

# **ΟΜΑΔΙΚΗ ΕΡΓΑΣΙΑ ΣΤΟ ΜΑΘΗΜΑ BIG DATA / ΑΝΑΛΥΣΗ ΔΕΔΟΜΕΝΩΝ ΜΕΓΑΛΟΥ ΟΓΚΟΥ**

Μέλη ομάδας:

Ζέρβας Μιχάλης (ics20015)

Δρίκος Χρήστος(ics20046)

# Περιεχόμενα

1. Αλγόριθμος που χρησιμοποιήσαμε για να επιλύσουμε το πρόβλημα	3
2. Περιορισμοί και θεωρήσεις της λύσης	3
3. Επιπλέον πακέτα που χρησιμοποιήθηκαν	4
4. Παράμετροι του συστήματος	4
5. Γραφικές παραστάσεις και πίνακες χρόνων	4
6. Σχολιασμός χρόνων	13

Για την υλοποίηση χρησιμοποιήθηκε ένας κύκλος MapReduce.

## 1. Αλγόριθμος που χρησιμοποιήσαμε για να επιλύσουμε το πρόβλημα

1. Ο mapper λαμβάνει ως είσοδο μία παραγγελία
2. Δημιουργεί όλους τους δυνατούς μοναδικούς συνδυασμούς προϊόντων από τη συγκεκριμένη παραγγελία. Η σειρά των προϊόντων δεν παίζει κάποιο ρόλο. Παράδειγμα ο συνδυασμός F1 F3 θεωρείται ίδιος με τον F3 F1.
3. Ο Mapper περνάει ως κλειδί κάθε δυνατό συνδυασμό προϊόντων που έχει παράχθηκε και ως value για τον καθένα από αυτούς την τιμή 1. Παράδειγμα  $\langle [F1\ F3], 1 \rangle$
4. Οι Combiners λαμβάνουν ως input το output των Mappers και παράγουν ως έξοδο: για key, κάθε συνδυασμό προϊόντων (προφανώς έναν τη φορά) και για value το άθροισμα των values που αντιστοιχούν σε κάθε συνδυασμό. Π.χ.  $\langle [F1\ F3], 300 \rangle$
5. Ο Reducer, παίρνοντας ως είσοδο τα αποτελέσματα του Combiner, αθροίζει συνολικά για κάθε συνδυασμό προϊόντων όλα τα values που έχει πάρει από τους combiners και αν το άθροισμα είναι μεγαλύτερο ή ίσο από το threshold, γράφεται ως έξοδος το αποτέλεσμα με key τον συγκεκριμένο συνδυασμό και value το πόσες φορές εμφανίστηκε. Αν είναι μικρότερο από το threshold το αγνοεί. Π.χ.  $\langle [F1\ F3], 6000 \rangle$

## 2. Περιορισμοί και θεωρήσεις της λύσης

Περιορισμοί: Ο κώδικας του combiner έπρεπε να είναι λίγο διαφορετικός από του reducer λόγω ελέγχου threshold. Επιλέξαμε να κάνουμε συμπίεση των αποτελεσμάτων που παράγει ο mapper, διότι έτσι μειώνεται ο όγκος των δεδομένων που παράγει. Συνεπώς μεταφέρονται γρηγορότερα στο δίκτυο και μειώνονται τα I/O στον δίσκο. Επιλέξαμε την συμπίεση Snappy για να έχουμε την γρηγορότερη δυνατή συμπίεση, ώστε να μην ζημιωθούμε εν τέλη από το τη διαδικασία της συμπίεσης.

Θεώρηση της αρχικής μας λύσης: Η αρχική μας ιδέα ήταν ότι αν κάποιο μεμονωμένο προϊόν εμφανίζεται πιο λίγες φορές από το threshold να

αφαιρεθεί αυτό καθώς και όλοι οι άλλοι συνδυασμοί στους οποίους εμφανίζεται, γλυτώνοντας έτσι περιττούς ελέγχους. Αυτό όμως απορρίφθηκε διότι μας δυσκόλεψε στην υλοποίηση. Θεωρήσαμε ότι θα διάβαζε τουλάχιστον 2 φορές τα αρχεία εισόδου και θα έγραφε τα ενδιάμεσα αποτελέσματα άλλες τόσες καθώς θα έπρεπε να μεταφερθούν στο δίκτυο. Αυτό το θεωρήσαμε χρονοβόρο αφού το πρόβλημα μπορούσαμε να το υλοποιήσουμε και με έναν κύκλο MapReduce.

