

ΑΝΩΤΑΤΟ ΤΕΧΝΟΛΟΓΙΚΟ ΕΚΠΑΙΔΕΥΤΙΚΟ ΙΔΡΥΜΑ ΠΕΙΡΑΙΑ
ΤΜΗΜΑ ΗΛΕΚΤΡΟΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΙΚΩΝ ΣΥΣΤΗΜΑΤΩΝ
ΤΟΜΕΑΣ ΑΡΧΙΤΕΚΤΟΝΙΚΗΣ Η/Υ, ΠΛΗΡΟΦΟΡΙΚΗΣ & ΔΙΚΤΥΩΝ

Εργ. Τεχνολογίας Λογισμικού & Υπηρεσιών

S²E Lab

ΠΤΥΧΙΑΚΗ ΕΡΓΑΣΙΑ

Σπουδάστρια: Κοντζοπούλου Παναγιώτα

Θέμα:

***Μελέτη-σχεδίαση εφαρμογής σε υπολογιστικό νέφος
(cloud computing) με έμφαση στην κατασκευή δέντρων.***



Εισηγητής:

Δ. Ν. Καλλέργης, MSc.

Εργ. Συνεργάτης

Πνευματικά δικαιώματα

Τα πνευματικά δικαιώματα χρησιμοποίησης του μη πρωτότυπου υλικού της εργασίας ανήκουν στο/στη φοιτητή/-τρια και τον επιβλέποντα εις ολόκληρον, δηλαδή εκάτερος μπορεί να κάνει χρήση αυτών χωρίς τη συναίνεση του άλλου. Τα πνευματικά δικαιώματα χρησιμοποίησης του πρωτότυπου μέρους διπλωματικής εργασίας ανήκουν στο/στη φοιτητή/-τρια και τον επιβλέποντα από κοινού, δηλαδή δεν μπορεί ο ένας από τους δύο να κάνει χρήση αυτού χωρίς τη συναίνεση του άλλου. Κατ' εξαίρεση, επιτρέπεται η δημοσίευση του πρωτότυπου μέρους της εργασίας σε επιστημονικό περιοδικό ή πρακτικά συνεδρίου από τον ένα εκ των δύο, με την προϋπόθεση να αναφέρονται τα ονόματα και των δύο ως συν-συγγραφών. Στην περίπτωση αυτή, προηγείται γραπτή ενημέρωση του μη συμμετέχοντα στη συγγραφή του επιστημονικού άρθρου. Δεν επιτρέπεται η κατά οποιοδήποτε τρόπο δημοσιοποίηση υλικού το οποίο έχει δηλωθεί εγγράφως ως απόρρητο.

Οι υπογράφοντες

Περίληψη

Η παρούσα εργασία αφορά την μελέτη της απόδοσης ενός κατανεμημένου συστήματος βασισμένου στο μοντέλο του υπολογιστικού νέφους (Cloud Computing), υποδομή ως υπηρεσία (Infrastructure As a Service). Γίνεται χρήση του προγραμματιστικού μοντέλου Hadoop MapReduce υλοποιώντας τη παράλληλη κατασκευή δέντρων επιθεμάτων σύμφωνα με τον αλγόριθμο Ukkonen. Τέλος πραγματοποιούνται μετρήσεις οι οποίες δείχνουν τη συμπεριφορά του αλγόριθμου σε παράλληλη εκτέλεση.

Abstract

The aim of this Thesis was to study the performance of a distributed system based on the model of Cloud Computing, Infrastructure as a Service. It uses the Hadoop MapReduce programming model by implementing the parallel construction of suffix trees according to Ukkonen's algorithm. Finally, experiments were conducted to evaluate the behavior of the algorithm with parallel execution.

