

# Minimum Spanning Tree: Prim's Algorithm

1. Initialize a tree with a single vertex, chosen arbitrarily from the graph.
2. Grow the tree by one edge: of the edges that connect the tree to vertices not yet in the tree, find the minimum-weight edge, and transfer it to the tree.
3. Repeat step 2 (until all vertices are in the tree).

# Complexity Analysis

<b>Minimum edge weight data structure</b>	<b>Time complexity (total)</b>
<b>adjacency matrix, searching</b>	<b><math>O(V*V)</math></b>
<b>binary heap and adjacency list</b>	<b><math>O((V + E) \log(V)) = O(E \log(V))</math></b>
<b>Fibonacci heap and adjacency list</b>	<b><math>O(E + V \log(V))</math></b>

# Parallel Formulation

## Prim's Algorithm for Finding the Minimum Spanning Tree

1. Initialize  $S$  with the start vertex,  $s$ , and  $V-S$  with the remaining vertices.
2. **for** all  $v$  in  $V-S$
3.     Set  $p[v]$  to  $s$ .
4.     **if** there is an edge  $(s, v)$
5.         Set  $d[v]$  to  $w(s, v)$ .
6.     **else** Set  $d[v]$  to  $\infty$ .
7.     **while**  $V-S$  is not empty
8.         **for** all  $u$  in  $V-S$ , find the smallest  $d[u]$ .
9.         Remove  $u$  from  $V-S$  and add it to  $S$ .
10.         Insert the edge  $(u, p[u])$  into the spanning tree.
11.         **for** all  $v$  in  $V-S$
12.             **if**  $w(u, v) < d[v]$
13.                 Set  $d[v]$  to  $w(u, v)$ .
14.                 Set  $p[v]$  to  $u$ .

Δύσκολα παραλληλοποιείται  
(δεν μπορούμε να βάλουμε δύο  
κορυφές ταυτόχρονα στο MST)

# Parallel Formulation

## Prim's Algorithm for Finding the Minimum Spanning Tree

1. Initialize  $S$  with the start vertex,  $s$ , and  $V-S$  with the remaining vertices.
2. **for** all  $v$  in  $V-S$
3.     Set  $p[v]$  to  $s$ .
4.     **if** there is an edge  $(s, v)$
5.         Set  $d[v]$  to  $w(s, v)$ .
6.     **else**
7.         Set  $d[v]$  to  $\infty$ .
8. **while**  $V-S$  is not empty
9.     **for** all  $u$  in  $V-S$ , find the smallest  $d[u]$ .
10.     Remove  $u$  from  $V-S$  and add it to  $S$ .
11.     Insert the edge  $(u, p[u])$  into the spanning tree.
12.     **for** all  $v$  in  $V-S$
13.         **if**  $w(u, v) < d[v]$
14.             Set  $d[v]$  to  $w(u, v)$ .
15.             Set  $p[v]$  to  $u$ .

Εύκολη  
παραλληλοποίηση!

P = ο αριθμός των διεργασιών

N = ο αριθμός των κόμβων του γράφου

$V_i$  = το σύνολο των κόμβων που ανατίθενται στην διεργασία  $P_i$ ,  $P_i = 0, \dots, P-1$

Ιδέα: κάθε διεργασία θα υπολογίζει παράλληλα τα

$$di[u] = \min\{di[v] \mid v \in (V - V_T) \cap V_i\}$$

Το ολικό ελάχιστο βρίσκεται με reduction πάνω σε μερικά ελάχιστα

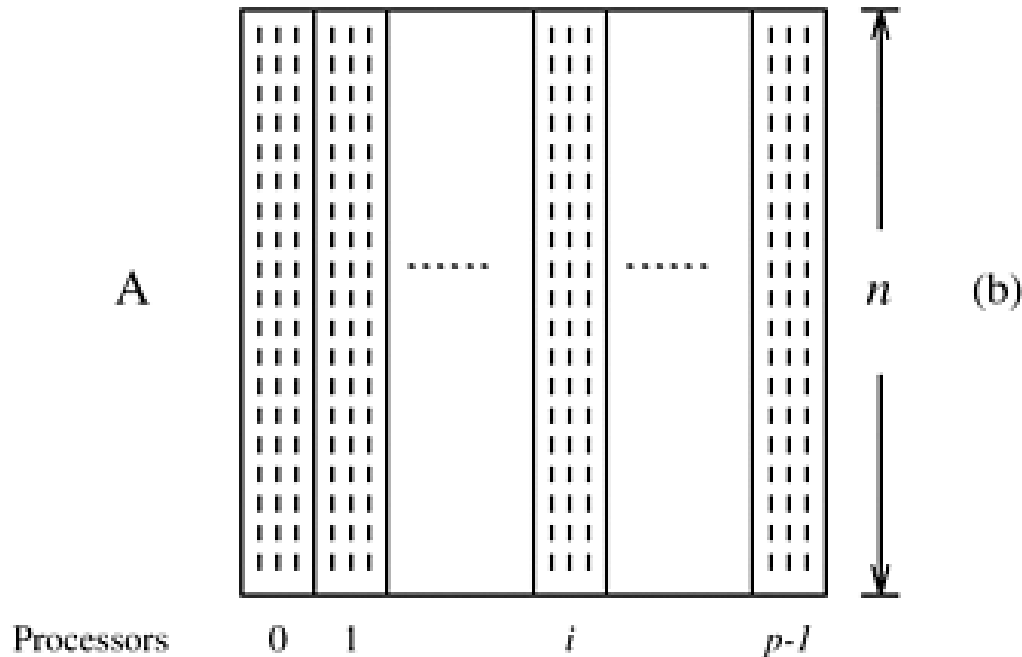
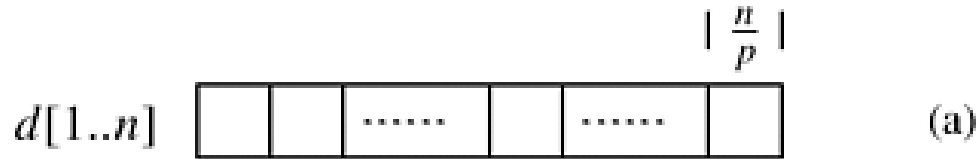
Σε κάθε βήμα η βασική διεργασία P0 θα έχει το αποτέλεσμα του reduce operation, δηλαδή ποιος κόμβος θα προστεθεί στο δέντρο.

Η διεργασία P0 θα κάνει Broadcast σε όλους το ποιος κόμβος προστέθηκε .

Κάθε διεργασία ανανεώνει σε κάθε επανάληψη τις τιμές για τις αποστάσεις από το νέο MST για τους κόμβους που της αντιστοιχούν .