### 3. Επιπλέον πακέτα που χρησιμοποιήθηκαν

Τα πακέτα αυτά χρησιμοποιήθηκαν στον Mapper για την δημιουργία όλων των δυνατών συνδυασμών των προϊόντων μιας παραγγελίας.

- `com.google.common.collect.Sets`
- `com.google.common.collect.ImmutableSet`

### 4. Παράμετροι του συστήματος

Οι επιπλέον παράμετροι είναι το `threshold` που το δίνουμε ως `args`. Επίσης δεν χρειάζεται να δίνεται το όνομα της `main`.

### 5. Γραφικές παραστάσεις και πίνακες χρόνων

#### 2 Nodes

#### Threshold 5000:

1. 00hrs, 07mins, 45sec job\_1670238997489\_0011

<b>Average Map Time</b>	1mins, 56sec
<b>Average Shuffle Time</b>	2mins, 30sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 29sec

2. 00hrs, 08mins, 12sec job\_1670238997489\_0019

<b>Average Map Time</b>	1mins, 35sec
-------------------------	--------------

<b>Average Shuffle Time</b>	2mins, 56sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 32sec

3. 00hrs, 09mins, 06sec job\_1670238997489\_0020

<b>Average Map Time</b>	1mins, 39sec
<b>Average Shuffle Time</b>	3mins, 41sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 31sec

4. 00hrs, 08mins, 08sec, job\_1670238997489\_0021

<b>Average Map Time</b>	1mins, 47sec
<b>Average Shuffle Time</b>	2mins, 49sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 34sec

5. 00hrs, 12mins, 47sec job\_1670238997489\_0022

<b>Average Map Time</b>	1mins, 51sec
<b>Average Shuffle Time</b>	7mins, 23sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 31sec

### Threshold 10000:

1. 00hrs, 12mins, 45sec job\_1670238997489\_0024

<b>Average Map Time</b>	1mins, 51sec
<b>Average Shuffle Time</b>	7mins, 18sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 32sec

2. 00hrs, 10mins, 34sec job\_1670238997489\_0025

<b>Average Map Time</b>	1mins, 44sec
<b>Average Shuffle Time</b>	5mins, 7sec
<b>Average Merge Time</b>	48sec

**Average Reduce Time** 2mins, 34sec

3. 00hrs, 08mins, 02sec job\_1670238997489\_0027

<b>Average Map Time</b>	1mins, 37sec
<b>Average Shuffle Time</b>	2mins, 14sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 33sec

4. 00hrs, 08mins, 12sec job\_1670238997489\_0028

<b>Average Map Time</b>	1mins, 35sec
<b>Average Shuffle Time</b>	2mins, 53sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 33sec

5. 00hrs, 08mins, 02sec job\_1670366157694\_0006

<b>Average Map Time</b>	1mins, 47sec
<b>Average Shuffle Time</b>	2mins, 43sec
<b>Average Merge Time</b>	52sec
<b>Average Reduce Time</b>	2mins, 31sec

**Threshold 50000:**

1. 00hrs, 12mins, 27sec job\_1670238997489\_0029

<b>Average Map Time</b>	1mins, 51sec
<b>Average Shuffle Time</b>	6mins, 57sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 34sec

2. 00hrs, 08mins, 06sec job\_1670238997489\_0030

<b>Average Map Time</b>	1mins, 46sec
<b>Average Shuffle Time</b>	2mins, 46sec
<b>Average Merge Time</b>	49sec
<b>Average Reduce Time</b>	2mins, 33sec

3. 00hrs, 08mins, 02sec job\_1670238997489\_0032

<b>Average Map Time</b>	1mins, 38sec
<b>Average Shuffle Time</b>	2mins, 13sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 33sec

4. 00hrs, 08mins, 16sec job\_1670238997489\_0033

<b>Average Map Time</b>	1mins, 35sec
<b>Average Shuffle Time</b>	2mins, 58sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 32sec