Περιεχόμενα

ΠΕΡΙΛΗΨΗ	1
ABSTRACT	1
ΚΑΤΑΛΟΓΟΣ ΕΙΚΟΝΩΝ	5
ΚΑΤΑΛΟΓΟΣ ΠΙΝΑΚΩΝ	6
ΚΕΦΑΛΑΙΟ 1 ΕΙΣΑΓΩΓΗ	7
1.1 ΑΝΤΙΚΕΙΜΕΝΟ ΤΗΣ ΠΤΥΧΙΑΚΗΣ ΕΡΓΑΣΙΑΣ	7
1.2 ΣΤΟΧΟΣ	8
1.3 ΔΟΜΗ	8
ΚΕΦΑΛΑΙΟ 2 ΘΕΩΡΗΤΙΚΟ ΥΠΟΒΑΘΡΟ	10
2.1 ΥΠΟΛΟΓΙΣΤΙΚΟ ΝΕΦΟΣ	10
2.1.1. Ιστορική Αναδρομή	11
2.1.2 Πρόδρομοι του Υπολογιστικού Νέφους	13
2.1.3 Χαρακτηριστικά του Υπολογιστικού Νέφους	16
2.1.4 Ομοιότητες και Διαφορές μεταξύ Υπολογιστικού Πλέγματος και Νέφους	17
2.1.5 Γενική Εικόνα του Υπολογιστικού Νέφους	20
2.1.6 Μοντέλα Υπολογιστικού Νέφους	21
2.1.7 Μοντέλα Ανάπτυξης Υπολογιστικού Νέφους	27
2.1.8 Σύμβαση Παροχής Υπηρεσιών (Service Level Agreement SLA)	28
2.1.9 Ασφάλεια στο Υπολογιστικό Νέφος	29
2.1.10 Εφαρμογές Υπολογιστικού Νέφους	31
2.2 ΑΡΧΗ ΗΑΔΟΟΡ	44
2.2.1 Hadoop HDFS	46
2.2.2 Το Προγραμματιστικό Μοντέλο Hadoop MapReduce	50
ΚΕΦΑΛΑΙΟ 3 ΕΡΓΑΛΕΙΑ ΚΑΙ ΜΕΘΟΔΟΙ	63
3.1 ΠΕΡΙΒΑΛΛΟΝ ΑΝΑΠΤΥΞΗΣ - ΕΡΓΑΛΕΙΑ	63
3.1.1 Eclipse	63
3.1.2 Java	64
3.1.3 Maven	66
3.1.4 VMware	66
3.1.5 Okeanos	68
3.2 ΠΕΡΙΓΡΑΦΗ ΑΛΓΟΡΙΘΜΟΥ	70
3.3 ΜΕΘΟΔΟΙ ΚΑΙ ΚΛΑΣΕΙΣ	74
3.3.1 Φάση Map	74
3.3.2 Φάση Reduce	78
3.3.3 Η Κλάση MyAppsMRDriver	79

3.3.4 Διαγράμματα Κλάσεων	81
ΚΕΦΑΛΑΙΟ 4 ΠΕΙΡΑΜΑΤΙΚΗ ΑΞΙΟΛΟΓΗΣΗ	82
4.1 ΕΠΙΛΟΓΗ ΚΑΙ ΧΑΡΑΚΤΗΡΙΣΤΙΚΑ ΤΟΥ CLUSTER	82
4.2 ΕΓΚΑΤΑΣΤΑΣΗ ΕΝΟΣ HADOOP MAPREDUCE CLUSTER	84
4.3 ΠΕΙΡΑΜΑΤΑ	88
4.4 ΣΥΜΠΕΡΑΣΜΑΤΑ	91
ΚΕΦΑΛΑΙΟ 5 ΕΠΙΛΟΓΟΣ	92
5.1 ΜΕΛΛΟΝΤΙΚΕΣ ΕΡΓΑΣΙΕΣ	92
ΒΙΒΛΙΟΓΡΑΦΙΑ	94
ΠΑΡΑΡΤΗΜΑ	98
ΚΑΤΑΛΟΓΟΣ ΚΩΔΙΚΑ	98