Διαμοιρασμός του πίνακα γειτνίασης και του πίνακα αποστάσεων σε 1D blocks.

Γιατί ;;



# Run time formulation

## Computation cost

- Σε κάθε βήμα το update των τιμών του πίνακα αποστάσεων σε κάθε διεργασία είναι  $\Theta(n/p)$  , αφού κάθε διεργασία έχει αναλάβει  $n/p$  κόμβους και το update με adjacency matrix είναι  $O(1)$
- Ο υπολογισμός του ελάχιστου πάνω στους  $p$  κόμβους είναι πάλι για κάθε διεργασία  $\Theta(n/p)$
- Για τις  $n$  επαναλήψεις  
computation cost=  $\Theta(n * n/p)$



# Communication cost για υλοποίηση σε αρχιτεκτονική κατανεμημένης μνήμης

Σε κάθε επανάληψη έχουμε

- all-to-one-reduction για υπολογισμό της ελάχιστης ακμής προς εισαγωγή ( $\Theta(\log p)$ )
- one-to-all broadcast για την γνωστοποίηση σε όλους της ακμής που τελικά εισάχθηκε στο MST ( $\Theta(\log p)$ )

Κόστος:  $\Theta(N \log p)$  για τις  $N$  επαναλήψεις

# Speedup and Efficiency

$$S = \frac{\Theta(n^2)}{\Theta(n^2/p) + \Theta(n \log p)}$$

$$E = \frac{1}{1 + \Theta((p \log p)/n)}$$

# MPI: Message Passing Interface

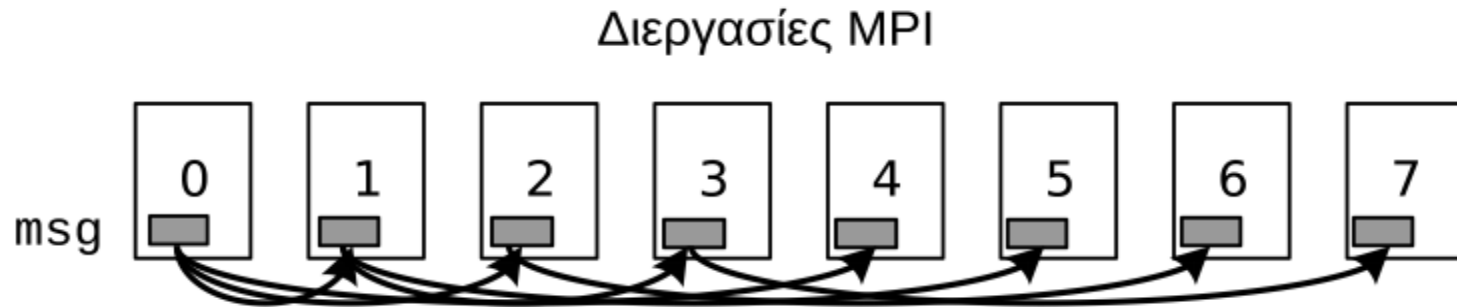
- Όλες οι διεργασίες εκτελούν το ίδιο πρόγραμμα
- Κάθε διεργασία επεξεργάζεται υποσύνολο των δεδομένων ή διαφοροποιεί τη ροή εκτέλεσης της με βάση το βαθμό (rank) που της αποδίδει το MPI

# Συλλογική Επικοινωνία στο MPI

Παράδειγμα:

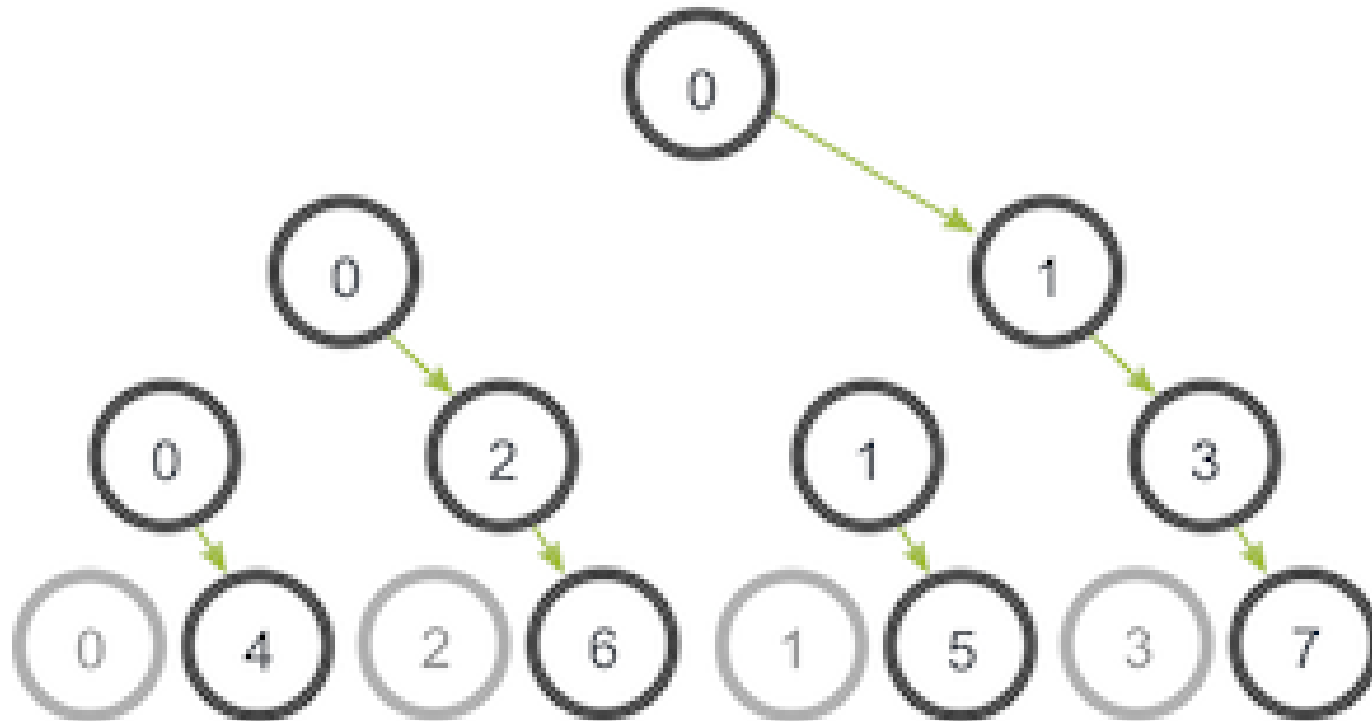
Αποστολή του msg στις διεργασίες 1-7 από τη 0

```
MPI_Bcast(msg,count,MPI_FLOAT,0,MPI_COMM_WORLD);
```

