηρούμε την μεγάλη σημασία που έχει στα συστήματα κατανεμημένης μνήμης το δίκτυο διασύνδεσης καθώς σε όλες τις περιπτώσεις ο χρόνος επικοινωνίας είναι ο μεγαλύτερος (μη ξεχνάμε ότι έχουμε 64 MPI διεργασίες που σημαίνει επικοινωνία μεταξύ 8 nodes).

Συνολικά, για την επίλυση του προβλήματος στο συγκεκριμένο σύστημα κατανεμημένης μνήμης θα επιλέγαμε την μέθοδο RedBlackSOR. Ωστόσο, στην γενική περίπτωση, η απάντηση δεν είναι τόσο απόλυτη. Σε περίπτωση για παράδειγμα που υπάρχει υψηλή κίνηση στο δίκτυο διασύνδεσης ίσως προτιμούσαμε την μέθοδο GaussSeidel. Αυτό γιατί ο RedBlackSOR έχει υψηλότερο ρυθμό επικοινωνίας και η απόδοσή του θα επηρεαζόταν περισσότερο από την απόδοση του GaussSeidel.

## 4.2 Μετρήσεις χωρίς έλεγχο σύγκλισης

Εδώ, πραγματοποιήσαμε μετρήσεις, απενεργοποιώντας τον έλεγχο σύγκλισης, για σταθερό αριθμό επαναλήψεων ( $T=256$ ), για μεγέθη πίνακα 2048x2048, 4096x4096, 6144x6144, για 1,2,4,8,16,32 και 64



MPI διεργασίες. Αρχικά, παρουσιάζεται ένα διάγραμμα επιτάχυνσης για κάθε μέγεθος πίνακα (άξονας x: MPI διεργασίες, άξονας y: speedup) που απεικονίζει τις τρεις μεθόδους (Jacobi, Gauss-Seidel SOR και Red-Black SOR).