Κατάλογος εικόνων

Εικόνα 2.1.5-1 Σχέσεις μεταξύ των μοντέλων του υπολογιστικού νέφους.....	20
Εικόνα 2.1.9-1 Το μοντέλο CIA Triad.....	29
Εικόνα 2.1.9-2 Προϋποθέσεις ασφαλείας χωρισμένες σύμφωνα με το μοντέλα υπολογιστικού νέφους.....	31
Εικόνα 2.1.10.1-1 Επίσημο Λογότυπο του Google Apps.....	31
Εικόνα 2.1.10.1-2 Επίσημο Λογότυπο του Microsoft Office 365.....	32
Εικόνα 2.1.10.2-1 Επίσημο Λογότυπο του Google App Engine	32
Εικόνα 2.1.10.2-2 Επίσημο Λογότυπο του Heroku	34
Εικόνα 2.1.10.3-1 Επίσημο Λογότυπο του Amazon Web Services	36
Εικόνα 2.1.10.3-2 Επίσημο Λογότυπο του Windows Azure.....	40
Εικόνα 2.1.10.3-3 Επίσημο Λογότυπο του Eucalyptus.....	41
Εικόνα 2.1.10.3-4 Επίσημο Λογότυπο του OpenStack.....	42
Εικόνα 2.2-1 Επίσημο Λογότυπο του Hadoop	44
Εικόνα 2.2-2 Διάγραμμα του Οικοσυστήματος του Hadoop.....	45
Εικόνα 2.2.1.1-1 Αρχιτεκτονική HDFS	47
Εικόνα 2.2.1.2-1 Πολιτική τοποθεσίας αντιγράφων του HDFS.....	50
Εικόνα 2.2.2-1 Δημιουργία νέας λίστας εξόδου	51
Εικόνα 2.2.2-2 Παίρνει ως είσοδο μια λίστα τιμών και επιστρέφει μια τιμή.....	52
Εικόνα 2.2.2.1-1 Ομαδοποίηση δεδομένων σύμφωνα με το κλειδί.	53
Εικόνα 2.2.2.2-1 Διάγραμμα Ροής Εργασιών Υψηλού επιπέδου	56
Εικόνα 2.2.2.2-2 Λεπτομερής Περιγραφή της Ροής Εργασιών ενός Hadoop MapReduce προγράμματος	57
Εικόνα 2.2.2.3-1 Διάγραμμα Ροής Εργασιών προσθέτοντας το στάδιο της combiner διεργασίας	60
Εικόνα 3.1.1-1 Επίσημο Λογότυπο του Eclipse.....	63
Εικόνα 3.1.2-1 Επίσημο Λογότυπο της Java.....	64
Εικόνα 3.1.3-1 Επίσημο Λογότυπο της Maven	66
Εικόνα 3.1.4-1 Επίσημο Λογότυπο της VMware.....	67
Εικόνα 3.1.5-1 Γενική εικόνα της αρχιτεκτονικής του Synnefo	69
Εικόνα 3.2-1 Αριστερά το δέντρο επιθεμάτων της ακολουθίας xabxax και δεξιά το πεπλεγμένο δέντρο επιθεμάτων της ίδιας ακολουθίας.	72
Εικόνα 3.3.3-1 Διάγραμμα κλάσεων με βάση τη MyAppsMRDriver.	81
Εικόνα 3.3.3-2 Συνολικό διάγραμμα κλάσεων	81
Εικόνα 4.1-1 Απεικόνιση των διαθέσιμων μηχανημάτων/χαρακτηριστικών/κατάστασης	83
Εικόνα 4.1-2 Απεικόνιση του δικτύου στο οποίο ανήκουν τα μηχανήματα.	84
Εικόνα 4.2-1 Εκτέλεση των εντολών για την εκκίνηση του Hadoop στο cluster	86
Εικόνα 4.2-2 Απεικόνιση της κατάστασης του cluster.....	86
Εικόνα 4.2-3 Απεικόνιση μίας ολοκληρωμένης εργασίας	87
Εικόνα 4.2-4 Εκτέλεση των εντολών για τη διακοπή του cluster	87
Εικόνα 4.3-1 Γράφημα απεικόνισης του χρόνου εκτέλεσης σύμφωνα με τον αριθμό των κόμβων για τα τρία πακέτα δεδομένων.....	90