5. 00hrs, 08mins, 09sec job\_1670366157694\_0007

<u><b>Average Map Time</b></u>	<u>1mins, 38sec</u>
<u><b>Average Shuffle Time</b></u>	<u>2mins, 16sec</u>
<u><b>Average Merge Time</b></u>	<u>47sec</u>
<u><b>Average Reduce Time</b></u>	<u>2mins, 34sec</u>

## 1 Node

### Threshold 5000:

1. 00hrs, 08mins, 36sec job\_1670362094916\_0001

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 4sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 34sec

2. 00hrs, 08mins, 35sec job\_1670362094916\_0002

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 6sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 31sec

3. 00hrs, 08mins, 37sec job\_1670362094916\_0003

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 8sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 31sec

4. 00hrs, 08mins, 39sec job\_1670362094916\_0004

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 7sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 35sec

5. 00hrs, 08mins, 36sec job\_1670362094916\_0005

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 2sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 34sec

#### Threshold 10000:

1. 00hrs, 08mins, 43sec job\_1670366157694\_0001

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	3mins, 11sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 32sec

2. 00hrs, 08mins, 39sec job\_1670366157694\_0002

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	3mins, 4sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 34sec

3. 00hrs, 08mins, 20sec job\_1670366157694\_0003

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	1mins, 52sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 30sec

4. 00hrs, 08mins, 32sec job\_1670366157694\_0004



<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 1sec
<b>Average Merge Time</b>	46sec
<b>Average Reduce Time</b>	2mins, 32sec

5. 00hrs, 08mins, 34sec job\_1670366157694\_0005

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	3mins, 7sec
<b>Average Merge Time</b>	46sec
<b>Average Reduce Time</b>	2mins, 31sec

### Threshold 50000:

1. 00hrs, 08mins, 25sec job\_1670366157694\_0008

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	1mins, 57sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 34sec

2. 00hrs, 08mins, 39sec job\_1670366157694\_0009

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	3mins, 10sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 32sec

3. 00hrs, 08mins, 38sec job\_1670366157694\_0010

<b>Average Map Time</b>	1mins, 27sec
<b>Average Shuffle Time</b>	3mins, 6sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 32sec

4. 00hrs, 08mins, 40sec job\_1670366157694\_0011

<b>Average Map Time</b>	1mins, 26sec
<b>Average Shuffle Time</b>	3mins, 8sec
<b>Average Merge Time</b>	48sec
<b>Average Reduce Time</b>	2mins, 33sec

5. 00hrs, 08mins, 28sec job\_1670366157694\_0012

<b>Average Map Time</b>	1mins, 27sec
-------------------------	--------------

<b>Average Shuffle Time</b>	1mins, 58sec
<b>Average Merge Time</b>	47sec
<b>Average Reduce Time</b>	2mins, 35sec

*Average Elapsed Time*

- 2 Nodes:

<b>threshold 5000:</b>	00:08:29
<b>threshold 10000:</b>	00:08:56
<b>threshold 50000:</b>	00:08:10
- 1 Node:

<b>threshold 5000:</b>	00:08:36
<b>threshold 10000:</b>	00:08:35
<b>threshold 50000:</b>	00:08:35

*Average Map Time*

- 2 Nodes:

<b>threshold 5000:</b>	00:01:40
<b>threshold 10000:</b>	00:01:39
<b>threshold 50000:</b>	00:01:40
- 1 Node:

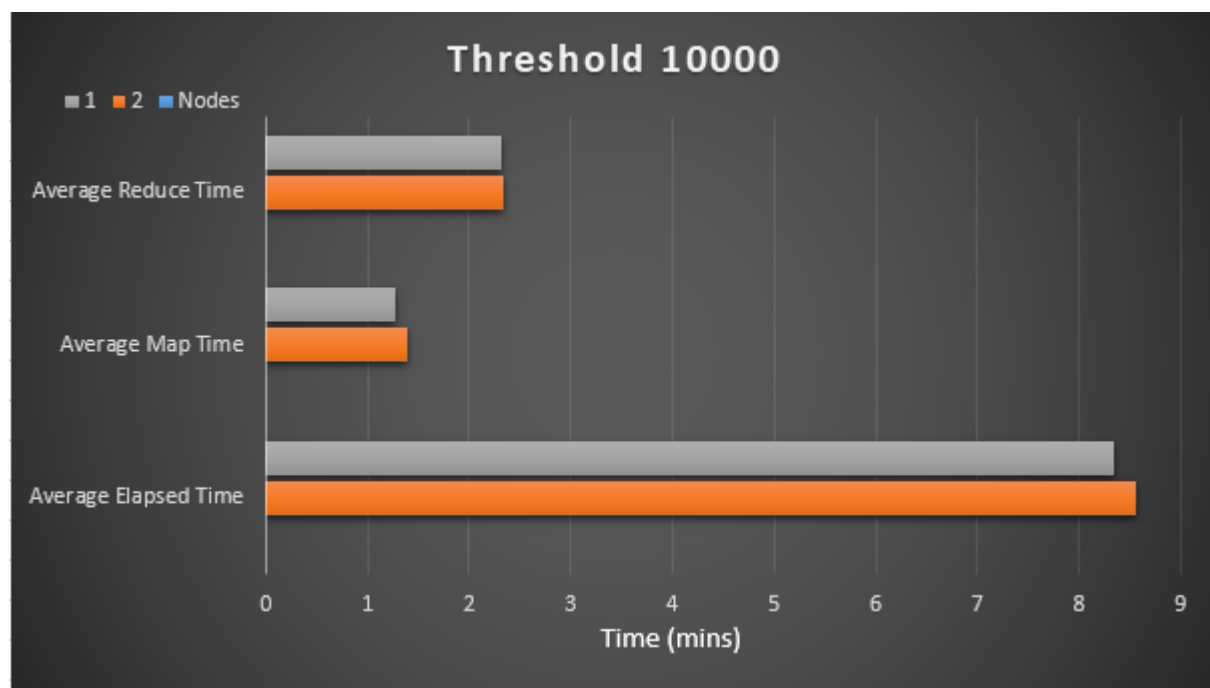
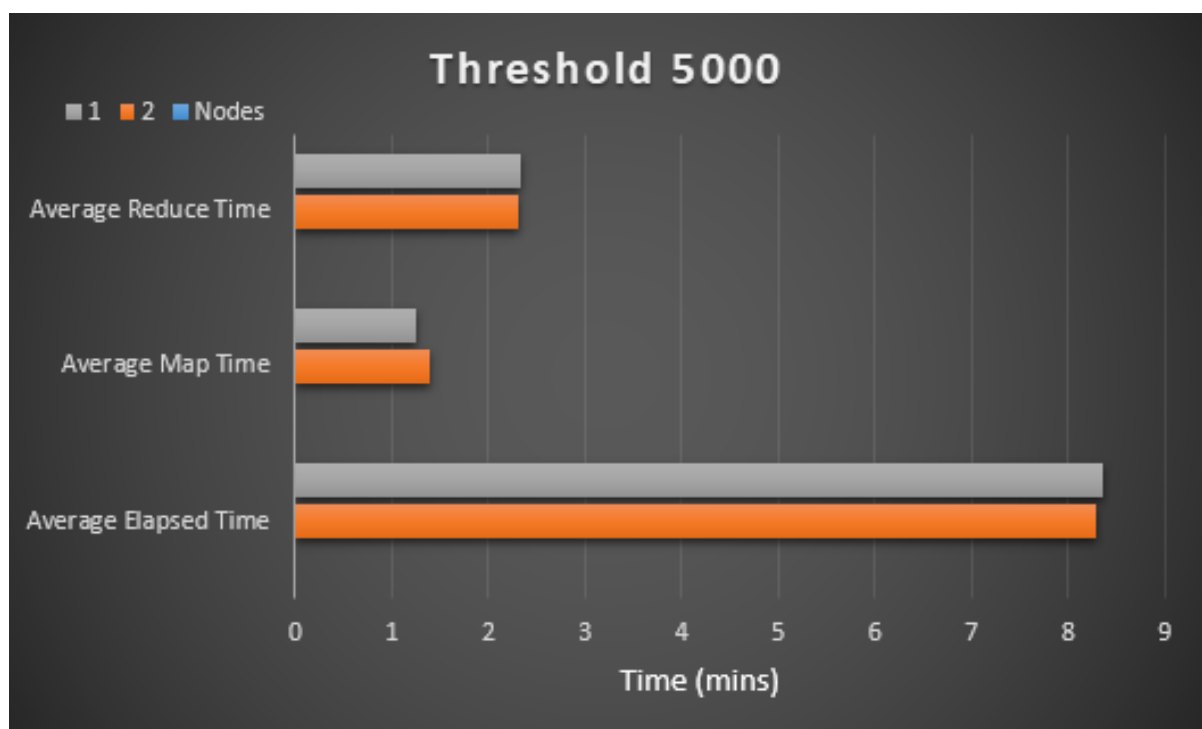
<b>threshold 5000:</b>	00:01:26
<b>threshold 10000:</b>	00:01:27
<b>threshold 50000:</b>	00:01:27