Figure 3: Speedup σε 2048x2048

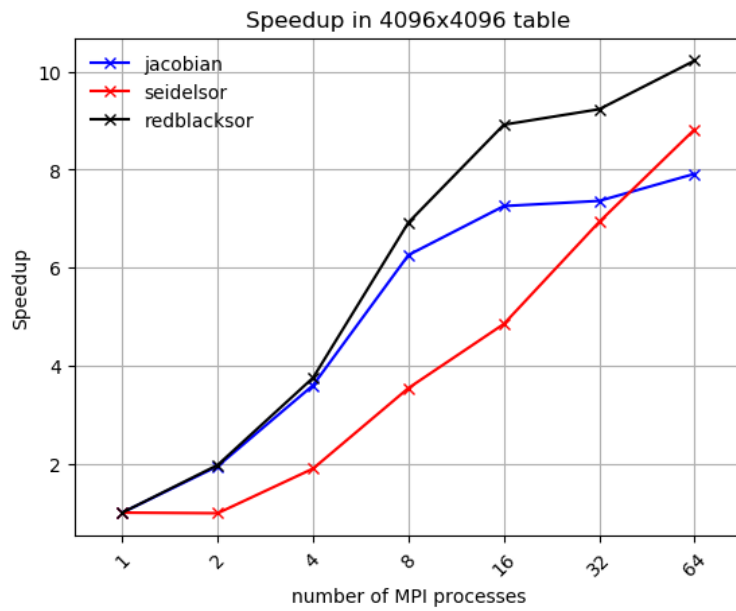


Figure 4: Speedup σε 4096x4096

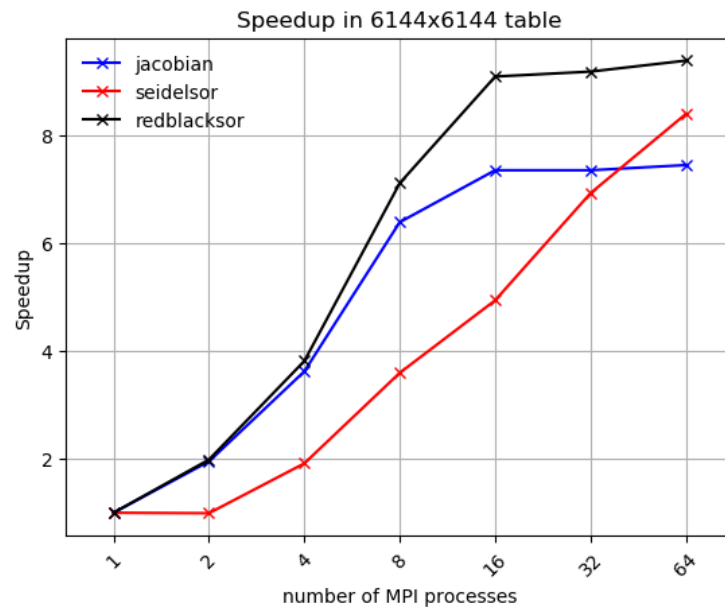
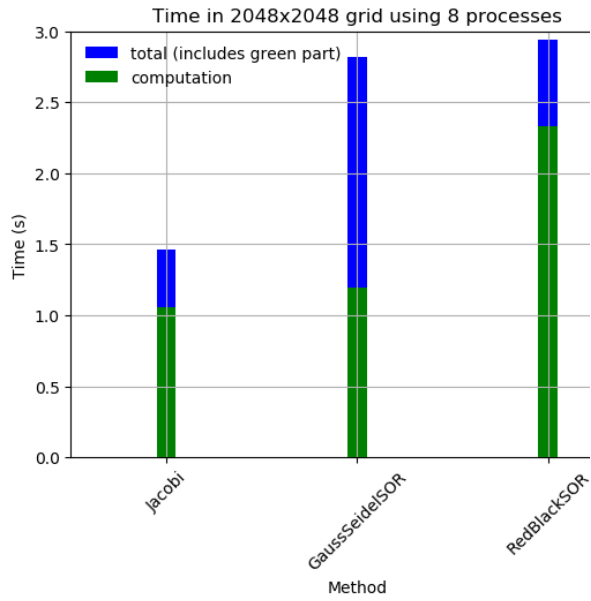


Figure 5: Speedup σε 6144x6144

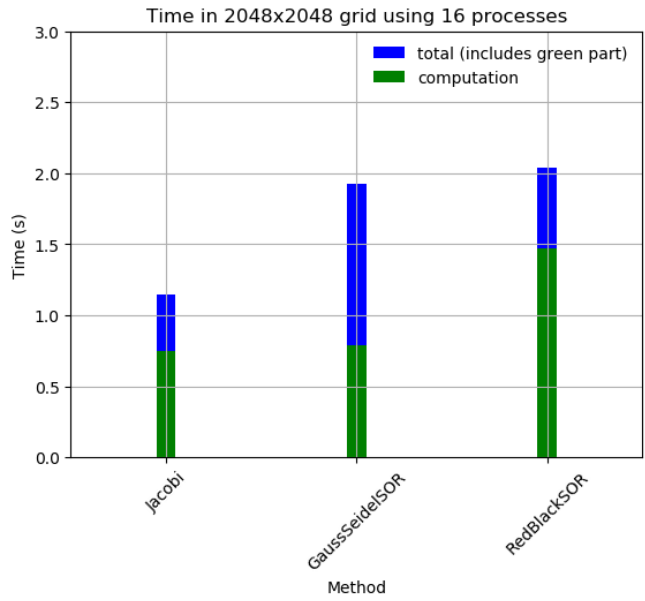
Επίσης, για καλύτερη ανάλυση των μετρήσεων, κατασκευάζουμε διαγράμματα με μπάρες (1 για κάθε

μέγεθος πίνακα και αριθμό επεξεργαστών) που θα απεικονίζουν το συνολικό χρόνο εκτέλεσης και το χρόνο υπολογισμού για κάθε μία από τις τρεις μεθόδους (άξονας x: μέθοδος, άξονας y: χρόνος). Για λόγους καλύτερης εποπτείας, κατασκευάσαμε διαγράμματα μόνο για 8, 16, 32 και 64 MPI διεργασίες, κρατώντας κοινή κλίμακα στον άξονα y ανά μέγεθος πίνακα.