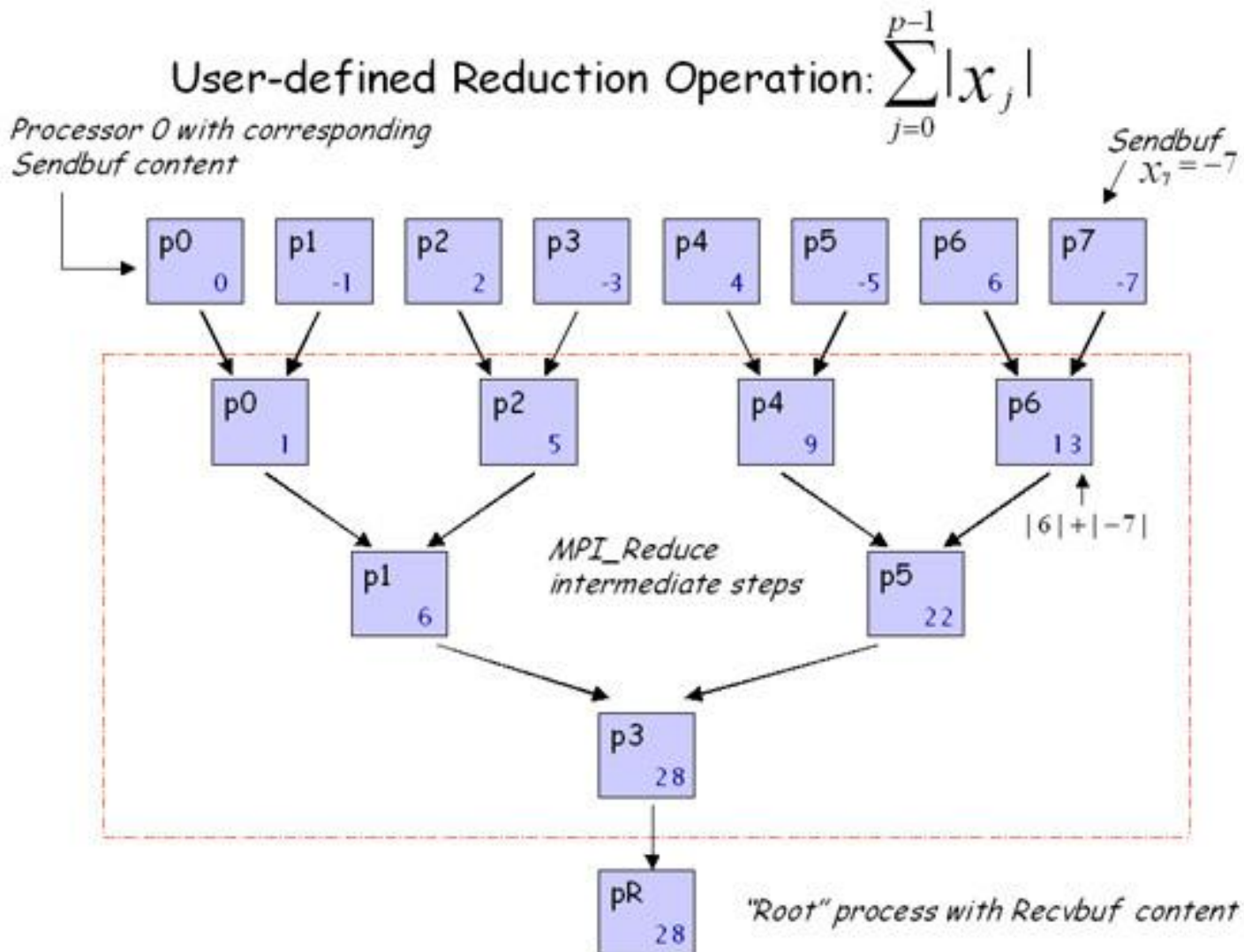


Γενικά: Για  $p$  διεργασίες έχουμε  $\log_2 p$  βήματα επικοινωνίας

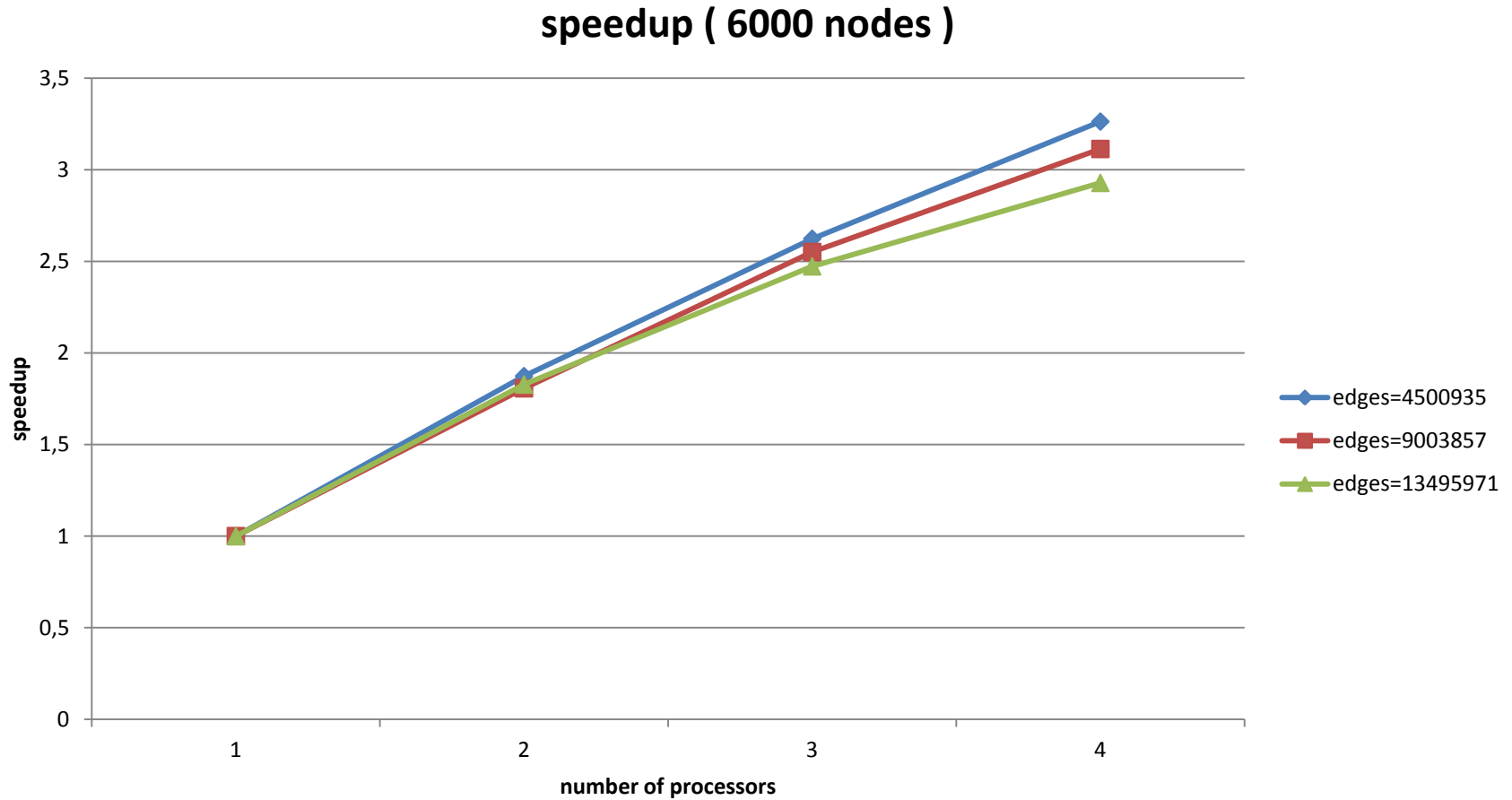