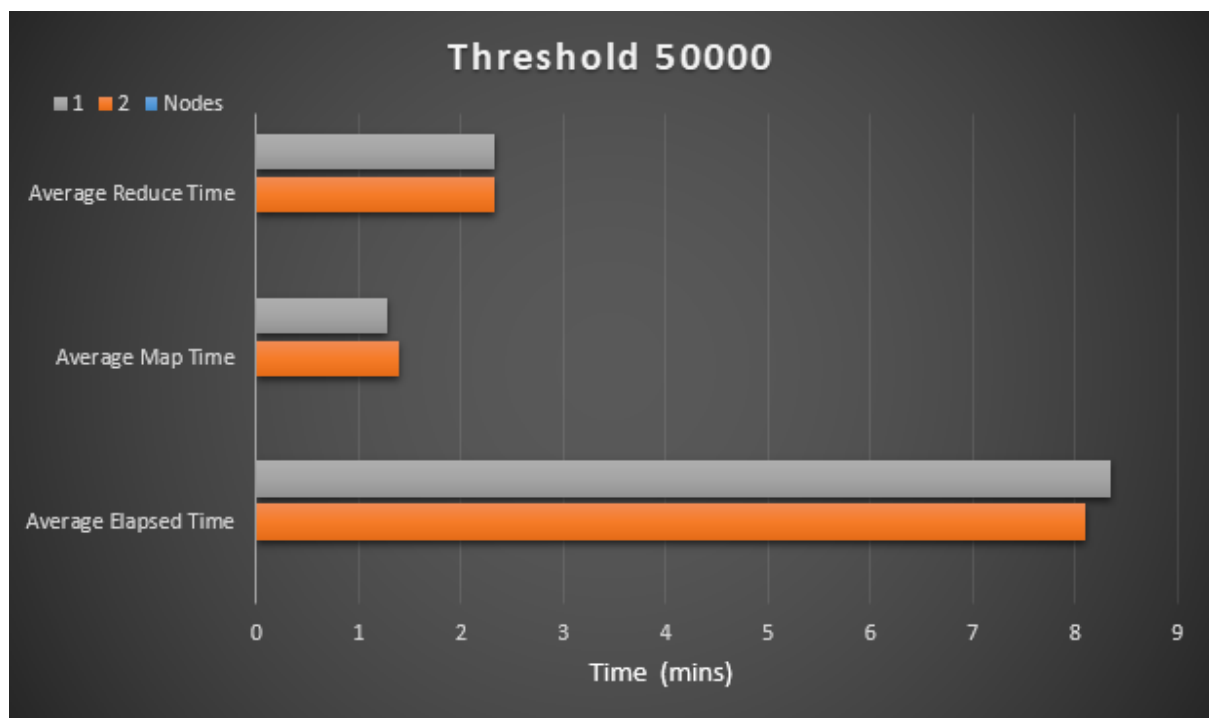
*Average Reduce Time*

- 2 Nodes:

<b>threshold 5000:</b>	00:02:32
<b>threshold 10000:</b>	00:02:33
<b>threshold 50000:</b>	00:02:33
- 1 Node:

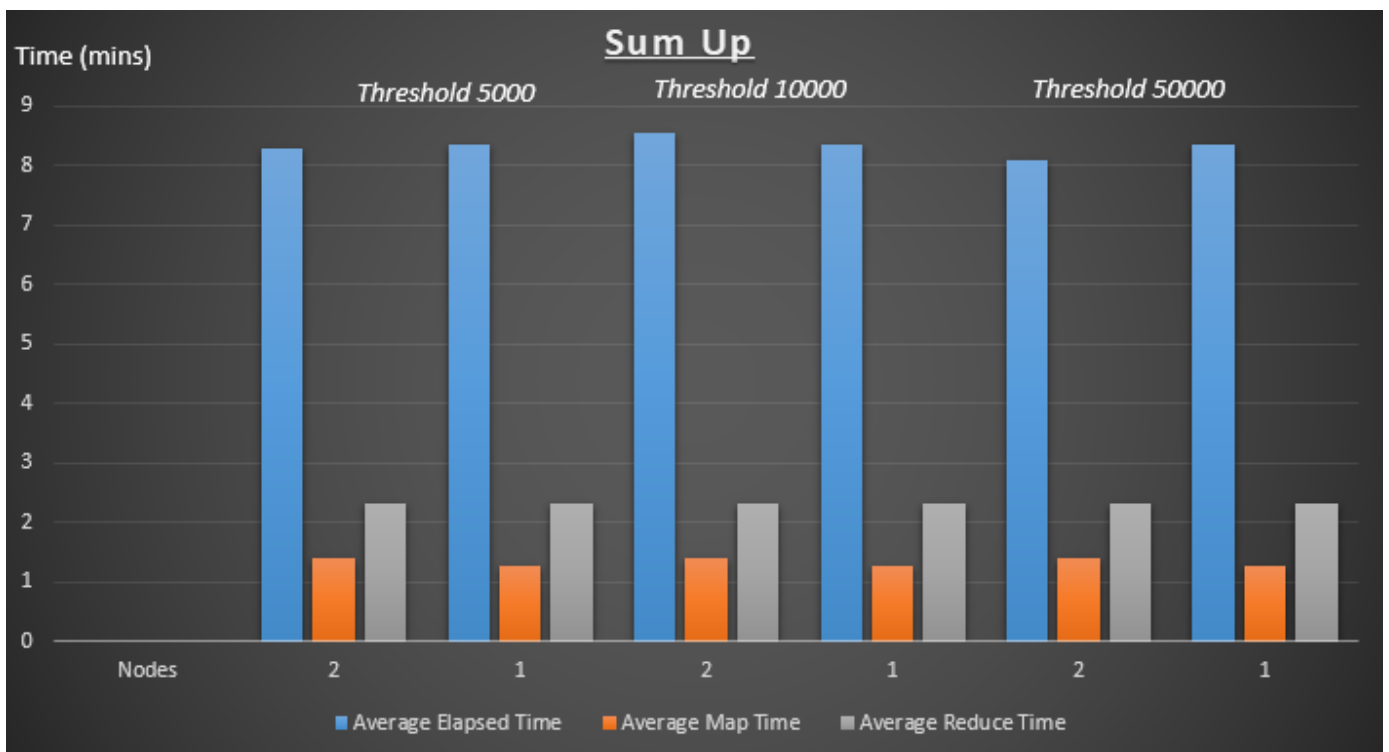
<b>threshold 5000:</b>	00:02:33
<b>threshold 10000:</b>	00:02:32
<b>threshold 50000:</b>	00:02:33





### Συνολικά

Nodes	Average Elapsed Time	Average Map Time	Average Reduce Time	Threshold
2	00:08:29	00:01:40	00:02:32	5000
1	00:08:36	00:01:26	00:02:33	5000
2	00:08:56	00:01:39	00:02:33	10000
1	00:08:35	00:01:27	00:02:32	10000
2	00:08:10	00:01:40	00:02:33	50000
1	00:08:35	00:01:27	00:02:33	50000



## 6. Σχολιασμός χρόνων

Γενικά παρατηρούμε πως δεν υπάρχουν μεγάλες διακυμάνσεις στους χρόνους. Ίσως θα περιμέναμε με 2 κόμβους να είναι αρκετά πιο γρήγορο, αλλά τελικά δεν φάνηκε κάτι τέτοιο. Επιπλέον παρατηρούμε πως το Average Map Time βγήκε λίγο μικρότερο κατά μέσο όρο με 1 κόμβο. Αυτό ίσως οφείλεται στο γεγονός ότι με 2 κόμβους απαιτείται κάποιος επιπλέον χρόνος για να μεταφερθούν τα δεδομένα μέσω δικτύου. Επίσης το Average Reduce Time είναι σταθερό για οποιεσδήποτε παραμέτρους. Τέλος οι μετρήσεις σε 1 Node παρουσιάζουν μικρότερες αποκλίσεις μεταξύ τους, κάτι που οφείλεται στο ότι δεν χρησιμοποιείται το δίκτυο.