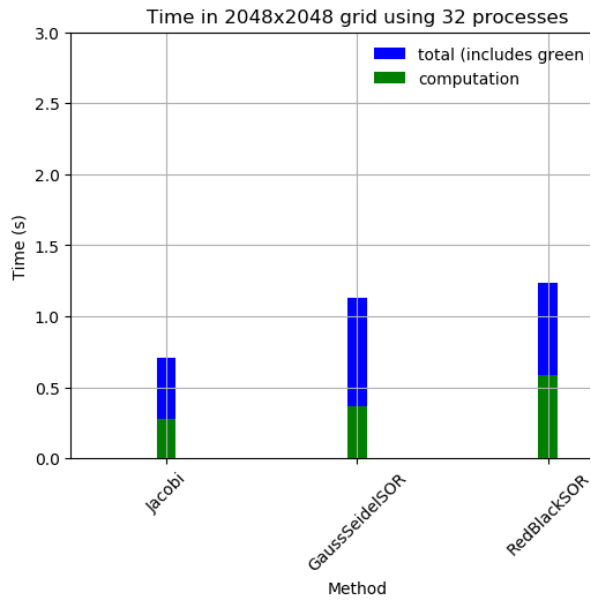
• 2048x2048



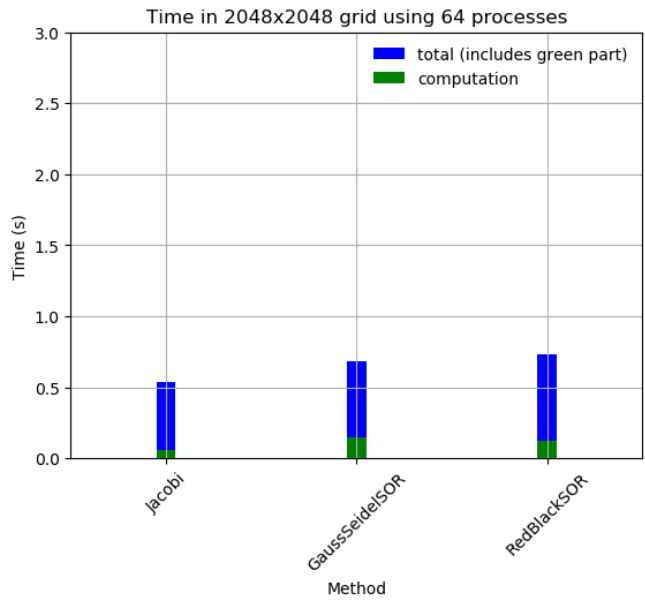
(a) Barplot για 8 processes



(b) Barplot για 16 processes

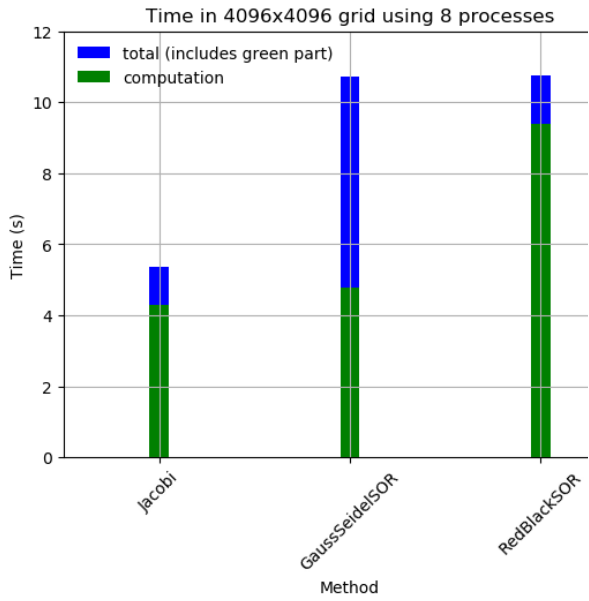


(a) Barplot για 32 processes

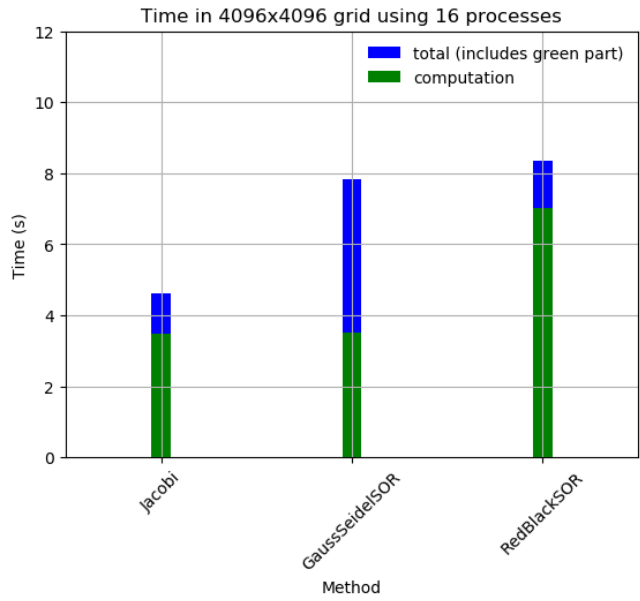


(b) Barplot για 64 processes