Κατάλογος πινάκων

Πίνακας 2.1.4-1 Κοινά χαρακτηριστικά υπολογιστικού πλέγματος και νέφους.....	17
Πίνακας 2.1.4-2 Κοινά χαρακτηριστικά με διαφορετική προσέγγιση.....	18
Πίνακας 2.1.4-3 Διαφορές μεταξύ υπολογιστικού πλέγματος και υπολογιστικού νέφους.....	19

ΚΕΦΑΛΑΙΟ 1 ΕΙΣΑΓΩΓΗ

1.1 Αντικείμενο της Πτυχιακής Εργασίας

Η ανάγκη για αποθήκευση και επεξεργασία δεδομένων έχει αυξηθεί κατακόρυφα τα τελευταία χρόνια. Καθημερινά δημιουργείται τεράστιος όγκος πληροφορίας ενώ μέσω του διαδικτύου διακινούνται μεγάλες ποσότητες δεδομένων. Αυτή την ανάγκη έρχονται να καλύψουν τα καταναμημένα συστήματα παράλληλης επεξεργασίας και το υπολογιστικό νέφος (cloud computing).

Το υπολογιστικό νέφος είναι η παροχή υπολογιστικού χώρου υπό τη μορφή υπηρεσίας και όχι προϊόντος, όπου κοινόχρηστοι πόροι, λογισμικό και πληροφορίες παρέχονται σε υπολογιστές και άλλες συσκευές διαμέσου ενός δικτύου. Στη πτυχιακή εργασία, αρχικά, γίνεται εκτενής ανάλυση σε θεωρητικό επίπεδο του υπολογιστικού νέφους έτσι ώστε να γίνει κατανοητή η έννοια, τα διαθέσιμα μοντέλα, τα χαρακτηριστικά κάθε μοντέλου καθώς και τα πλεονεκτήματα και τα μειονεκτήματα του καθενός. Τα μοντέλα χωρίζονται σε δύο κατηγορίες ως προς το είδος της υπηρεσίας και ως προς τον τρόπο ανάπτυξης. Βασιζόμενοι στο είδος της υπηρεσίας, τα διαθέσιμα μοντέλα είναι τα λογισμικό ως υπηρεσία (Software as a service), πλατφόρμα ως υπηρεσία (Platform as a service) και υποδομή ως υπηρεσία (Infrastructure as a service). Τα μοντέλα ανάπτυξης είναι τα δημόσιο υπολογιστικό σύννεφο (public cloud), ιδιωτικό υπολογιστικό σύννεφο (private cloud), κοινοτικό υπολογιστικό σύννεφο (community cloud) και υβριδικό υπολογιστικό σύννεφο (hybrid cloud).

Ακολουθεί μελέτη για το ανοιχτό λογισμικό Apache Hadoop που χρησιμοποιήθηκε για την υλοποίηση του προγράμματος και του μοντέλου υποδομή ως υπηρεσία του υπολογιστικού νέφους. Το Hadoop είναι ένα λογισμικό γραμμένο σε Java που υποστηρίζει την επεξεργασία μεγάλων συνόλων δεδομένων σε ένα καταναμημένο υπολογιστικό σύστημα. Προσφέρει ένα επίπεδο απλοποίησης στη διαδικασία ανάπτυξης παράλληλων προγραμμάτων αναλαμβάνοντας να διαχειριστεί το διαμοιρασμό των δεδομένων, την συγκέντρωση των αποτελεσμάτων, τις πιθανές αποτυχίες κόμβων και

άλλα θέματα. Το προγραμματιστικό μοντέλο Hadoop MapReduce είναι υπεύθυνο για την κατασκευή και την εκτέλεση παράλληλων προγραμμάτων. Με τη βοήθειά του υλοποιείται η παράλληλη κατασκευή δέντρων επιθεμάτων (suffix trees) σύμφωνα με τον αλγόριθμο του Ukkonen.