# Broadcast communication tree



# Reduce communication tree

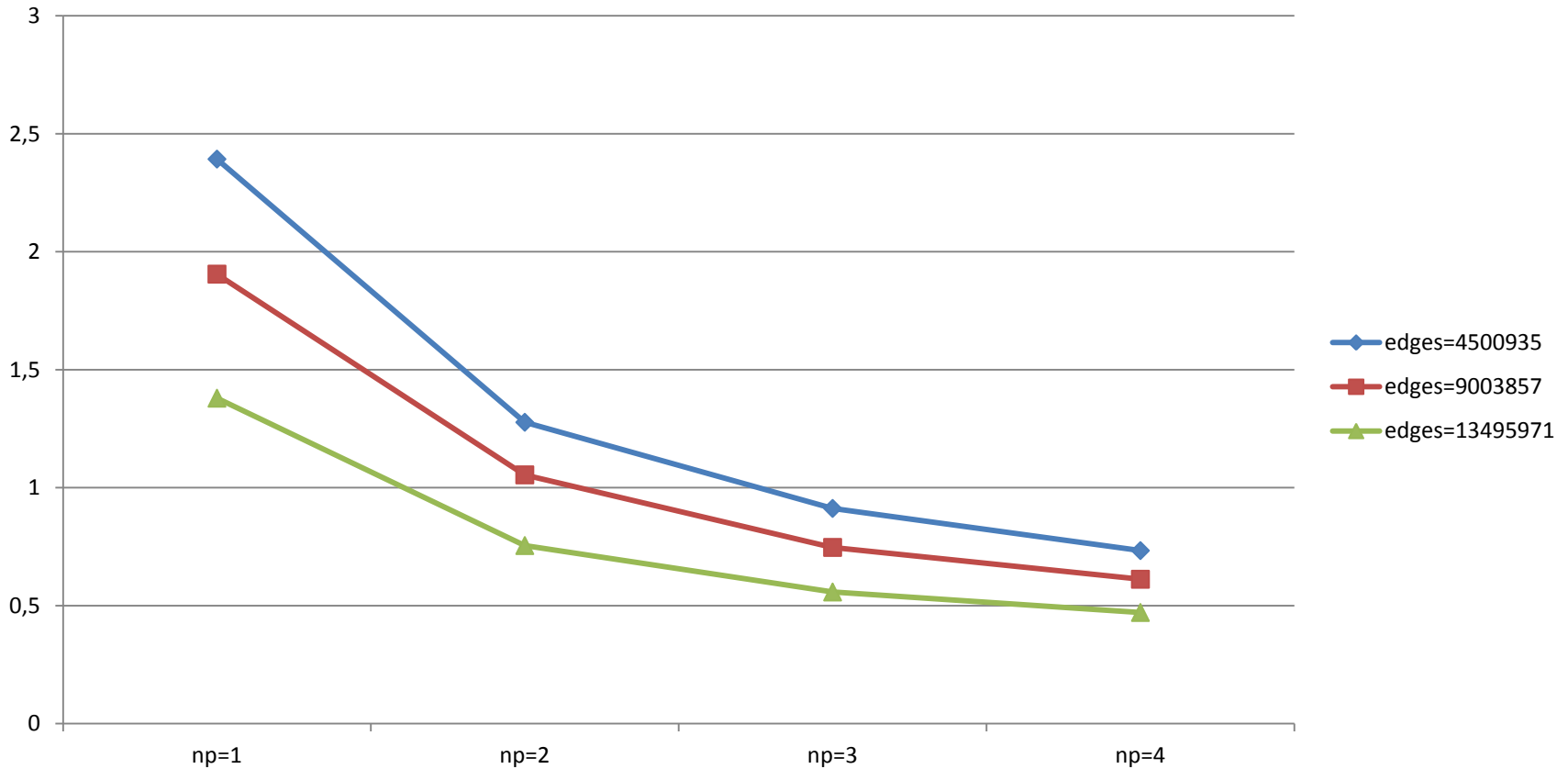


# Αποτελέσματα



# Αποτελέσματα