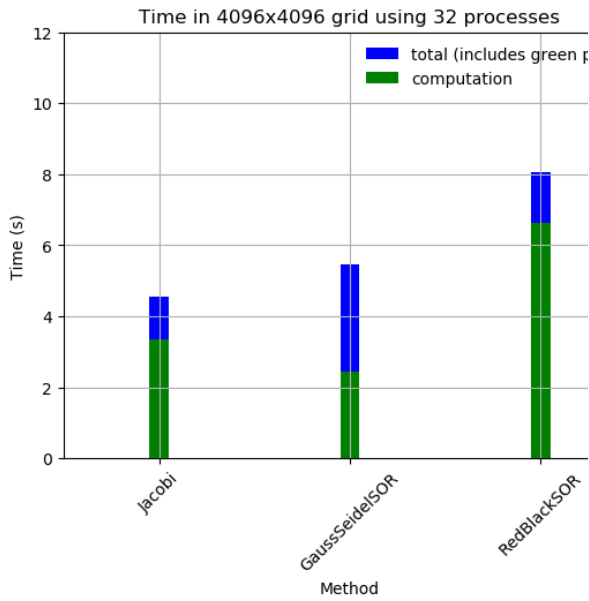
• 4096x4096



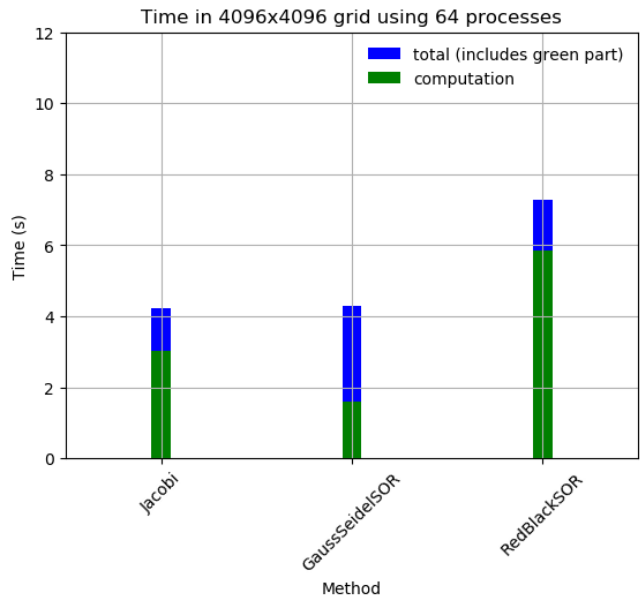
(a) Barplot για 8 processes



(b) Barplot για 16 processes

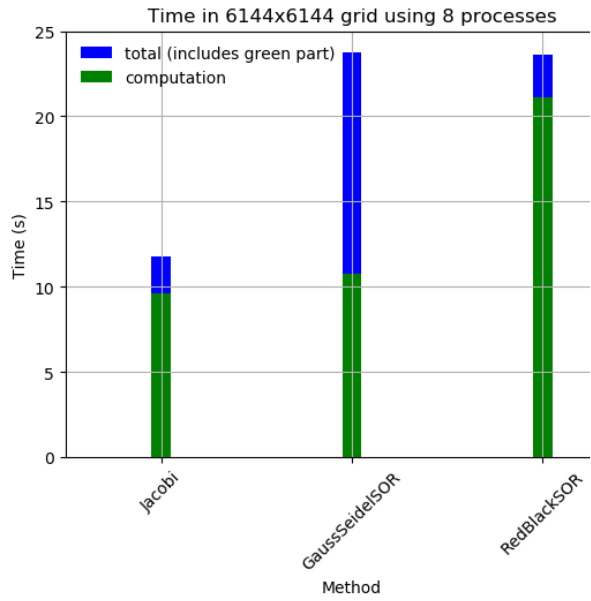


(a) Barplot για 32 processes

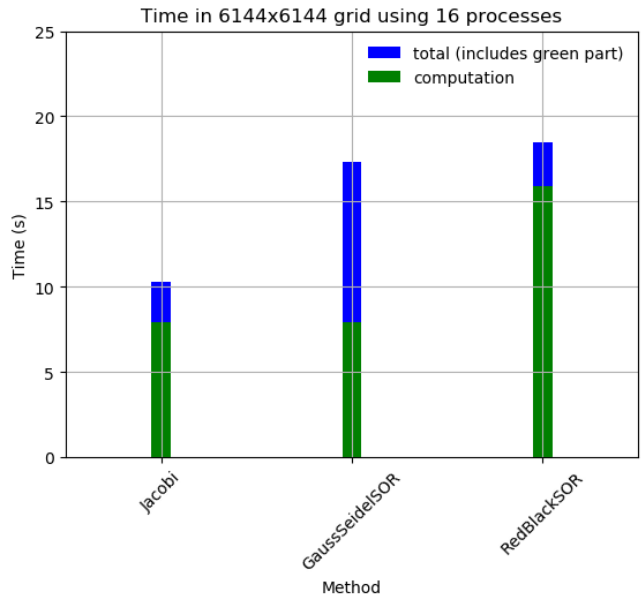


(b) Barplot για 64 processes