Το δέντρο επιθεμάτων είναι δομή δεδομένων που αναπαριστά ένα επίθεμα που μπορεί να ανήκει σε μία ή παραπάνω συμβολοσειρές επιτρέποντας τη γρήγορη υλοποίηση πολλών διεργασιών συμβολοσειράς. Συγκεκριμένα για την εφαρμογή χρησιμοποιήθηκε ο αλγόριθμος Ukkonen ο οποίος είναι ένας γραμμικός, online αλγόριθμος που προτάθηκε από τον Esko Ukkonen το 1995.

1.2 Στόχος

Η παρούσα πτυχιακή, αρχικά, δίνει στον αναγνώστη μια γενική εικόνα για την έννοια του λεγόμενου υπολογιστικού νέφους. Οι βασικοί στόχοι της όμως είναι δύο:

- Η υλοποίηση ενός παράλληλου προγράμματος κατασκευής δέντρων επιθεμάτων βασισμένο στον παραπάνω αλγόριθμο και στο προγραμματιστικό μοντέλο του MapReduce για την επίλυση του προβλήματος εύρεσης της μεγαλύτερης κοινής υπό-συμβολοσειράς (Longest common substring problem).
- Η ανάλυση της απόδοσης του παραπάνω προγράμματος σε πραγματικό δίκτυο υπολογιστών που περιλαμβάνει ένα αυξανόμενο πλήθος κόμβων.

1.3 Δομή

Το κείμενο της πτυχιακής αποτελείται από πέντε κεφάλαια. Το παρόν κεφάλαιο είναι το πρώτο και σε αυτό γίνεται αναφορά στο αντικείμενο της πτυχιακής και στον τρόπο οργάνωσης του υπόλοιπου κειμένου.

Στο κεφάλαιο 2 επιχειρείται η παροχή του απαραίτητου υπόβαθρου που χρειάζεται ο αναγνώστης για την κατανόηση του συνόλου της εργασίας. Συγκεκριμένα, γίνεται θεωρητική ανάλυση του υπολογιστικού νέφους, των μοντέλων του και των χαρακτηριστικών του. Στη συνέχεια παρουσιάζονται οι πιο σύγχρονες τάσεις παροχής

δικτυακών υπηρεσιών που βασίζονται σε υπολογιστικά νέφη. Τέλος εισάγεται στον αναγνώστη η έννοια της τεχνολογία του framework Hadoop MapReduce καθώς περιγράφεται ο τρόπος λειτουργίας του Hadoop MapReduce και του κατακευματισμένου συστήματος αρχείων HDFS.

Στο κεφάλαιο 3 παρουσιάζεται η συνεισφορά της εργασίας σε επίπεδο υλοποιήσεων. Αναλυτικότερα, παρουσιάζονται τα εργαλεία που χρησιμοποιήθηκαν, γίνεται μια θεωρητική προσέγγιση στον αλγόριθμο του Ukkonen και έπειτα αναλύονται σημαντικά κομμάτια κώδικα του προγράμματος που υλοποιήθηκε.

Το κεφάλαιο 4 περιλαμβάνει τον τρόπο εγκατάστασης και πραγματοποίησης ενός Hadoop MapReduce cluster και τις μετρήσεις που έγιναν για τη παραπάνω υλοποίηση με σκοπό την κατανόηση της συμπεριφοράς του και της απόδοσης του σε πραγματικές συνθήκες.

Ακολουθεί το κεφάλαιο 5 που αποτελεί την πρόταση για μελλοντικές εργασίες. Συγκεκριμένα αναφέρονται διάφορες εκδοχές υλοποίησης που δίνουν τη δυνατότητα περαιτέρω έρευνας.