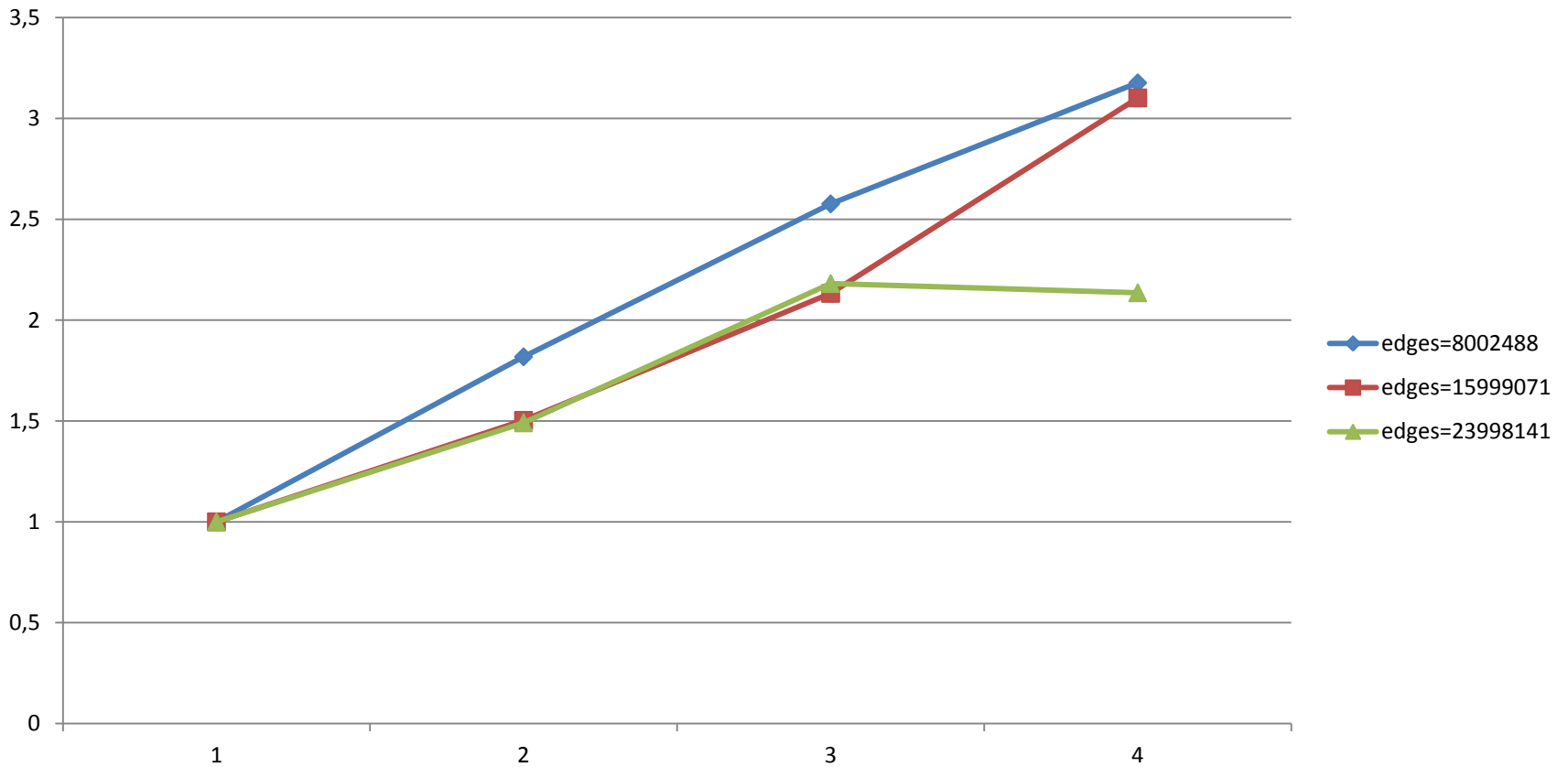
time (6000 nodes)





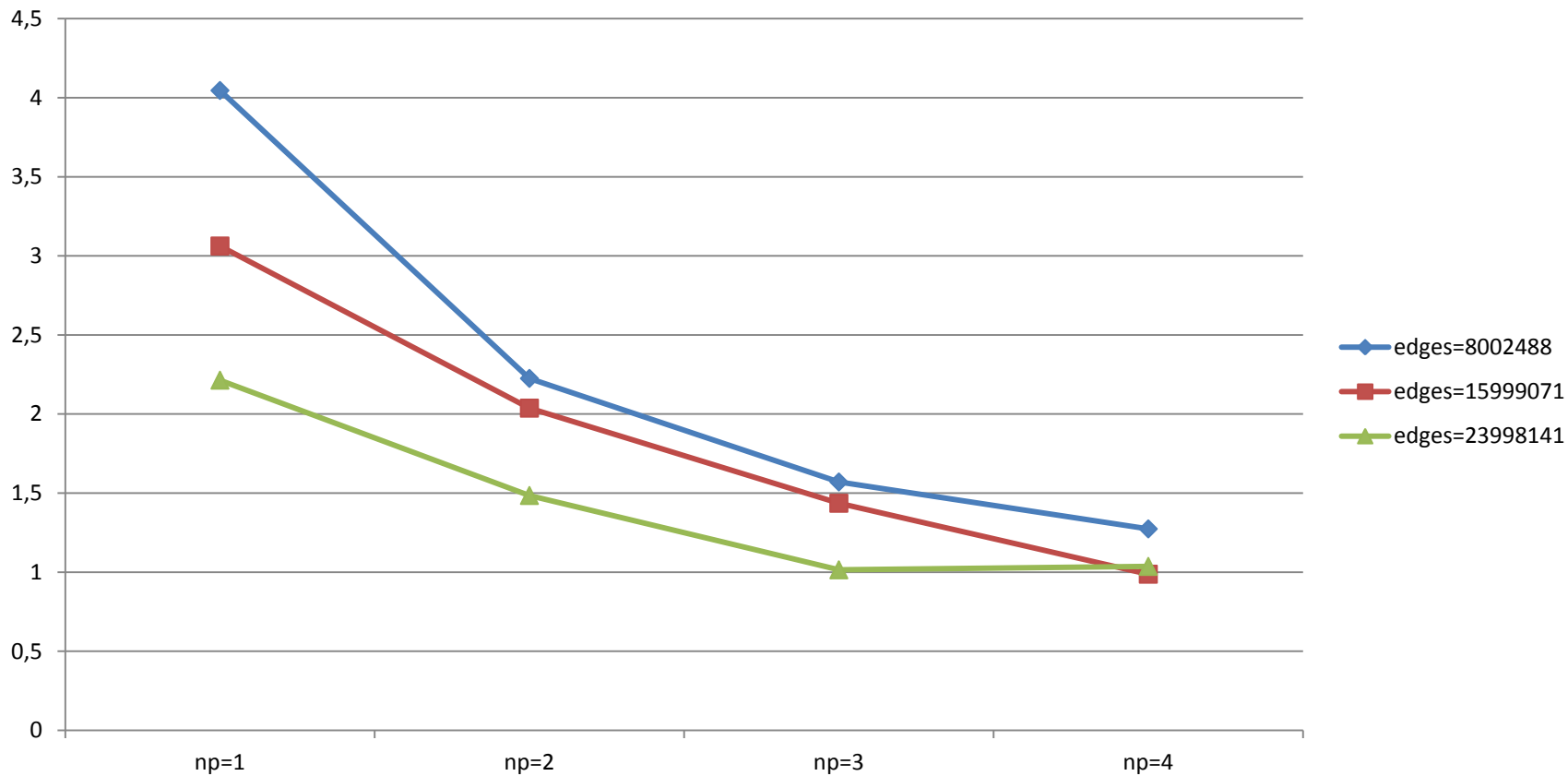
# Αποτελέσματα

speedup (8000 nodes)