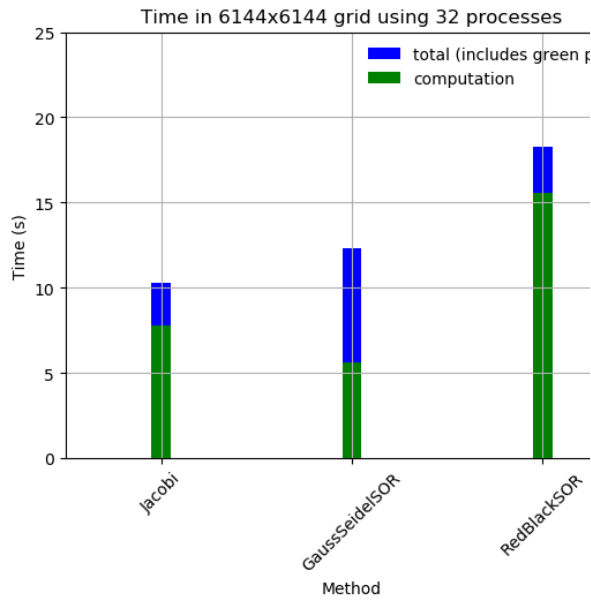
• 6144x6144



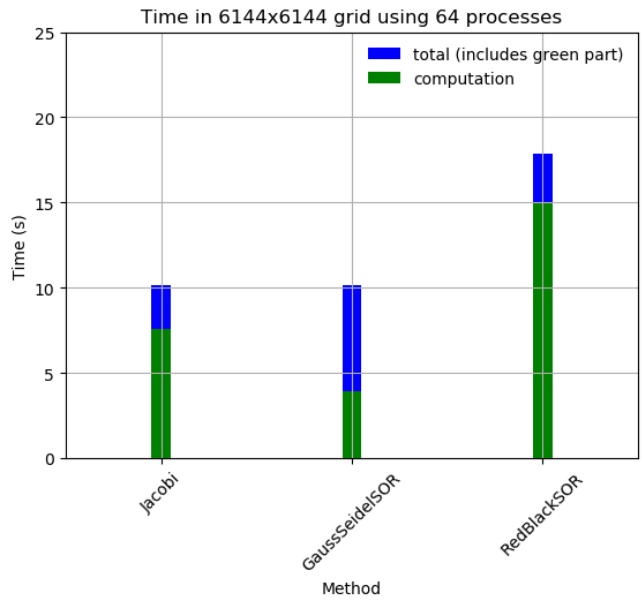
(a) Barplot για 8 processes



(b) Barplot για 16 processes



(a) Barplot για 32 processes



(b) Barplot για 64 processes

Έχουμε τις εξής παρατηρήσεις ως προς την κλιμακωσιμότητα των αλγορίθμων μας:

- ✓ Για σταθερό αριθμό επαναλήψεων, ο Jacobi είναι πιο γρήγορος από τους δύο άλλους αλγορίθμους καθώς έχει λιγότερο κόστος επικοινωνίας από τον Gauss και λιγότερο υπολογιστικό κόστος από τον RedBlack.
- ✓ Αν εξαιρέσουμε το μικρότερο μέγεθος πίνακα, όπου κλιμακώνουν όλοι οι αλγόριθμοι, αυτός που κλιμακώνει καλύτερα για τα μεγαλύτερα μεγέθη είναι ο GaussSeidel. Ο RedBlackSOR μπορεί να έχει συνεχώς μεγαλύτερο speedup από τους άλλους δύο αλλά δεν πρέπει να παραβλέψουμε ότι ο χρόνος για ένα πυρήνα του RedBlackSOR είναι σχεδόν διπλάσιος από τους άλλους δύο, το οποίο οφείλεται στο διπλάσιο υπολογιστικό κόστος. Έτσι, μπορεί ο RedBlack να έχει μεγάλο speedup αλλά δεν κλιμακώνει από ένα σημείο και μετά λόγω του υπολογιστικού κόστους που πρέπει να καλύπτει κάθε φορά.
- ✓ Η καλύτερη κλιμάκωση του Gauss φαίνεται και στα τελευταία barplot όπου καταφέρνει να πετύχει ίσους χρόνους με τον Jacobi, παρά το γεγονός ότι προφανώς ο χρόνος επικοινωνίας είναι μεγαλύτερος.
- ✓ Συνολικά, ο GaussSeidelSOR, παρά το κόστος επικοινωνίας που έχει, κλιμακώνει καλύτερα μειώνοντας συνεχώς το συνολικό του χρόνο, το οποίο οφείλεται σίγουρα και στην non-blocking διαχείριση της επικοινωνίας. Από την άλλη, ο RedBlackSOR και ο Jacobi δεν κλιμακώνουν τόσο καλά διότι τα blocking SendRecv εμποδίζουν την επιπλέον μείωση του χρόνου επικοινωνίας.

## 5 Σημειώσεις

- Όλες οι υλοποιήσεις μεταγλωτίστηκαν με όρισμα -O3 για compiler optimizations.
- Η ορθότητα κάθε υλοποίησης εξετάστηκε συγκρίνοντας το output που έδινε με το output του αντίστοιχου σειριακού αλγορίθμου.
- Μια πιθανή βελτιστοποίηση στο Jacobi που δεν προλάβαμε να υλοποιήσουμε είναι, καθώς γίνεται non-blocking επικοινωνία με τους γείτονες, να υπολογίζονται οι νέες τιμές στο εσωτερικό του πίνακα όπου δεν απαιτείται η γνώση τιμών που ανήκουν σε άλλες διεργασίες.
- Όπου αναφέρεται χρόνος επικοινωνίας, έχει προκύψει από την διαφορά του συνολικού χρόνου από τον χρόνο σύγκλισης και τον χρόνο υπολογισμού. Συνεπώς, περιλαμβάνει τον χρόνο για το swap το οποίο θεωρείται αμελητέο.
- Όλες οι μετρήσεις έγιναν από 3 φορές και κρατήσαμε κάθε φορά τον μέσο όρο τους για μεγαλύτερη ακρίβεια.



## Αναφορές

- [1] "Σημειώσεις του μαθήματος" <http://www.cslab.ntua.gr/courses/pps/notes.go>
- [2] "MPI documentation" <https://www.open-mpi.org/doc/>
- [3] "Heat equation" [https://en.wikipedia.org/wiki/Heat\\_equation](https://en.wikipedia.org/wiki/Heat_equation)