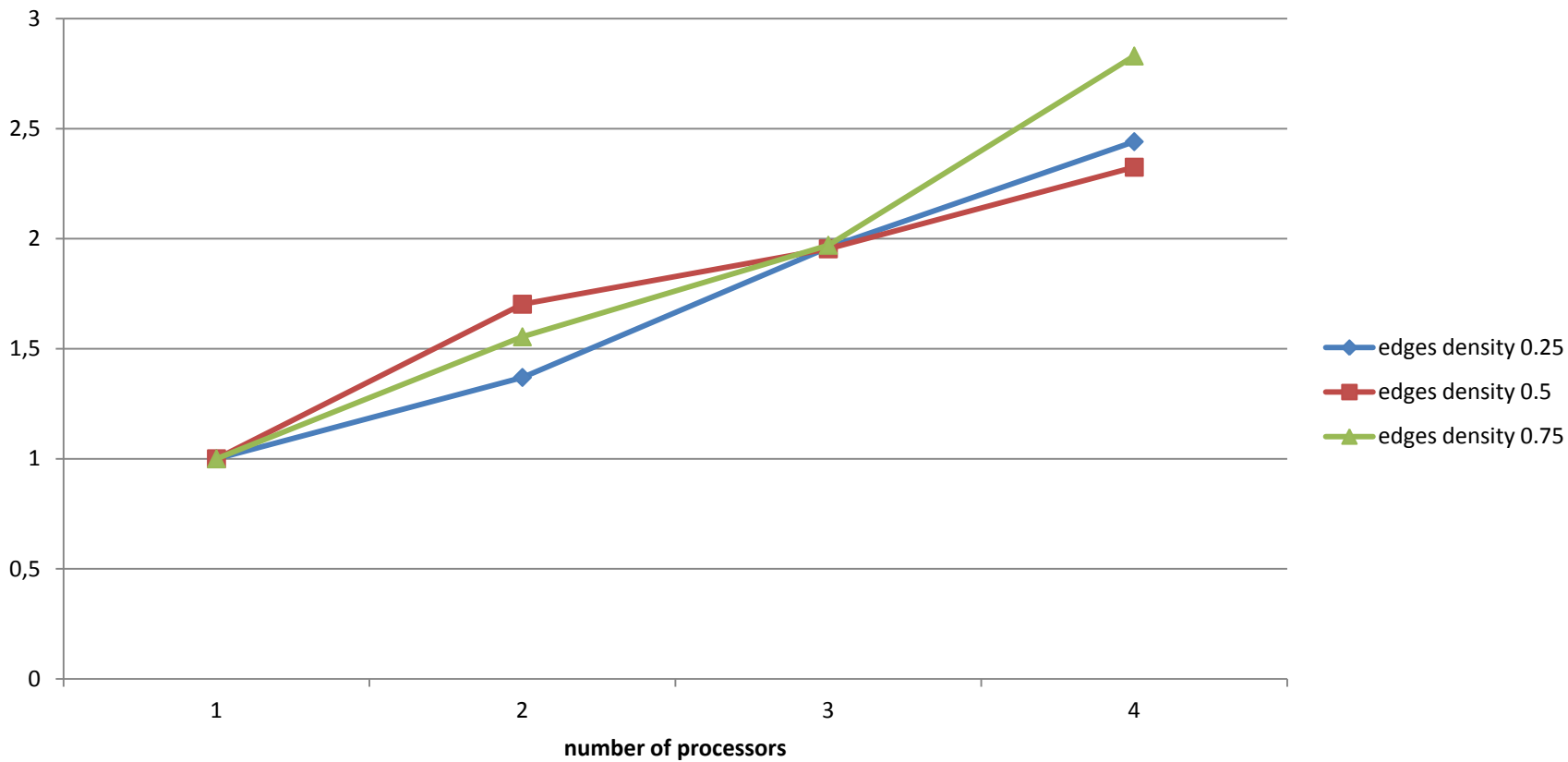
# Αποτελέσματα

time (8000 nodes)



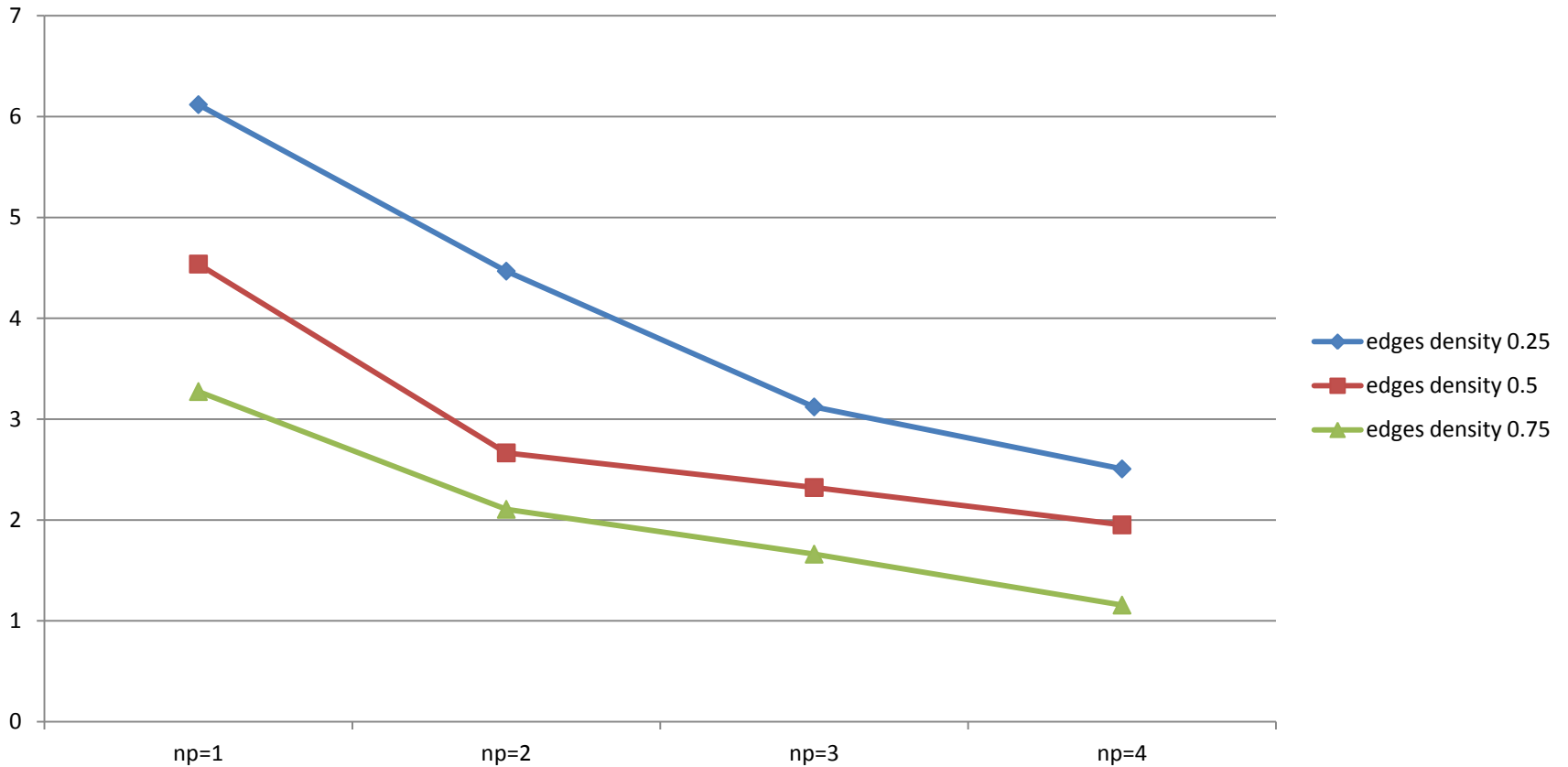
# Αποτελέσματα

speedup (10000 nodes)

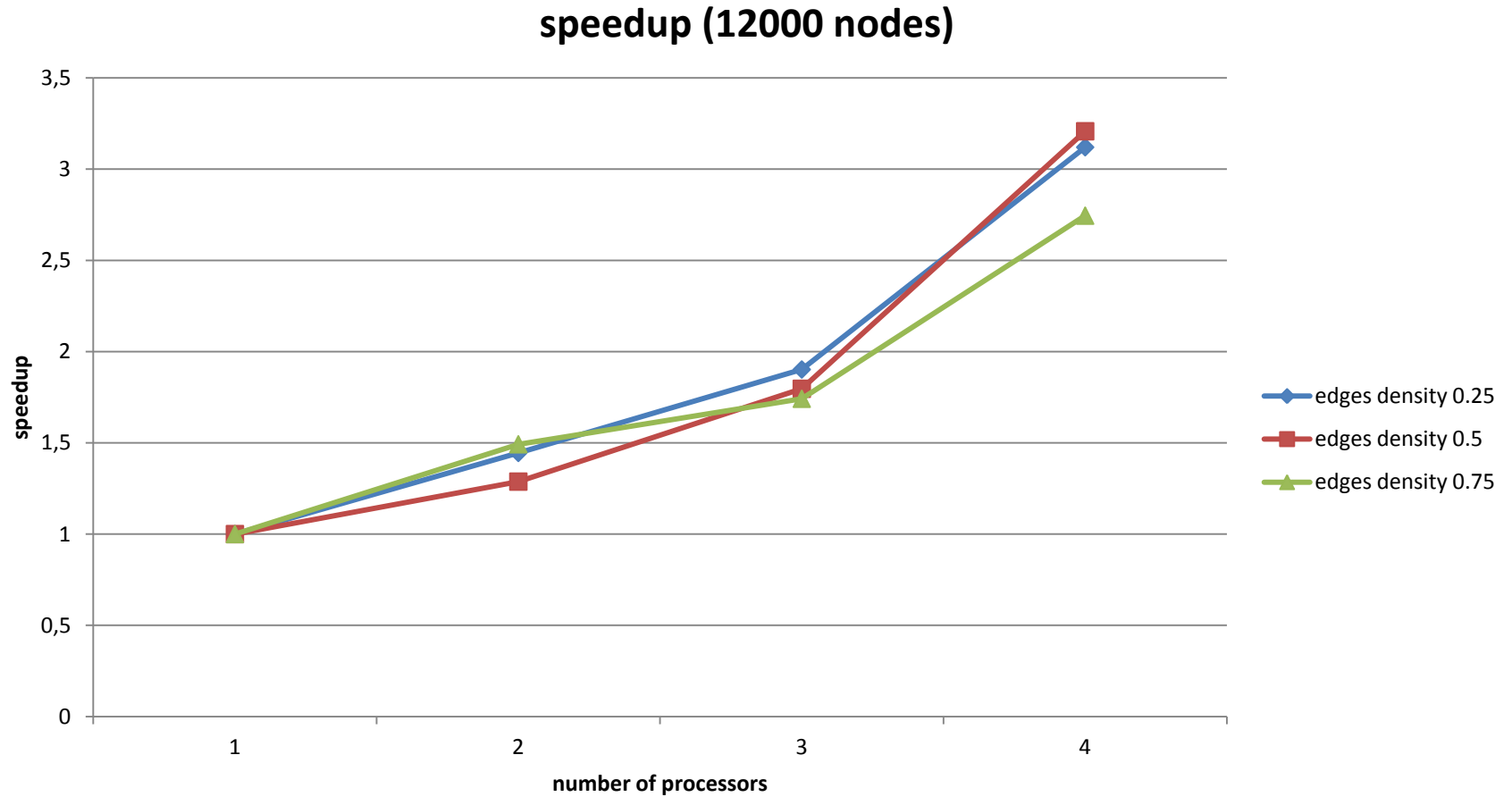


# Αποτελέσματα

time (10000 nodes)

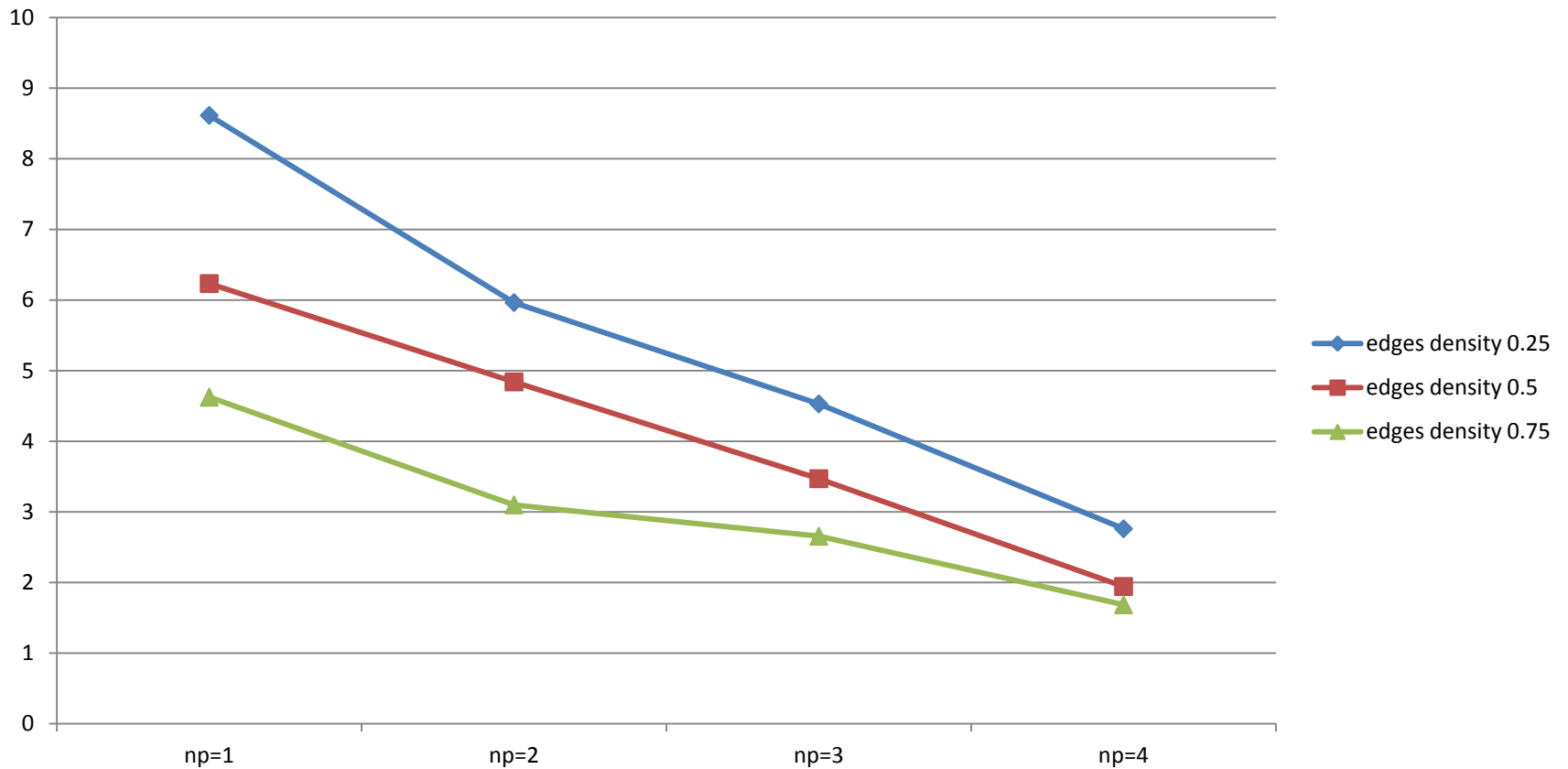


# Αποτελέσματα

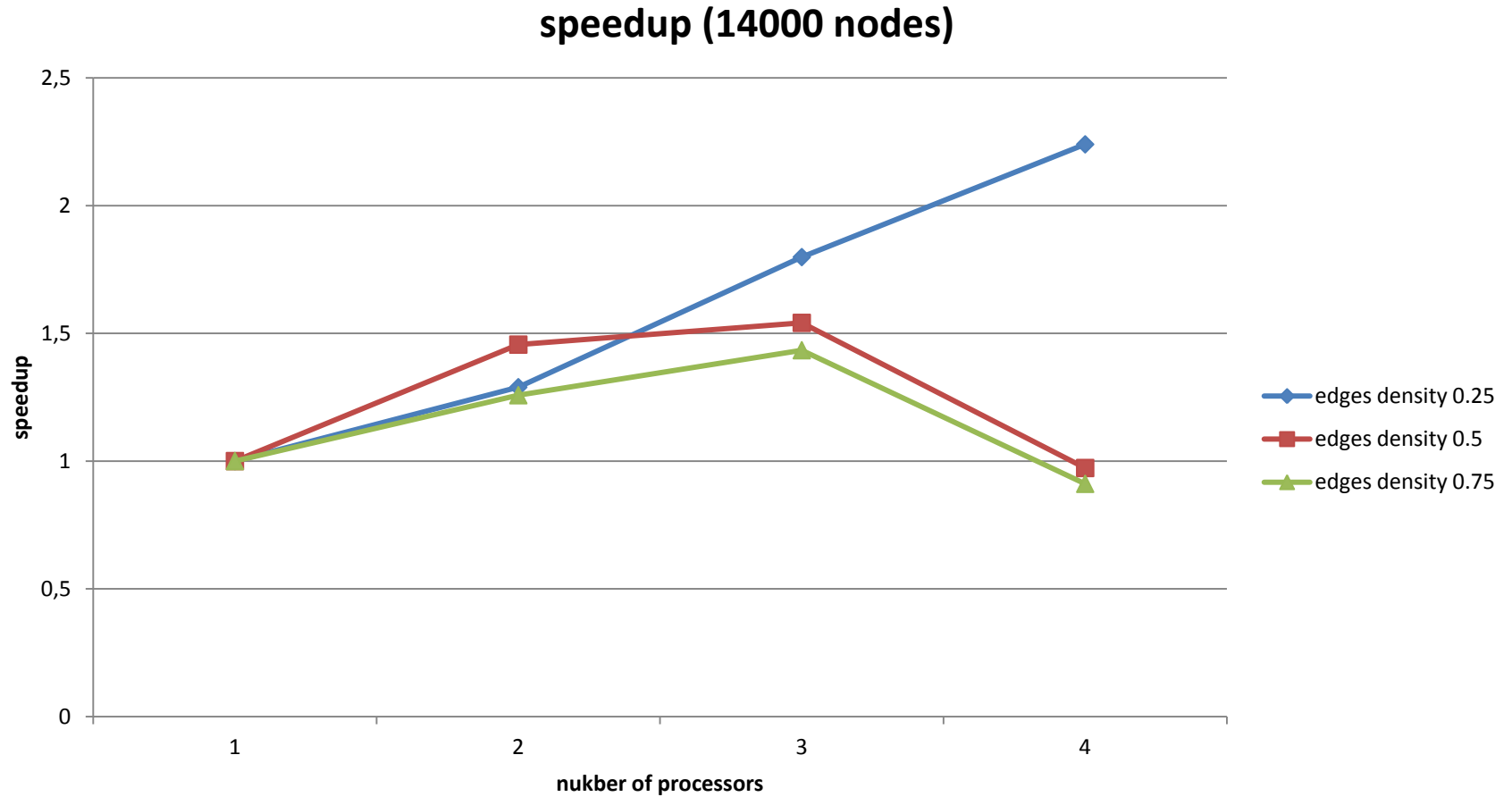


# Αποτελέσματα

time (12000 nodes)



# Αποτελέσματα



# Αποτελέσματα

time (14000 nodes)

